

Université de Lille, faculté des Sciences et Technologies

ÉCOLE DOCTORALE MADIS

**EFFICIENT MULTI-OBJECTIVE PURE EXPLORATION:
PARETO SET IDENTIFICATION AND APPLICATIONS TO
CLINICAL TRIALS**

**EXPLORATION PURE À OBJECTIFS MULTIPLES :
IDENTIFICATION DU FRONT DE PARETO ET APPLICATIONS
AUX ESSAIS CLINIQUES DE PHASE PRÉCOCE**

Présentée par

CYRILLE KONE

THÈSE DE DOCTORAT

Spécialité **Informatique**

sous la direction d'Émilie Kaufmann et de Laura Richert

Soutenue publiquement à Villeneuve d'Ascq, le 09/12/2025 devant un jury composé de

M. Vianney Perchet	Professeur, CREST/Ensaë Paris	Rapporteur
M. Wouter M. Koolen	Professeur, CWI	Rapporteur
M. Antoine Chambaz	Professeur, Université Paris Cité	Président du jury
M. Peter Auer	Professeur, University of Leoben	Examineur
M. Clayton Scott	Professeur, University of Michigan	Examineur
M ^{me} Sarah Zohar	Directrice de recherche, INSERM	Invitée
M ^{me} Émilie Kaufmann	Chargée de recherche, CNRS	Directrice de thèse
M ^{me} Laura Richert	Professeure, Université de Bordeaux	Co-directrice de thèse

Centre de Recherche en Informatique, Signal et Automatique de Lille (CRISTAL),
UMR 9189 Équipe Scool, 59650, Villeneuve d'Ascq, France



Remerciements

I would like to express my sincere and profound gratitude to my advisors, Émilie and Laura, for their guidance throughout my doctoral studies. Their intellectual rigour, availability, and steady support have had a decisive influence on this work. I am particularly grateful for their high standards, careful feedback, and continued encouragement, which were essential to both the direction and the completion of this thesis.

À toi, Émilie, merci pour ta générosité et pour la confiance que tu m'as accordée. Travailler à tes côtés au quotidien a été un vrai plaisir. Au fil des années, j'ai énormément appris auprès de toi, tant sur le plan scientifique que dans la manière de faire de la recherche avec rigueur. Ta maîtrise des sujets, ainsi que ton aisance à les rendre accessibles, m'ont permis de développer une compréhension plus profonde et une exigence accrue dans mon travail. Merci aussi pour ta bienveillance, tes conseils, et pour l'autonomie que tu m'as progressivement donnée.

À toi, Laura, merci pour tous ces séjours à Bordeaux et pour ton accueil toujours chaleureux. Tu as su rendre accessibles des notions clés liées aux essais cliniques, tout en nous partageant avec clarté les défis concrets qui ont nourri nos problématiques en bandits. Je te suis très reconnaissant pour cette expérience pluridisciplinaire particulièrement enrichissante, et pour la qualité de nos échanges, à la fois exigeants, ouverts et stimulants. De surcroît, vous avez été toutes deux d'un immense apport dans l'élaboration de mon projet post-thèse. Merci pour votre implication constante, vos conseils, et votre soutien dans cette étape importante. Tout cela montre que vous êtes non seulement d'excellentes scientifiques, mais aussi de superbes encadrantes.

I would also like to thank the reviewers of this thesis, Prof. Vianney Perchet and Prof. Wouter Koolen for the time they devoted to evaluating my work. I am grateful for their careful reading, insightful comments, and constructive suggestions, which helped improve both the manuscript and its presentation. I also thank all the members of the examination committee for their participation and for the stimulating discussions.

I am grateful to my colleagues and friends at Inria, whose presence and collaboration made this PhD an enriching scientific and human experience: Hector, Brahim, Waris, Marc, Achraf, Riad, Timothée, Adrienne, Sabrine, Rémy, Tuan, Yann, Mahdi, Udvas, Thomas, Sumit, Fabien, Patrick, Dorian, Matheus, Ayoub, Anthony, Deb, Julien, Amélie, Odalric, Naafi, Alena, Hernan, Riccardo, Mohamed Yacine, Reda, Nathan, Redouane, Nicolas, Hadrien, Adrien, Juliette, Lorenzo, Hatim, Guillaume, Tanguy, Brell, Mericel, Mickaël, and Philippe. I thank them for the many discussions, their support, and the constructive working environment that accompanied this research. Beyond the scientific exchanges, I particularly appreciated the Scool team's atmosphere: a rare combination of intellectual rigour, generosity, and genuine camaraderie. Thank you for the unforgettable pétanque parties, badminton sessions, and football matches, which provided welcome opportunities to unwind and recharge during these years. I feel fortunate to have been surrounded by such talented researchers and colleagues; your curiosity, energy, and professionalism have contributed greatly to my development throughout this journey. Finally, I would like to thank Philippe, as team

director, for fostering and sustaining this exceptional environment.

I would like to offer special thanks to Riad, Brahim, Waris, Hector, and Tuan for the many moments shared outside the laboratory, including our culinary explorations in Lille, which remain among my most memorable experiences during these years.

Je suis également reconnaissant aux membres de l'équipe SISTM à Bordeaux. Merci en particulier à Rodolphe pour son implication constante et la qualité de nos échanges lors de mes séjours à Bordeaux. Merci aussi à Quentin, Mélanie, Boris, Myrtille et Linda pour ces discussions stimulantes et enrichissantes autour des problématiques liées aux essais cliniques, qui ont nourri ma compréhension de ces sujets et donné une profondeur supplémentaire à ce travail.

I would also like to thank Arnab, Nicolas and Kevin at the University of Washington for the many fruitful exchanges during my visit to Seattle.

À mes amis de longue date — Abdoulaye, Ibrahim, Paul, Arthur, Milacet, Christian et Simplicie — je vous remercie chaleureusement pour votre présence constante, vos encouragements et votre soutien. Je repense avec gratitude à tous les beaux moments qui ont ponctué ces trois années de thèse: vos messages, nos discussions, les retrouvailles, les rires, et cette manière si naturelle que vous avez de remettre les choses en perspective. Votre amitié m'a apporté stabilité et réconfort dans les périodes exigeantes, et a rendu ce parcours bien plus léger et plus joyeux.

À ma famille, je tiens à exprimer ma profonde reconnaissance, en particulier envers mes frères et sœurs — Julie, Estelle, Anicet — pour leur affection, leur patience et leurs encouragements tout au long de ce parcours.

Je souhaite adresser un remerciement tout particulier à mon oncle Nanontia Koné, dont le soutien indéfectible, depuis mon enfance, a été déterminant. Cette thèse n'aurait pas été possible sans sa générosité, sa confiance et sa présence constante. C'est grâce à l'ordinateur familial qu'il a acquis en 2008 que j'ai découvert l'informatique ; depuis, mon attachement à cette discipline n'a jamais cessé de grandir. Je lui en suis infiniment reconnaissant.

Enfin, je souhaite rendre hommage à mes parents, trop tôt disparus, Émilie Touré et Koné Niéna. Leur mémoire a été, durant toutes ces années, une source de motivation, et je leur dédie ce travail. Je ne pourrais ouvrir ce manuscrit sans rappeler combien je le dois à ma mère. Elle m'a transmis très tôt le goût des études et la fierté d'apprendre. Je me souviens avec émotion de la joie qui était la sienne lorsque je rapportais mes premiers résultats à l'école primaire ; c'est le souvenir de cette joie qui m'a porté tout au long de ce parcours, dans les moments d'élan comme dans les périodes de doute. Et malgré mille détours, les épreuves et les turbulences, je garde l'espoir que ce travail soit, à sa manière, le prolongement de ce que vous avez semé en moi.

Abstract

In this dissertation, we study a class of *multi-objective pure exploration problems* in stochastic bandits. Our goal is to design algorithms that identify the set of *Pareto-optimal* arms—those that cannot be improved in one objective without deteriorating another—while minimizing the number of samples required to achieve a prescribed level of confidence.

We begin by reviewing the literature on classical and multi-objective bandits, outlining the key theoretical tools that motivate the study of *Pareto Set Identification* (PSI). Throughout the thesis, we progressively move from simple algorithms based on confidence sequences with finite-time guarantees to more advanced approaches that achieve asymptotic optimality, while maintaining computational tractability. All proposed methods are systematically evaluated on synthetic and real-world-inspired datasets to provide a comprehensive understanding of their empirical performance.

In the fixed-budget regime, we develop EGE (*Empirical Gap Elimination*), the first algorithm for PSI in this setting. In the fixed-confidence regime, we introduce APE, a δ -correct procedure based on adaptive confidence intervals. Applications to multi-endpoint vaccine trials illustrate how PSI can guide early-stage clinical decision-making by identifying treatment strategies that jointly optimize multiple immunological indicators.

Chapter 4 extends PSI to the *structured* case, where arm means are linearly related through an unknown parameter matrix θ and known feature representations. We derive problem-dependent lower bounds on the sample complexity of adaptive algorithms and we propose efficient algorithms that exploit this structure to significantly reduce the number of required samples.

Chapter 5 focuses on algorithms that achieve the *information-theoretic lower bound*. We introduce PSIPS (*Pareto Set Identification with Posterior Sampling*), the first computationally efficient and asymptotically optimal algorithm for PSI with correlated objectives. PSIPS uses posterior sampling for both the sampling and stopping rules, eliminating the need to predefine the confidence parameter δ as a parameter of the sampling rule. It achieves theoretical optimality and demonstrates excellent empirical performance in extensive experiments, including on the *COV-BOOST* vaccine dataset, while being several orders of magnitude faster than gradient-based approaches.

Finally, Chapter 6 explores PSI under *linear constraints*, where the goal is to identify the set of Pareto-optimal arms among those that satisfy feasibility conditions. We propose e-CAPE, an extension of APE, and characterize the information-theoretic complexity of constrained PSI.

Overall, this work advances the theoretical and algorithmic foundations of *multi-objective pure exploration*. By bridging sequential decision-making, information theory, optimal experimental design, and Bayesian principles, it demonstrates that one can achieve both theoretical optimality and computational efficiency. Beyond bandits, the insights and methods developed here open promising avenues toward a general theory of *multi-objective sequential learning*, with applications to adaptive clinical trials, reinforcement learning, and multi-criteria decision-making under uncertainty.

Keywords: Multi-armed bandits, Pure exploration, Multi-objective, Pareto, Sequential learning, Clinical trials

Résumé

Cette thèse porte sur l'exploration pure à objectifs multiples dans le cadre des bandits stochastiques. L'objectif est de concevoir des algorithmes capables d'identifier l'ensemble des bras Pareto-optimaux; ceux qu'il est impossible d'améliorer sur un objectif sans en détériorer un autre, tout en minimisant le nombre d'échantillons nécessaires pour atteindre un niveau de confiance donné.

Après une revue de la littérature sur les bandits classiques et multiobjectifs, nous introduisons le problème d'Identification de l'Ensemble de Pareto (IEP) et les outils théoriques qui le sous-tendent. Nous progressons des méthodes basées sur des intervalles de confiance offrant des garanties en temps fini vers des approches asymptotiquement optimales, tout en préservant l'efficacité computationnelle. Les algorithmes proposés sont évalués de manière systématique sur des jeux de données synthétiques et réels, afin d'analyser finement leurs performances empiriques.

La première partie présente l'algorithme EGE, qui identifie efficacement l'ensemble de Pareto lorsque le budget d'échantillons est fixé. Revisitant les travaux sur l'IEP à confiance fixée, nous introduisons APE, une procédure adaptative δ -correcte reposant sur des intervalles de confiance. APE fournit de fortes garanties en nombre fini d'échantillons et de très bonnes performances pratiques. Des applications à des jeux de données issus d'essais vaccinaux multicritères illustrent comment EGE et APE peuvent guider la sélection de stratégies thérapeutiques optimisant simultanément plusieurs marqueurs immunologiques.

Nous étendons ensuite l'IEP au cas structuré, où les moyennes des bras dépendent linéairement de descripteurs connus via une matrice de régression inconnue. L'algorithme GESE exploite cette structure pour réduire la complexité d'échantillonnage et atteint des bornes inférieures spécifiques à l'instance de bandit.

Le chapitre suivant identifie les limites de ces approches et introduit PSIPS, premier algorithme à la fois asymptotiquement optimal et à faible coût de calcul pour l'IEP avec objectifs corrélés. Basé sur un échantillonnage a posteriori pour le choix des bras et la règle d'arrêt, PSIPS élimine la nécessité de fixer à l'avance le paramètre de confiance δ . Il atteint la borne inférieure théorique sur la complexité d'échantillonnage et se montre très compétitif, y compris sur des données réelles comme COV-BOOST, tout en étant bien plus rapide que les méthodes à base de gradient.

Enfin, le dernier chapitre traite du cas contraint, où seuls les bras respectant certains critères doivent être considérés pour le calcul de l'ensemble de Pareto. Un algorithme doit donc apprendre simultanément les bras viables et déterminer leur optimalité. Nous proposons e-CAPE, un algorithme qui généralise celui du Chapitre 3 et établit des bornes théoriques pour l'IEP sous contraintes.

Dans l'ensemble, cette thèse fait progresser les fondements théoriques et algorithmiques de l'exploration pure à objectifs multiples. En combinant les outils issus de statistiques, apprentissage séquentiel, conception expérimentale optimale et principes bayésiens, elle montre qu'il est possible de concilier optimalité théorique et efficacité pratique. Les idées développées ouvrent la voie à une théorie générale de l'apprentissage séquentiel à objec-

tifs multiples, avec des applications aux essais cliniques adaptatifs, à l'apprentissage par renforcement et plus globalement à l'apprentissage multicritère sous incertitude.

Mots-clés : Bandits manchots, Exploration pure, Multi-objectif, Pareto, Apprentissage séquentiel, Essais cliniques

Table of Contents

1	Introduction	1
1.1	Bandit Algorithms for Early-Stage Clinical Trials	2
1.2	Active Testing (Pure Exploration)	5
1.2.1	Stopping and Recommendation Rules	8
1.2.2	Sampling Rules	10
1.2.3	Lower Bounds	11
1.3	Background	11
1.3.1	Best Arm Identification and Beyond	12
1.3.2	Multi-objective Pure Exploration	14
1.4	Contributions of This Dissertation	19
1.4.1	Fixed-Budget PSI	19
1.4.2	Adaptive Fixed-Confidence Algorithms for PSI	21
1.4.3	PSI with Linear Structure	22
1.4.4	Asymptotically Optimal Algorithms for Fixed-Confidence PSI	23
1.4.5	PSI with Linear Feasibility Constraints	25
1.5	Outline of the Dissertation	26
1.6	Publications	27
2	Pareto Set Identification: The Fixed-Budget Setting	29
2.1	Introduction	30
2.2	Algorithmic contributions	33
2.2.1	Empirical Gap Elimination	33
2.2.2	Particular instances	35
2.2.3	Beyond exact PSI	36
2.3	Main theoretical results	37
2.3.1	Upper bound on error probability	37
2.3.2	Lower bound on the error probability	39
2.3.3	Sketch of proof of Theorem 2.3.1	40
2.4	Numerical study and discussion	45
2.5	Additional proofs	47
2.5.1	Analysis of EGE	48
2.5.2	Analysis of EGE-SR- k	53
2.5.3	Simplifying the sub-optimality gaps	58
3	Adaptive Algorithms for Fixed-Confidence Pareto Set Identification	63
3.1	Introduction	64
3.2	Adaptive Pareto Exploration	66
3.2.1	Generic algorithm(s)	67
3.2.2	Our instantiation	71
3.3	Main theoretical results	73

3.3.1	Sketch of proof of Theorem 3.3.1	74
3.4	Numerical study and discussion	76
3.5	Additional proofs	80
3.5.1	Probability of the good event	80
3.5.2	Sample complexity	81
4	Linear Model for Pareto Set Identification	87
4.1	Introduction	88
4.1.1	Complexity measures for Pareto Set Identification	90
4.1.2	Least-squares estimation and optimal designs	91
4.2	Algorithmic contribution	92
4.2.1	Optimal designs and gap estimation	92
4.2.2	Fixed-budget algorithm	96
4.3	Main theoretical results	96
4.3.1	Fixed-budget	96
4.3.2	Fixed-confidence	97
4.3.3	Sketch of proofs	99
4.4	Numerical study and discussion	100
4.5	Additional proofs	102
4.5.1	Empirical gaps: Proof of Proposition 4.3.6	103
4.5.2	Probability of error: Proof of Theorem 4.3.1	108
4.5.3	Sample complexity: Proof of Theorem 4.3.3	111
4.5.4	Lower bounds	116
4.5.5	Concentration lemmas	118
5	Posterior Sampling for Pareto Set Identification	121
5.1	Introduction	122
5.2	PSI with Posterior Sampling	125
5.2.1	The Posterior Sampling (PS) Stopping Rule	125
5.2.2	Game-based sampling rule	128
5.3	Main theoretical results	129
5.3.1	Sketch of proofs	130
5.4	Numerical study and discussion	134
5.5	Additional proofs	137
5.5.1	Stopping rule	137
5.5.2	Sample complexity	146
5.5.3	Posterior convergence	149
6	Pareto Set Identification with Linear Constraints	155
6.1	Introduction	156
6.2	On the complexity of constrained PSI	159
6.2.1	Constrained PSI without explainability (cPSI)	159
6.2.2	Constrained PSI with Explainability (e-cPSI)	160
6.3	Constrained Adaptive Pareto Exploration	161
6.4	Main theoretical results	165
6.4.1	Sample complexity	167
6.5	Numerical study and discussion	170

6.6	Additional proofs	174
6.6.1	Stopping time	174
6.6.2	Sample complexity: Proof of Theorem 6.4.1	182
6.6.3	Technical results	185
6.6.4	Lower bound of e-cPSI	186
6.6.5	Complexity of Best response for cPSI	192
7	Conclusion and Perspectives	197
	List of Datasets	201
	List of Notation	205
	List of Figures	216
	List of Tables	217
	List of Algorithms	219
	Bibliography	219

Chapter 1

Introduction

This doctoral research was conducted between November 2022 and October 2025 within the Scool team at Inria Lille, hosted by Inria and the CRISAL Computer Science Laboratory of the University of Lille. It was carried out under the supervision of Dr. Émilie Kaufmann (CNRS/Inria Lille) and Prof. Laura Richert (Inserm Bordeaux / University of Bordeaux), thanks to joint Inria–Inserm funding. As part of this collaboration, I made regular visits to Prof. Richert’s group in Bordeaux, which provided valuable opportunities for interdisciplinary exchanges between machine learning and clinical-trial design.

The manuscript is organized as follows. Each chapter studies a particular aspect of Pareto set identification, beginning with confidence-based algorithms offering finite-time guarantees, and gradually extending to structured and constrained formulations. The work culminates in the design of asymptotically optimal algorithms based on posterior sampling, combining theoretical guarantees with computational efficiency. Finally, the empirical sections validate these approaches on both synthetic and real-world-inspired datasets, and the thesis concludes with perspectives on future research in multi-objective pure exploration.

1.1	Bandit Algorithms for Early-Stage Clinical Trials	2
1.2	Active Testing (Pure Exploration)	5
1.2.1	Stopping and Recommendation Rules	8
1.2.2	Sampling Rules	10
1.2.3	Lower Bounds	11
1.3	Background	11
1.3.1	Best Arm Identification and Beyond	12
1.3.2	Multi-objective Pure Exploration	14
1.4	Contributions of This Dissertation	19
1.4.1	Fixed-Budget PSI	19
1.4.2	Adaptive Fixed-Confidence Algorithms for PSI	21
1.4.3	PSI with Linear Structure	22
1.4.4	Asymptotically Optimal Algorithms for Fixed-Confidence PSI	23
1.4.5	PSI with Linear Feasibility Constraints	25
1.5	Outline of the Dissertation	26
1.6	Publications	27

1.1 Bandit Algorithms for Early-Stage Clinical Trials

Vaccination remains one of the most effective strategies to prevent and control infectious diseases. However, the clinical development of vaccines—like that of most experimental drugs—remains a long, costly, and uncertain process. The accelerated development of first-generation SARS-CoV-2 vaccines during the COVID-19 pandemic was an exceptional case of global mobilization and investment; in most contexts, vaccine development proceeds under much tighter resource and time constraints. At the beginning of a vaccine’s clinical development plan, early to intermediate trials (phases I and II) are conducted to evaluate safety and explore immunogenicity in a limited number of participants. These early-stage trials play a crucial role: they must identify, among a potentially large set of vaccine candidates, those most promising for further efficacy testing in large phase IIb/III studies, which are extremely costly and logistically demanding.

From clinical motivation to bandit models. Adaptive clinical trial designs are one of the most promising approaches to accelerate and optimize vaccine development, as they offer a higher degree of adaptivity and statistical efficiency than traditional methods. The stochastic multi-armed bandit (MAB) problem was introduced as a stylized abstraction of an *adaptive clinical trial* (Thompson 1933; H. E. Robbins 1952), where a physician must sequentially allocate patients to competing treatments while learning from observed outcomes. Each treatment, or *arm*, may correspond, for instance, to a particular combination of dose, adjuvant, injection interval, or delivery platform. Patient responses are modeled as independent and identically distributed (*i.i.d.*) samples from the arm received by the patient.

Upon observing the responses of treated participants, the algorithm adaptively updates its allocation strategy, typically to balance exploration of uncertain options with exploitation of the most promising ones. This abstraction directly mirrors adaptive trial design: the aim is to identify promising arms as efficiently as possible, while respecting sample size and ethical constraints.

Initially, this abstraction was developed for *phase III* clinical trials, where treatments are tested for efficacy on cohorts of patients, often including patients affected by the target disease. In such settings, ethical considerations are paramount: since enrolled patients may derive therapeutic benefit from the experimental intervention, one seeks to minimize the number of individuals exposed to suboptimal treatments. This perspective aligns naturally with the *regret-minimization* formulation in bandit theory (H. E. Robbins 1952; T. Lai & H. Robbins 1985; Auer 2003; Garivier & Cappé 2011; Agrawal & Goyal 2012; Kaufmann, Korda, et al. 2012; Lattimore & Szepesvari 2020; Russo et al. 2020), where the goal is to balance learning (reducing uncertainty about treatment effects) and earning (maximizing the number of patients cured).

Early-stage vaccine development. In contrast, *early-phase trials* (typically phases I and II) pursue a different goal: rather than maximizing immediate patient outcomes, they aim to



Figure 1.1: Illustration of the multi-arm trial model

collect enough evidence to assess the safety, tolerability, or immunogenicity of candidate vaccines or treatments, often on healthy volunteers. Maximizing cumulative outcomes during the trial is secondary, and the focus shifts to accurate final recommendations. These trials involve small cohorts and limited budgets.

Adaptive allocation is particularly useful in this setting as it allows learning from early immunological responses to guide subsequent patient assignments, improving both statistical efficiency and ethical acceptability. This dissertation focuses on this side of the problem, which corresponds to the *pure exploration* setting in the literature on multi-armed bandits (Bubeck, Munos, et al. 2009; Degenne 2019). Here, treatments are allocated to maximize information gain rather than patient outcomes, and success is measured by the probability of correctly identifying effective treatments at the end of the trial.

Moreover, vaccine strategies for novel pathogens often yield a large number of experimental arms: multiple formulations, dose levels, injection intervals, booster strategies, and adjuvant combinations may all be evaluated concurrently. The COVID-19 booster-dose study (Munro et al. 2021), for example, simultaneously compared several booster vaccines across hundreds of participants—a design that naturally resembles a large-scale multi-armed bandit problem.

The multi-objective challenge. Most adaptive algorithms and trial designs treat treatment quality as one-dimensional. However, at the early stages of development, no validated *correlate of protection*, that is a single predictive immune marker of efficacy, is typically available (Plotkin & Gilbert 2012). Instead, several immunological readouts are monitored simultaneously (e.g., antibody titers, neutralizing activity, T-cell response, and reactogenicity), and these may not align. In such settings, the notion of a single “best” treatment is ill-defined: one candidate may induce stronger antibodies but higher reactogenicity (Mark C et al. 2013), while another elicits broader T-cell responses but weaker neutralizing titers (Munro et al. 2021). Such trade-offs imply that no single vaccine candidate dominates all others. Instead, the goal we pursue in this dissertation is to identify the subset of *non-dominated* treatments forming the *Pareto front* or *Pareto set*: those that are not strictly worse across all measured criteria.

Formally, if $\mu_i := (\mu_i^1, \dots, \mu_i^d)$ denotes the mean response vector for arm i , we say that arm i is *dominated* by j (denoted $i \preceq j$ or $\mu_i \preceq \mu_j$) if $\mu_i^c \leq \mu_j^c$ for all $c \in \{1, 2, \dots, d\}$, with strict inequality in at least one dimension. The goal is then to identify the set of treatments that are not dominated by any other: the *Pareto-optimal set*.

The COVID-19 example. For instance, in the early stages of COVID-19 vaccine development, several immune correlates—such as binding antibody titers, neutralizing activity, and T-cell responses—were simultaneously monitored to characterize immunogenicity before a definitive correlate of protection was established (Munro et al. 2021). The study of Munro et al. 2021 simultaneously evaluated multiple COVID-19 booster strategies, while measuring several immunological endpoints. In such cases, each vaccine candidate is associated with multiple response dimensions, and no single metric fully captures efficacy. Such designs naturally give rise to many arms (representing distinct vaccine strategies), a limited number of participants, and multiple partially conflicting measures of efficacy. This example illustrates precisely the kind of design challenges that motivate this work.

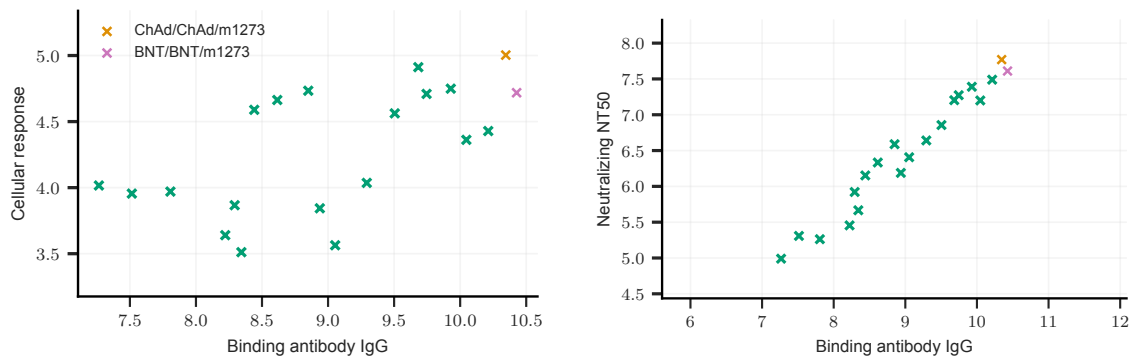


Figure 1.2: Empirical (arithmetic) average of the log-transformed immune response for three immunogenicity indicators reported by Munro et al. 2021. Each acronym corresponds to a vaccine. There are two groups of arms corresponding to the first 2 doses: one with prime BNT/BNT (BNT as first and second dose) and the second with prime ChAd/ChAd (ChAd as first and second dose). Combining the first two doses received and the third dose investigated in the trial, there were $K = 20$ arms.

Formal problem statement. We model this setting as a stochastic multi-armed bandit with vector-valued outcomes. At each time $t = 1, 2, \dots, T$, a treatment A_t is adaptively assigned to the arriving participant, and a random vector $Z_t = (Z_t^1, \dots, Z_t^d)$ is observed. Each arm k corresponds to an unknown distribution ν_k over \mathbb{R}^d with mean vector $\mu_k \in \mathbb{R}^d$; the outcome Z_t is modeled as a sample from the distribution ν_{A_t} . These distributions belong to some family \mathcal{P} . The goal is to identify the Pareto-optimal subset $\mathcal{S}^*(\mu) := \{i \in [K] : \nexists j, \mu_i \preceq \mu_j\}$ as efficiently as possible: for example by minimizing the probability of misidentification of optimal vaccine strategies using only a few patients (formalized in Section 1.2). In the example of Figure 1.2, two vaccine strategies are Pareto-optimal.

Dissertation scope and objectives. This dissertation develops adaptive algorithms tailored to the design of early-stage, multi-objective clinical trials. The focus is on the *pure exploration* setting, where the goal is to identify the best (or non-dominated) treatments at the end of the trial rather than to maximize patient benefit during the trial. Such designs are especially relevant to vaccine development, where the number of participants is small, the arms are many, and multiple endpoints must be considered simultaneously.

Throughout the manuscript, in particular in the experiments, we will frequently refer to the COVID-19 booster-dose trial (Munro et al. 2021) as a motivating application. It epitomizes the challenges addressed in this work: large numbers of candidate arms, multiple immunological endpoints, limited resources. We build upon and extend recent advances in the theory of bandit algorithms to handle:

- multi-dimensional outcomes and Pareto set identification,
- linear or feasibility constraints between arms, and
- Bayesian formulations and randomized algorithms based on posterior sampling.

The resulting algorithms provide statistical guarantees (e.g., controlled error probability, sample efficiency) and their relevance is shown through simulations, datasets inspired by real-world vaccine studies and other benchmarks from the literature on multi-objective bandits.

Positioning with respect to the literature. For phase IIb/III trials where the patients should be cured during the trial, the field of regret minimization in bandits is now well developed, with extensions to non-stationary environments, contextual and adversarial settings, and structured action spaces (Lattimore & Szepesvari 2020). Upper Confidence Bound (UCB) algorithms, introduced by Auer 2003, and their variants, KL-UCB (Garivier & Cappé 2011), as well as randomized Bayesian algorithms such as Thompson Sampling (Thompson 1933), were shown to combine strong empirical performance with rigorous guarantees (Agrawal & Goyal 2012; Kaufmann, Korda, et al. 2012; Russo et al. 2020).

For early vaccine trials that focus on identifying the most immunogenic or balanced strategies, existing work on bandit algorithms provides a rich theoretical foundation, ranging from fixed-confidence best-arm identification (Jamieson & Nowak 2014; Garivier & Kaufmann 2016; Shang et al. 2020; Jourdan, Degenne, Baudry, et al. 2022) to fixed-budget pure exploration (Audibert & Bubeck 2010; Karnin et al. 2013; Locatelli et al. 2016). Yet, few studies have addressed the case of vector-valued outcomes or Pareto set identification, besides the work of (Auer et al. 2016). This dissertation builds precisely on that gap by proposing algorithms tailored to multi-objective exploration problems.

The next section introduces the general *pure exploration* framework underlying this work, detailing sampling, stopping, and recommendation principles.

1.2 Active Testing (Pure Exploration)

Pure exploration, also called *active testing*, models learning scenarios in which an algorithm is judged only by its final recommendation (Bubeck, Munos, et al. 2009; Degenne & W. Koolen 2019). Unlike regret minimization, interim outcomes are irrelevant; the goal is to collect just enough information to answer a prespecified query with high confidence.

Bandit model and query task. A bandit instance is described by $\nu := (\nu_1, \dots, \nu_K)$ with unknown mean vectors $\mu := (\mu_1, \dots, \mu_K) \in \mathcal{I}^K$ (scalar outcomes: $d = 1$; multi-objective: $\mu_k \in \mathbb{R}^d$). A *query task* is a (possibly set-valued) mapping

$$\mathcal{S}^* : \mathcal{I}^K \rightarrow 2^{[K]},$$

whose value $\mathcal{S}^*(\mu)$ collects all acceptable answers (e.g., $\{\arg \max_k \mu_k\}$ in BAI; the Pareto set in PSI). Any element of $\mathcal{S}^*(\mu)$ is correct.

Example 1.1 (Best Arm Identification). For best arm identification, the means are scalar-valued, $\mu = (\mu_1, \dots, \mu_K) \in \mathbb{R}^K$, and the query task associates each $\mu \in \mathbb{R}^K$ to the index of the best arm: $\mathcal{S}^*(\mu) = \{k \in [K] : \mu_k \geq \max_{i \neq k} \mu_i\}$.

Example 1.2 (Pareto Set Identification). For Pareto set identification, the mean of each arm is vector-valued, $\mu = (\mu_1, \dots, \mu_K) \in (\mathbb{R}^d)^K$, and the goal is to identify the set of arms that are not uniformly dominated on all d objectives by another arm. Formally, $\mathcal{S}^*(\mu) = \{k \in [K] : \nexists i \in [K] \setminus \{k\}, \mu_k \preceq \mu_i\}$, where $\mu_i \preceq \mu_j$ means that for all objectives $c \in \{1, \dots, d\}$, $\mu_i^c \leq \mu_j^c$, with strict inequality for at least one objective. This is further illustrated in Figure 1.3.

In this section, we focus on problems with a unique correct answer, *i.e.*, for a given μ , the task operator $\mathcal{S}^*(\mu)$ is a singleton-valued (its unique element may itself be a set, as in PSI). We then write $\mathcal{S}^*(\mu) \neq S$ to denote without ambiguity $S \notin \mathcal{S}^*(\mu)$, and similarly, we write $S = \mathcal{S}^*(\mu)$ to denote the unique element of $\mathcal{S}^*(\mu)$.

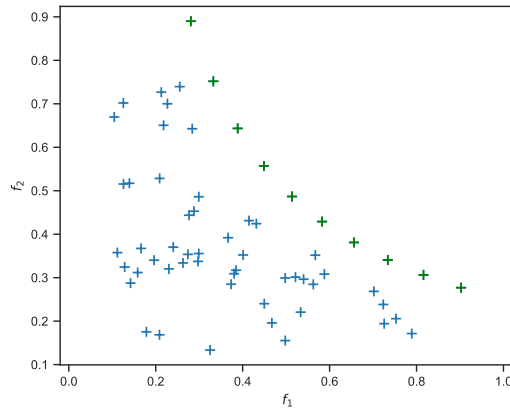


Figure 1.3: Pareto set in a bi-objective setting. Green points are Pareto-optimal, while blue points are sub-optimal.

At round $t \geq 1$, the learner chooses an arm $A_t \in [K]$, observes $Z_t \sim \nu_{A_t}$, and updates based on the history $\mathcal{H}_t = \sigma(A_1, Z_1, \dots, A_t, Z_t)$.

Policies and performance criteria. A pure exploration policy (or learner, or agent, or algorithm) has three main components: (i) a *sampling rule* $(A_t)_{t \geq 1}$ that dictates the arm to sample at each round, (ii) a *recommendation rule* \hat{S}_t (a random element or subset of $[K]$) representing a guess for the correct answer, and (iii) a *stopping rule* τ (a stopping time *w.r.t.*

\mathcal{H}_t), which dictates when to stop collecting new samples from the environment. The final output is \widehat{S}_τ .

For risk $\delta \in (0, 1)$, a policy is δ -PAC (or δ -correct) if, for every instance with mean μ ,

$$\mathbb{P}_\nu(\tau < \infty, \widehat{S}_\tau \neq \mathcal{S}^*(\mu)) \leq \delta. \quad (1.1)$$

Two classical formulations exist:

- **Fixed-confidence (FC):** the goal is to minimize the *sample complexity*, i.e., the (instance-specific) number of samples used before stopping, while satisfying (1.1) (Even-Dar et al. 2002; Kalyanakrishnan et al. 2012; Jamieson & Nowak 2014; Degenne, Ménard, et al. 2020). When the algorithm draws exactly one sample per round (or a fixed number per round), this is equivalent to minimizing the stopping time τ (or its expected value $\mathbb{E}[\tau]$).
- **Fixed-budget (FB):** Audibert & Bubeck 2010; Locatelli et al. 2016; Katz-Samuels & Scott 2018; Roy Chaudhuri & Kalyanakrishnan 2019 for a given horizon T , stop at $\tau := T$ and minimize $e_T(\nu) := \mathbb{P}_\nu(\widehat{S}_T \neq \mathcal{S}^*(\mu))$; interest centers on the decay rate of $e_T(\nu)$.

In short, the fixed-budget setting emphasizes optimal *allocation* (stopping is fixed), whereas the fixed-confidence setting couples *adaptive sampling* with *adaptive stopping*.

Algorithm 1.1: Generic Pure Exploration Bandit Algorithm

Require: Arms $[K] = \{1, \dots, K\}$; risk δ (FC) or budget T (FB)

- 1 Initialize counts $N_k \leftarrow 0$, for all $k \in [K]$; $t \leftarrow 0$.
 - 2 **repeat**
 - 3 $t \leftarrow t + 1$.
 - 4 // Sampling rule (policy)
 - 4 Select arm $A_t \leftarrow \mathbf{SamplingRule}(\mathcal{H}_{t-1})$.
 - 5 Observe outcome $Z_t \sim \nu_{A_t}$.
 - 6 Update $N_{A_t} \leftarrow N_{A_t} + 1$ and estimates $\hat{\mu}_{t,A_t}$ (and confidence intervals/posterior as needed).
 - 7 // Recommendation rule (current guess)
 - 7 $\widehat{S}_t \leftarrow \mathbf{RecommendationRule}(\{\hat{\mu}_{t,k}\}_{k=1}^K, \mathcal{H}_t)$ (e.g. empirical best, empirical Pareto set, etc.).
 - 8 **until** $\mathbf{StoppingRule}(\mathcal{H}_t, \delta)$ is true (FC) or $t = T$ (FB)
 - 9 **return** final recommendation \widehat{S}_τ .
-

1.2.1 Stopping and Recommendation Rules

The recommendation rule maps the observed data to an element of $\{\mathcal{S}^*(\lambda) : \lambda \in \mathcal{I}^K\}$. The most frequent recommendation rules are:

- Plug-in estimate: use empirical data to estimate the answer, e.g., $\widehat{S}_t = \arg \max_k \widehat{\mu}_{t,k}$ in BAI, or $\widehat{S}_t = \{i : \nexists j, \widehat{\mu}_{t,i} \preceq \widehat{\mu}_{t,j}\}$ in PSI.
- Confidence-based: return the most “optimistic” answer e.g., the arm with the largest upper confidence bound (Gabillon et al. 2012).
- Bayesian: greedy response w.r.t. posterior sample, most probable answer as measured by the posterior (Russo 2016; Shang et al. 2020)

In practice, the recommendation is coupled with stopping: the algorithm halts only when the chosen recommendation is certified correct by the stopping rule.

Stopping rules are (Markov) stopping times w.r.t. the filtration $\{\mathcal{H}_t\}_{t=1}^\infty$ and determine when the accumulated evidence suffices to recommend a correct answer with the desired confidence. Designing efficient and computationally tractable *stopping rules* is a central challenge in the construction of bandit algorithms. Below, we unify frequentist and Bayesian constructions. Before introducing them, let $\text{Alt} : 2^{[K]} \rightarrow (\mathbb{R}^d)^K$ defined as

$$\text{Alt}(\mathcal{S}^*(\mu)) := \{\lambda \in \mathcal{I}^K : \mathcal{S}^*(\lambda) \neq \mathcal{S}^*(\mu)\}$$

be the set of parameters that have a different correct response than μ .

An *anytime confidence region* $(\mathcal{R}_t(\delta))_{t \geq 1}$ at level δ satisfies $\mathbb{P}(\forall t \geq 1, \mu \in \mathcal{R}_t(\delta)) \geq 1 - \delta$. With an estimator $\widehat{\mu}_t$ of μ , consider the generic stopping rule

$$\tau := \inf\{t \geq 1 : \mathcal{R}_t(\delta) \cap \text{Alt}(\widehat{S}_t) = \emptyset\}. \quad (1.2)$$

This yields δ -PAC correctness, since

$$\begin{aligned} \mathbb{P}(\tau < \infty, \widehat{S}_\tau \neq \mathcal{S}^*(\mu)) &\leq \mathbb{P}(\exists t : \mu \in \text{Alt}(\widehat{S}_t), \mathcal{R}_t(\delta) \cap \text{Alt}(\widehat{S}_t) = \emptyset) \\ &\leq \mathbb{P}(\exists t : \mu \notin \mathcal{R}_t(\delta)) \leq \delta. \end{aligned}$$

In words, on the event $\{\forall t \geq 1, \mu \in \mathcal{R}_t(\delta)\}$, the true parameter is never excluded before stopping, and at τ all parameters consistent with data agree on the same answer.

Depending on the geometry of the confidence regions, various forms of stopping rules can be derived from Eq. (1.2).

Individual confidence sequences. For scalar-valued problems, if each arm distribution ν_k is σ -subgaussian, i.e., $\log \mathbb{E}_{X \sim \nu_k}[e^{\lambda(X - \mu_k)}] \leq \sigma^2 \lambda^2 / 2$ for all $\lambda \in \mathbb{R}$, then the empirical mean $\widehat{\mu}_{t,k}$ of arm k after $N_{t,k}$ samples admits an anytime *Hoeffding-type* confidence interval:

$$[L_k(t, \delta), U_k(t, \delta)] = \left[\widehat{\mu}_{t,k} - \sqrt{\frac{2\sigma^2 \log(Ct^2 K / \delta)}{N_{t,k}}}, \widehat{\mu}_{t,k} + \sqrt{\frac{2\sigma^2 \log(Ct^2 K / \delta)}{N_{t,k}}} \right],$$

valid with probability at least $1 - \delta/K$ for a universal constant $C > 0$. The Cartesian product $\mathcal{R}_t(\delta) = \prod_{k=1}^K [L_k(t, \delta), U_k(t, \delta)]$ defines an *anytime confidence region* for μ .

For the best-arm identification (BAI) query, substituting \mathcal{R}_t into the generic stopping rule (1.2) yields the classical LUCB-style stopping time:

$$\tau_\delta = \inf \left\{ t \geq 1 : L_{\hat{S}_t}(t, \delta) \geq U_k(t, \delta), \forall k \neq \hat{S}_t \right\},$$

where $\hat{S}_t = \arg \max_k \hat{\mu}_{t,k}$ is the empirical best arm. At stopping, every $\lambda \in \mathcal{R}_t(\delta)$ satisfies $\lambda_{\hat{S}_t} \geq \max_k \lambda_k$, ensuring that the empirical best arm is also the true best arm for all mean vectors consistent with the collected data.

When the arms belong to a one-parameter exponential family, tighter bounds are obtained by replacing Hoeffding radii with Kullback-Leibler (KL) constraints. Let $d(x, y)$ denote the KL divergence between ν^x and ν^y for distributions in a parametric family indexed by x, y . The *KL-based confidence interval* for an arm k is defined implicitly by

$$\begin{aligned} U_k(t, \delta) &= \sup \{ q \geq \hat{\mu}_{t,k} : N_{t,k} d(\hat{\mu}_{t,k}, q) \leq f(t, \delta) \}, \\ L_k(t, \delta) &= \inf \{ q \leq \hat{\mu}_{t,k} : N_{t,k} d(\hat{\mu}_{t,k}, q) \leq f(t, \delta) \}, \end{aligned}$$

where $f(t, \delta)$ is a time-dependent exploration function of order $\log(Ct/\delta)$. The corresponding product region $\mathcal{R}_t(\delta) = \prod_k [L_k(t, \delta), U_k(t, \delta)]$ plugged into (1.2) yields KL-LUCB-type stopping rules (Kaufmann & Kalyanakrishnan 2013), offering sharper, distribution-dependent confidence bounds and improved finite-time performance.

More refined confidence regions can be constructed from sums of KL or likelihood bounds, yielding Chernoff or GLR-type stopping rules such as those used in Garivier & Kaufmann 2016.

Joint KL confidence region/Chernoff stopping. A global (joint) confidence region aggregates evidence across arms:

$$\mathcal{R}_t(\delta) = \left\{ \lambda \in \mathcal{I}^K : \sum_{k=1}^K N_{t,k} d(\hat{\mu}_{t,k}, \lambda_k) \leq \beta(t, \delta) \right\},$$

with β chosen such that $\mathbb{P}(\forall t, \mu \in \mathcal{R}_t(\delta)) \geq 1 - \delta$. This region captures the collective statistical evidence about the parameter μ . Then Eq. (1.2) becomes

$$\begin{aligned} \tau &= \inf \left\{ t \geq 1 : \forall \lambda \in \text{Alt}(\hat{S}_t), \sum_{k=1}^K N_{t,k} d(\hat{\mu}_{t,k}, \lambda_k) > \beta(t, \delta) \right\} \\ &= \inf \left\{ t \geq 1 : \inf_{\lambda \in \text{Alt}(\hat{S}_t)} \left[\sum_{k=1}^K N_{t,k} d(\hat{\mu}_{t,k}, \lambda_k) \right] > \beta(t, \delta) \right\}. \end{aligned}$$

The latter form corresponds to the GLR stopping rule or *Chernoff stopping*, which can also be derived from a hypothesis-testing perspective (Garivier & Kaufmann 2021), and has been at the core of the most recent algorithms for pure exploration (Garivier & Kaufmann 2016; Degenne, Ménard, et al. 2020).

Remark. Confidence-sequence and generalized likelihood-ratio (GLR) constructions are two sides of the same coin: modern results (Ramdas et al. 2022) show equivalences between anytime valid confidence sets and likelihood-based tests (such as Chernoff stopping), which explains the parallel forms of the rules above. Deriving tight confidence sequences has played a central role in the design of bandit algorithms (T. L. Lai 1976; Howard et al. 2018; Kaufmann & W.-M. Koolen 2021; Emmenegger et al. 2023; Kirschner et al. 2025).

Bayesian stopping. Let P_0 be a prior over μ and P_t the posterior at round t . A basic principle is to stop when the posterior probability of error is small:

$$e_t(q) := \mathbb{P}_{\lambda \sim P_t}(\mathcal{S}^*(\lambda) \neq q \mid \mathcal{H}_{t-1}), \quad \hat{S}_t \in \arg \min_q e_t(q), \quad \tau_\delta = \inf\{t \geq 1 : e_t(\hat{S}_t) \leq f(t, \delta)\}.$$

This stopping time triggers when the posterior mass on incorrect answers falls below δ . Proposed by Russo 2016 and analyzed by Shang et al. 2020, these rules yield near-optimal sample complexity but are often computationally demanding in complex, structured, or multi-dimensional settings.

Summary. Across all approaches, the core idea remains identical: stop as soon as the collected evidence is sufficient to rule out all alternative hypotheses with high confidence. Stopping rules derived from Hoeffding-type confidence regions are simple and offer finite-time guarantees but can be conservative, while GLR-based rules achieve asymptotic optimality but are much more involved to compute. A bandit policy must then pull arms in order to reduce uncertainty as fast as possible.

1.2.2 Sampling Rules

Sampling determines where information is gathered. The most frequent sampling rules belong to the following generic families:

- Non-adaptive allocations (uniform, stratified) are simple and analytically convenient, but can be wasteful when the number of arms is large (this will be illustrated throughout the manuscript). They are often coupled with adaptive eliminations, for example, when some arms can be confidently ruled out as sub-optimal.
- Uncertainty-based (racing, LUCB-style) prioritizes arms whose status (optimal vs. suboptimal; feasible vs. infeasible) is most uncertain under current confidence sets. A typical example is the OFU (*Optimism in the Face of Uncertainty*) principle, which, as employed in Jamieson, Malloy, et al. 2014, pulls the arm with the largest optimistic estimate of its mean.
- Tracking aims at targeting instance-optimal proportions (solutions to max-min formulations of instance-dependent lower bounds) via tracking or posterior sampling (e.g., Top-Two variants (Russo 2016), Track-and-Stop (Garivier & Kaufmann 2016)). These naturally pair with GLR or Bayesian stopping to approach information-theoretic lower bounds.

1.2.3 Lower Bounds

At a deeper level, the limits of adaptivity are governed by *information-theoretic lower bounds*, which formalize the minimal number of samples required for reliable identification. These results show that, regardless of the sampling strategy, no algorithm can beat a threshold determined by the intrinsic difficulty of the problem instance. A key ingredient in such analyses is the *information contraction lemma*, popularized by Kaufmann, Cappé, et al. 2016 and later refined in Garivier, Ménard, et al. 2019; Garivier & Kaufmann 2021. Informally, it states that if two bandit models ν and ν' generate similar distributions over observation histories (in the sense of KL divergence), then no algorithm can reliably distinguish them without collecting a sufficient number of samples. This principle underlies all instance-dependent lower bounds developed in this dissertation.

Lemma 1.2.1 (Information contraction, Kaufmann, Cappé, et al. 2016; Garivier, Ménard, et al. 2019). *Let ν and ν' be two bandit models, and let \mathbb{P}_ν and $\mathbb{P}_{\nu'}$ denote the distributions over histories induced by an arbitrary (possibly adaptive) strategy up to stopping time τ . Then for any event \mathcal{E} measurable with respect to \mathcal{H}_τ ,*

$$\sum_{a=1}^K \mathbb{E}_\nu[N_{\tau,a}] \text{KL}(\nu_a \| \nu'_a) \geq \text{kl}(\mathbb{P}_\nu(\mathcal{E}), \mathbb{P}_{\nu'}(\mathcal{E})),$$

where $\text{kl}(\cdot, \cdot)$ denotes the binary relative entropy and $N_{\tau,a}$ is the number of samples drawn from arm a before stopping.

This result shows that the ability of an algorithm to discriminate between events under two instances is fully determined by its sampling frequencies and the KL divergence between arm distributions. It forms the cornerstone of all lower bound arguments used later in this work.

In the next sections, we instantiate these generic components for the multi-objective setting and, in particular, Pareto set identification.

1.3 Background

Originally inspired by adaptive clinical trial design, where treatments are sequentially allocated to balance learning and patient welfare, the bandit framework has since evolved into a general paradigm for sequential decision-making under uncertainty.

This chapter reviews prior work, with a focus on the transition from single- to multi-objective formulations. We first revisit core results in *Best Arm Identification* (BAI) and related single-objective problems, which form the methodological foundation for our study. We then turn to the multi-objective setting, reviewing existing approaches to *Pareto Set Identification* (PSI) and its related problems.

1.3.1 Best Arm Identification and Beyond

The rich literature on BAI has laid the theoretical foundation for pure exploration in bandits, with sophisticated tools such as change-of-measure arguments (Kaufmann, Cappé, et al. 2016), law of iterated logarithm (LIL) confidence bounds (Jamieson, Malloy, et al. 2014), and generalized likelihood ratio stopping rules (Garivier & Kaufmann 2021). These methods will play a central role as we extend pure exploration to multi-objective settings, the focus of this dissertation.

Understanding the complexity of pure exploration, particularly the BAI problem, has long been a central question in the stochastic bandit literature. Given K arms with means $\mu = (\mu_1, \dots, \mu_K)$, the learner sequentially samples arms until a stopping time τ , and must recommend $\widehat{S}_\tau \in [K]$ such that $\mathbb{P}(\widehat{S}_\tau \neq \mathcal{S}^*) \leq \delta$, where $\mathcal{S}^* = \operatorname{argmax}_a \mu_a$. The design of efficient algorithms depends on both the setting (fixed-confidence or fixed-budget) and the statistical assumptions on arm distributions.

From gap-based heuristics to information-theoretic analysis. Early fixed-confidence algorithms, such as *Successive Elimination*, *LUCB*, and *LIL-UCB* (Even-Dar et al. 2002; Gabillon et al. 2012; Kalyanakrishnan et al. 2012; Karnin et al. 2013; Jamieson, Malloy, et al. 2014), relied on subgaussian concentration inequalities to derive confidence intervals around empirical means. Their performance was characterized by the sub-optimality gaps, $\Delta_i = \mu_{\mathcal{S}^*} - \mu_i$ for $i \neq \mathcal{S}^*$ and $\Delta_{\mathcal{S}^*} = \min_{i \neq \mathcal{S}^*} \Delta_i$, leading to sample complexity guarantees of the form

$$\mathbb{E}_\nu[\tau] \leq \mathcal{O}(H(\nu) \log(H(\nu)/\delta)), \quad H(\nu) := \sum_{i=1}^K \frac{1}{\Delta_i^2}.$$

While these results captured key scaling laws, the guarantees were not instance-optimal, an important mismatch was observed between best known lower and upper bounds (Jamieson, Malloy, et al. 2014; Kaufmann, Cappé, et al. 2016).

A major refinement came from replacing subgaussian confidence bounds with information-theoretic quantities such as Kullback-Leibler divergences. Kaufmann & Kalyanakrishnan 2013 introduced *KL-LUCB*, which achieves sharper confidence intervals and tighter sample complexity bounds governed by the Chernoff information rather than mean gaps. The theoretical shift from difference-based to KL-based analysis eventually culminated in the *game-value characterization* of bandit complexity, introduced by Garivier & Kaufmann 2016 and defined as:

$$T^*(\nu)^{-1} = \max_{w \in \Delta_K} \inf_{\lambda \in \operatorname{Alt}(\mathcal{S}^*(\mu))} \sum_{i=1}^K w_i \operatorname{KL}(\nu_a \| \lambda_a), \quad (1.3)$$

where $\Delta_K := \{w \in \mathbb{R}_+^K : \sum_{i=1}^K w_i = 1\}$ is the probability simplex, ν is a bandit with means μ , and $\operatorname{Alt}(\mathcal{S}^*(\mu))$ denotes the set of alternative bandit instances with a different best arm. The quantity $T^*(\nu)$ is called the *characteristic time* of instance ν and, essentially $T^*(\nu) \log \frac{1}{\delta}$ represents the optimal sample complexity for which an algorithm can distinguish ν from its alternatives, in particular when δ is small.

We say that an algorithm is *asymptotically optimal* if its sample complexity matches that lower bound, namely if $\limsup_{\delta \rightarrow 0} \mathbb{E}_\nu[\tau_\delta] / \log \frac{1}{\delta} \leq T^*(\nu)$.

Asymptotically optimal algorithms. Garivier & Kaufmann 2016 established the first asymptotically optimal BAI algorithm, *Track-and-Stop*, achieving $\mathbb{E}_\nu[\tau_\delta] \sim_{\delta \rightarrow 0} T^*(\nu) \log \frac{1}{\delta}$ for any bandit with distributions in one-parameter exponential families. They were the first to exactly match the lower bound, in the asymptotic regime ($\delta \rightarrow 0$).

The idea of *Track-and-Stop* algorithms is first, to have a forced exploration procedure to ensure that at least a sublinear fraction of samples is allocated to every arm, e.g., ensure that at round t each arm is sampled at least \sqrt{t} times. This ensures, by the law of large numbers, that empirical estimates converge to the true means almost surely. Then, at each round, *Track-and-Stop* algorithms compute the $\hat{w}^*(\hat{\mu}_t) \in \Delta_K$, a maximizer of Eq. (1.3) for the empirical means, and following a tracking procedure, they will select the next arm to sample. This sampling rule is coupled with a generalized likelihood ratio (GLR) stopping rule ensuring δ -PAC correctness. Similar algorithmic designs further appeared in more general identification problems, linear bandits (Ménard 2019; Degenne, Ménard, et al. 2020; Jedra & Proutiere 2020a).

In parallel, Bayesian approaches such as *Top-Two Thompson Sampling* (Russo 2016) and related variants showed that posterior sampling can asymptotically realize the same optimal allocation without explicitly solving the optimization problem in Eq. (1.3), while admitting elegant analyses in terms of posterior contraction rates (Shang et al. 2020; Jourdan, Degenne, Baudry, et al. 2022; You et al. 2023).

Fixed-budget formulations. The fixed-budget regime, formalized by Audibert & Bubeck 2010, considers a fixed sampling horizon T and seeks to minimize the error probability $e_T(\nu) := \mathbb{P}(\hat{S}_T \neq S^*)$. Two influential algorithms, *Successive Rejects* (SR) and *Sequential Halving* (SH), were proposed by Audibert & Bubeck 2010 and Karnin et al. 2013, respectively.

These algorithms allocate the budget in phases, discarding empirically suboptimal arms. Their analysis shows that the error probability decreases at an exponential rate

$$e_T(\nu) \leq \exp\left(-c \frac{T}{H(\nu) \log K}\right),$$

and lower bounds (Carpentier & Locatelli 2016; Degenne 2023) show that these rates are essentially unimprovable up to logarithmic factors.

Beyond BAI. Beyond standard Best Arm Identification, several extensions have broadened the scope of pure exploration in single-objective settings. These include *thresholding bandits* (Locatelli et al. 2016), where the goal is to classify arms as above or below a fixed level.

Other extensions include structured formulations such as (*transductive*) *linear* and *combinatorial* bandits (Chen et al. 2014; Soare et al. 2014; Fiez et al. 2019; Degenne, Ménard, et al. 2020). Other formulations consider *multi-fidelity* (Poiani et al. 2024) or *cost-sensitive* sampling (Kanarios et al. 2024), where arms differ in evaluation cost or reliability, as well as *safe* or *constrained* bandits (Camilleri, A. Wagenmaker, et al. 2022), where exploration must respect safety or feasibility conditions throughout the learning process; good arm identification (Roy Chaudhuri & Kalyan Krishnan 2019), a relaxation of *thresholding bandits* where the goal is to identify an arm whose mean is above a given threshold.

Summary. The evolution from gap-based to information-theoretic designs has led to a unified theoretical understanding of pure exploration in single-objective bandits. Fixed-budget algorithms optimize allocation under a deterministic horizon, while fixed-confidence algorithms balance adaptive sampling and stopping. The resulting principles: optimal allocation, information-theoretic stopping, and posterior contraction, form the methodological backbone for the multi-objective extensions developed in this dissertation.

1.3.2 Multi-objective Pure Exploration

The extension of pure exploration to multi-objective settings, where arms are evaluated along multiple potentially conflicting criteria, introduces fundamental challenges absent from single-objective problems. Rather than identifying a single best arm, the learner must characterize a set of arms representing different trade-offs—such as the Pareto set of non-dominated alternatives, or arms satisfying feasibility and ranking constraints. This section surveys the landscape of multi-objective pure exploration, reviewing theoretical results and algorithmic approaches for Pareto set identification and related multi-objective formulations.

Pareto Set Identification (PSI). Given a multi-armed bandit with K arms, where each arm $i \in [K]$ has a mean vector $\mu_i \in \mathbb{R}^d$, we say that arm i is *Pareto dominated* by arm j (denoted $\mu_i \preceq \mu_j$) if $\mu_i^c \leq \mu_j^c$ for all objectives $c \in [d]$ with strict inequality for at least one objective. The strict domination (strict inequality for all objectives) is denoted $\mu_i \prec \mu_j$. The *Pareto set* is the set of non-dominated arms:

$$\mathcal{S}^*(\mu) = \{i \in [K] : \nexists j \in [K] \setminus \{i\}, \mu_i \preceq \mu_j\}.$$

The goal of Pareto set identification is to identify $\mathcal{S}^*(\mu)$ with high confidence. When $d = 1$, this reduces to the classical best-arm identification problem.

Early studies of PSI focused on uniform sampling and elimination schemes (Zuluaga, Sergeant, et al. 2013; Auer et al. 2016; Zuluaga, Krause, et al. 2016), formulated under the fixed-confidence setting. The seminal contribution of Auer et al. 2016 provided the first formal instance-dependent analysis of PSI algorithms, introducing key instance-dependent quantities to characterize problem difficulty and establishing worst-case matching upper and lower bounds on their sample complexity.

To quantify the hardness of PSI, Auer et al. 2016 defined arm-specific *suboptimality gaps* that measures how difficult it is to determine whether an arm is Pareto-optimal. For arms $i, j \in [K]$, define¹

$$\begin{aligned} m(i, j) &:= \min_{c \in [d]} [\mu_j^c - \mu_i^c], \\ M(i, j) &:= \max_{c \in [d]} [\mu_i^c - \mu_j^c]. \end{aligned} \tag{1.4}$$

¹The initial definition was $\tilde{m}(i, j) = \max\{m(i, j), 0\}$ and $\tilde{M}(i, j) = \max\{M(i, j), 0\}$. Some algebra shows that these expressions can be simplified to Eq. (1.4) while keeping the gaps' value.

If $\mu_i \not\prec \mu_j$, then $M(i, j)$ is the smallest uniform increase of j required to dominate i , while, if $\mu_i \prec \mu_j$, then $m(i, j)$ is the smallest coordinate-wise increase of i needed to make it non-strictly-dominated by j . From these primitives, one defines for each arm i the suboptimality gap

$$\Delta_i := \begin{cases} \max_{j \in \mathcal{S}^*} m(i, j), & \text{if } i \notin \mathcal{S}^*, \\ \min\{\delta_i^+, \delta_i^-\}, & \text{if } i \in \mathcal{S}^*, \end{cases} \quad (1.5)$$

where

$$\begin{aligned} \delta_i^+ &:= \min_{j \in \mathcal{S}^* \setminus \{i\}} \min\{M(i, j), M(j, i)\}, \\ \delta_i^- &:= \min_{j \notin \mathcal{S}^*} [(M(j, i))_+ + \Delta_j]. \end{aligned}$$

These quantities capture how close each arm lies to the Pareto frontier. The overall instance complexity is then measured by

$$H(\nu) := H(\mu) := \sum_{i=1}^K \frac{1}{\Delta_i^2}. \quad (1.6)$$

This term was shown by [Auer et al. 2016](#) to govern the finite-time instance-dependent sample complexity of fixed-confidence PSI algorithms through an upper bound on the sample complexity of their algorithm of order

$$\mathcal{O}(H(\mu) \log(dH(\mu)/\delta)),$$

combined with a lower bound, also scaling with $H(\mu)$ (recalled in [Theorem 1.3.1](#)). This result establishes a nearly-matching information-theoretic lower bound, confirming (up to logarithmic factors) the dependence of the sample complexity on $H(\mu)$ for some instances.

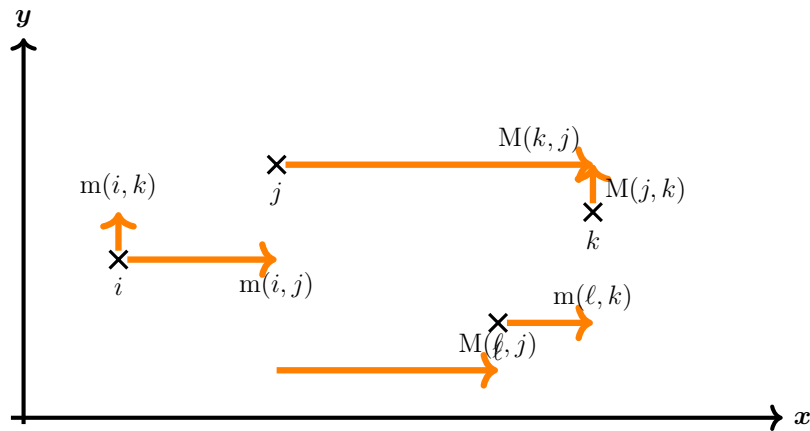


Figure 1.4: PSI gaps and distances

Theorem 1.3.1 (Theorem 17 of [Auer et al. 2016](#)). *For any collection of mean vectors $\mu_i \in [1/4, 3/4]^d$, there exist distributions $(\mathcal{D}_i)_{1 \leq i \leq K}$ such that, with probability at least $1 - \delta$, any δ -correct algorithm for PSI requires at least*

$$\Omega\left(\sum_{i=1}^K \frac{1}{\widetilde{\Delta}_i^2} \log \frac{1}{\delta}\right)$$

samples to identify the Pareto set, where for suboptimal arms $\widetilde{\Delta}_i := \Delta_i$ and for optimal arms $\widetilde{\Delta}_i := \delta_i^+$.

From an application perspective, PSI provides a rigorous abstraction of multi-criteria decision-making in which several performance indicators must be considered simultaneously. In clinical research, this is particularly relevant to early-phase vaccine or treatment studies, as motivated in Section 1.1.

In what follows, we review several extensions of PSI proposed to handle more general forms of partial orders, feasibility constraints, and structured models.

Cone-Ordered Pareto Set Identification. A generalization of the Pareto set identification problem was proposed by [Ararat & Tekin 2023](#), who extended the notion of dominance to arbitrary partial orders induced by convex cones. In this framework, a cone $\mathcal{C} \subset \mathbb{R}^d$ defines a preference relation between mean vectors, where $\mu_i \prec_{\mathcal{C}} \mu_j$ if and only if $\mu_j - \mu_i \in \mathcal{C}$. The standard Pareto order corresponds to the positive orthant \mathbb{R}_+^d , but other cones can encode domain-specific trade-offs between objectives. Another example of a cone is $\mathcal{C}_\theta = \{x : \theta^\top x \geq 0\}$, for a direction $\theta \in \mathbb{R}^d$, which results in a scalarization of the multi-objective problem ([Miettinen 1998](#)).

This cone-based formulation unifies PSI and several related problems, such as constrained optimization and scalarization-based selection, within a single mathematical framework. [Ararat & Tekin 2023](#) analyzed this problem in the fixed-confidence setting and proposed a general elimination algorithm whose sample complexity depends on cone-specific quantities generalizing the suboptimality gaps of [Auer et al. 2016](#). Their work showed that, although the underlying geometry of the cone affects the difficulty of the problem, similar logarithmic sample complexity dependencies can be achieved, extending PSI theory beyond the axis-aligned Pareto order. This line of work was later studied by [Shukla & Basu 2024](#); [Das et al. 2025](#) where, instead of identifying optimal arms, the goal is rather to identify a policy supported on optimal arms. The authors derived an information-theoretic lower bound and a Track-and-Stop style algorithm.

Feasible Arm Identification. The *Feasible Arm Identification* problem was introduced by [Katz-Samuels & Scott 2018](#) in the fixed-budget setting. In this formulation, the learner is given a convex polyhedron $P \subset \mathbb{R}^d$, $d \geq 1$ and by interacting with a multi-objective bandit model, must identify the set of arms whose mean vectors lie within P , given a prescribed exploration budget. This setting can be viewed as a generalization of the classical *bandit thresholding problem* ([Locatelli et al. 2016](#)), extended to multi-dimensional objectives. The

authors proposed MD-APT (Multi-Dimensional Anytime Parameter-Free Thresholding), an algorithm that achieves a near-optimal fixed-budget guarantee:

$$e_T \leq \mathcal{O}\left(e^{-c \cdot \frac{T}{H(\mu)}}\right),$$

where c is an absolute constant, $H(\mu) := \sum_i \Delta_i^{-2}$ is a complexity measure based on the arm-dependent gaps

$$\Delta_i := \begin{cases} \text{dist}(\mu_i, P^c), & \text{if } \mu_i \in P, \\ \text{dist}(\mu_i, \partial P), & \text{otherwise.} \end{cases}$$

The authors further showed that this rate is essentially optimal by establishing instance-dependent lower bounds for specific problem configurations. Beyond its theoretical appeal, the feasible-arm formulation has natural connections to adaptive clinical trial design: the polyhedron P can represent eligibility or safety constraints—such as acceptable toxicity-efficacy trade-offs, and identifying feasible arms corresponds to discovering treatments that satisfy predefined clinical criteria under limited experimental budget.

Top Feasible Arm Identification. The *Top Feasible Arm Identification* problem (Katz-Samuels & Scott 2019) extends the feasible-arm framework by incorporating preference rankings over multiple objectives.

In this setting, the learner is given a convex polyhedron $P \subset \mathbb{R}^d$, a weighting vector $w \in \mathbb{R}^d$, and an integer parameter $m \leq K$. Letting $F(\mu) := \{k : \mu_k \in P\}$ denote the set of feasible arms, the learner’s goal is to identify

$$\mathcal{O}^* := \left\{k : \mu_k \in P \text{ and } w^\top \mu_k \geq \max_{i \in F(\mu)}^m w^\top \mu_i\right\},$$

that is, the m best feasible arms according to the scalarized objective $k \mapsto w^\top \mu_k$. The authors formulated a fixed-confidence version of the problem, called δ -PAC-explanatory, where the learner outputs a triplet $(\widehat{O}, \widehat{S}, \widehat{I})$ forming a partition of $[K]$ such that

$$\widehat{O} = \mathcal{O}^*, \quad \widehat{I} \subset F(\mu)^c, \quad \widehat{S} \subset \{k : w^\top \mu_k \leq \max_{i \in \mathcal{O}^*} w^\top \mu_i\}.$$

This formulation captures both the feasibility constraint and the ranking objective within a unified PAC framework, and provides statistical guarantees on the correctness of the identified sets.

From a clinical trial perspective, this problem naturally models the selection of the most promising treatments among those satisfying pre-defined safety or eligibility criteria. The weight vector w encodes the relative importance of various clinical endpoints—such as immunogenicity, safety, or tolerability—while the polyhedron P specifies admissible safety-efficacy trade-offs. Hence, identifying \mathcal{O}^* corresponds to finding the m best treatment candidates that are both clinically acceptable and statistically supported within a fixed-budget or fixed-confidence experiment.

Constrained Bi-Objective Identification. A closely related special case of multi-objective pure exploration with constraints was introduced by [Faizal & Nair 2022](#). The authors considered a bi-objective ($d = 2$) bandit model where each arm $k \in [K]$ is associated with an unknown mean vector $\mu_k = (\mu_k^{(1)}, \mu_k^{(2)}) \in [0, 1]^2$. The learner aims to identify the arm maximizing the first objective $\mu_k^{(1)}$ under a feasibility constraint on the second, namely

$$\mathcal{S}_{\text{feas}}^* := \left\{ k \in [K] : \mu_k^{(2)} \leq b \text{ and } \mu_k^{(1)} = \max_{i: \mu_i^{(2)} \leq b} \mu_i^{(1)} \right\},$$

where b is a known threshold. This setting can be interpreted as a variant of *Top Feasible Arm Identification*, where the feasibility region is $\{\mu : \mu^{(2)} \leq b\}$.

In the fixed-budget regime, the authors proposed a Successive-Rejects-type algorithm, C-SR, which alternates between (i) feasibility testing to discard arms likely to violate the constraint and (ii) dominance elimination among the remaining feasible arms. They showed that the algorithm achieves an exponential rate

$$e_T(\nu) \leq \mathcal{O}\left(\exp\left(-c \frac{T}{H_{\text{feas}}(\mu) \log K}\right)\right),$$

for some constant $c > 0$, where the instance complexity is governed by the composite measure

$$H_{\text{feas}}(\mu) := \sum_{i=1}^K \frac{1}{\tilde{\Delta}_i^2},$$

where for $i \neq \mathcal{S}^*$, and assuming the feasible set is non-empty, $\Delta_i = \min(b - \mu_{\mathcal{S}^*}^{(2)}, \tilde{\delta}_i)$ with

$$\tilde{\delta}_i := \begin{cases} \mu_{\mathcal{S}^*}^{(1)} - \mu_i^{(1)} & \text{if } \mu_i^{(2)} \leq b \\ \mu_i^{(2)} - b & \text{if } \mu_i^{(2)} > b \text{ and } \mu_i^{(1)} \geq \mu_{\mathcal{S}^*}^{(1)} \\ \max(\mu_i^{(2)} - b, \mu_{\mathcal{S}^*}^{(1)} - \mu_i^{(1)}) & \text{else.} \end{cases}$$

This expression highlights the two competing sources of difficulty: arms close to the feasibility boundary (small $|\mu_k^{(2)} - b|$) and arms close to optimality within the feasible region (small $\mu_{\mathcal{S}^*}^{(1)} - \mu_k^{(1)}$). Although limited to $d = 2$, this work formally connected feasibility constraints and multi-objective identification, and provided one of the first fixed-budget guarantees for a constrained pure exploration problem. Such a setting is directly relevant to early-phase clinical trials, where treatments must simultaneously achieve sufficient efficacy and remain below pre-specified toxicity or reactogenicity thresholds.

Summary and open challenges. The studies reviewed above reveal that, before the onset of this work, the theory of multi-objective pure exploration in bandits remained largely fragmented. Existing algorithms were either restricted to low-dimensional cases (e.g., two-objective feasibility problems), limited to fixed-confidence formulations, or relied on uniform sampling strategies with conservative performance guarantees. Most of the existing algorithms were designed under the assumption of uncorrelated objectives, or were not able to exploit potential correlation that may arise for example when both toxicity and efficacy are recorded.

In contrast to the well-understood single-objective setting, several fundamental questions were still open: How to extend instance-dependent analyses to the fixed-budget regime? How to adaptively allocate samples across arms to solve large-scale, high-dimensional Pareto problems? How to design an asymptotically optimal and computationally-efficient algorithm for PSI? How to handle additional structural or feasibility constraints that naturally arise in clinical and engineering applications?

The next section introduces the methodological and theoretical contributions of this dissertation, which aim to fill these gaps by developing adaptive, efficient, and interpretable algorithms for multi-objective exploration.

1.4 Contributions of This Dissertation

Building on the gaps identified in the previous section, this chapter presents the main contributions of the dissertation. It develops new theoretical and algorithmic tools for *Pareto Set Identification* (PSI) and related *pure exploration* problems in stochastic bandits, with a focus on multi-objective and constrained formulations motivated by early-stage clinical trials. The section is organized around three broad themes: (i) the fixed-budget setting, where the sampling budget is prescribed in advance, (ii) the fixed-confidence setting, where sampling continues until sufficient evidence is collected, and (iii) structured and constrained extensions of PSI. Each of these directions addresses an open question and is detailed in the following paragraphs.

1.4.1 Fixed-Budget PSI

This section summarizes our work on the fixed-budget formulation of Pareto Set Identification (PSI) presented in [Kone, Kaufmann, et al. 2024](#). This constitutes the first theoretical and algorithmic analysis of PSI when the total number of samples is fixed in advance.

In many experimental settings, particularly in early-stage vaccine trials, the total number of enrolled subjects is predetermined. The learner must therefore allocate a fixed budget T across K arms and recommend an estimated Pareto set \hat{S}_T after exactly T samples. The goal is to minimize the probability of misidentification

$$e_T(\nu) = \mathbb{P}_\nu(\hat{S}_T \neq \mathcal{S}^*(\mu)) ,$$

and to characterize how quickly $e_T(\nu)$ decays with T .

While the fixed-confidence regime had been studied, *e.g.*, in [Zuluaga, Sergent, et al. 2013](#); [Auer et al. 2016](#); [Zuluaga, Krause, et al. 2016](#), no results were available in the fixed-budget setting. To bridge this gap, we propose *Empirical Gap Elimination* (EGE), a family of elimination-based algorithms that generalizes classical strategies such as Successive Rejects and Successive Halving to the multi-objective setting.

The design of EGE relies on an empirical reformulation of the suboptimality gaps introduced by [Auer et al. 2016](#). Recall that these gaps depend on the true Pareto set $\mathcal{S}^*(\mu)$ and therefore

were difficult to estimate directly. We show that these gaps (as recalled in Eq. (1.5)) can be equivalently rewritten as

$$\Delta_i = \max\{\Delta_i^*, \delta_i^*\},$$

where we have

$$\begin{aligned} \Delta_i^* &:= \max_{j \in [K]} m(i, j) \\ \delta_i^* &:= \min_{j \in [K] \setminus \{i\}} [M(i, j) \wedge (M(j, i)_+ + (\Delta_j^*)_+)], \end{aligned}$$

which makes them amenable to empirical estimation without prior knowledge of $\mathcal{S}^*(\mu)$. EGE maintains at each phase r an active set of arms \mathcal{A}_r and empirical means $(\hat{\mu}_{i,r})_{i \in \mathcal{A}_r}$, and computes the empirical gaps $\hat{\Delta}_{i,r}$ obtained by replacing μ with $\hat{\mu}$ and $[K]$ with \mathcal{A}_r . Arms with large empirical gaps are progressively eliminated with a particular tie-breaking rule.

The budget T is split across successive elimination rounds. In EGE-SR, the sampling budget is divided into $K - 1$ phases of increasing length, and one arm is removed per phase. In EGE-SH, the budget is divided geometrically and approximately half of the remaining arms are eliminated at each round. Both strategies share the same elimination criterion but differ in their allocation schedule.

Our analysis establishes the first information-theoretic lower bound for PSI in the fixed-budget setting. For any adaptive algorithm, there exists a universal constant $c > 0$ and a class of instances $\tilde{\mathcal{D}}^K$ such that

$$\forall \nu \in \tilde{\mathcal{D}}^K, \liminf_{T \rightarrow \infty} -\frac{1}{T} \log e_T(\nu) \leq \frac{c}{H_2(\mu)},$$

where

$$H_2(\nu) := H_2(\mu) := \max_{k \in [K]} k \Delta_{(k)}^{-2}, \quad (1.7)$$

and $k \mapsto (k)$ is the permutation such that $\Delta_{(1)} \leq \Delta_{(2)} \leq \dots \leq \Delta_{(K)}$. We further show that both EGE-SR and EGE-SH achieve

$$e_T(\nu) \leq \mathcal{O}\left(\exp\left(-c \frac{T}{H_2(\mu) \log K}\right)\right),$$

which nearly matches the lower bound up to logarithmic factors, thereby characterizing the fundamental error exponent for PSI under a fixed sampling budget.

Empirical evaluations confirm the theoretical predictions: both EGE-SR and EGE-SH substantially outperform uniform allocation, particularly on unbalanced instances where dominated arms are easily identifiable.

In practice, they offer principled and adaptive allocation strategies for early-phase clinical studies, where limited sample budgets must be distributed across multiple candidate treatments and efficacy indicators.

1.4.2 Adaptive Fixed-Confidence Algorithms for PSI

In this contribution, we introduce the first fully adaptive, LUCB-style algorithm for fixed-confidence Pareto Set Identification (PSI) with finite-time guarantees. While prior works such as [Auer et al. 2016](#) established PAC guarantees under uniform sampling, our goal is to design fully adaptive algorithms that pull one or two informative arms per round.

Formally, given K arms with d -dimensional mean vectors μ_1, \dots, μ_K , the learner must identify the Pareto set $\mathcal{S}^* = \{i \in [K] : \nexists j, \mu_i \prec \mu_j\}$ with probability at least $1 - \delta$, while minimizing the expected stopping time $\mathbb{E}[\tau]$. The central difficulty lies in the fact that dominance relations depend on all objectives simultaneously, making the problem inherently multi-dimensional and asymmetric. To capture this structure, we introduce an adaptive gap-based framework that extends the LUCB (Lower-Upper Confidence Bound) principle from Best Arm Identification to multi-objective bandits.

At each round t , the algorithm maintains $(1 - \delta)$ -level confidence boxes for each arm,

$$\mathcal{R}_{t,i}(\delta) = \prod_{c=1}^d [\hat{\mu}_{t,i}^c - \beta_i(t, \delta), \hat{\mu}_{t,i}^c + \beta_i(t, \delta)],$$

where $\hat{\mu}_{t,i}^c$ denotes the empirical mean on objective c , and $\beta_i(t, \delta)$ is a confidence radius calibrated via a law-of-iterated-logarithm bound. Arms whose confidence boxes are entirely dominated by those of others are eliminated, while sampling concentrates on ambiguous pairs (i, j) where the dominance relation remains uncertain. This adaptive allocation rule implicitly tracks the "most confusing" arms, in analogy to classical LUCB algorithms for single-objective identification.

We derive finite-time, instance-dependent guarantees showing that the proposed algorithms achieve near-optimal sample complexity:

$$\mathbb{E}_\nu[\tau_\delta] \leq \mathcal{O}(H(\mu) \log(dH(\mu)/\delta)),$$

where $H(\nu) := \sum_{i=1}^K \Delta_i^{-2}$ is the complexity term defined from the Pareto gaps in (1.6). This matches the information-theoretic lower bound up to logarithmic factors, establishing the first finite-time optimality guarantees for PSI. Our analysis further characterizes the adaptive evolution of the active set and provides confidence-based stopping rules that ensure δ -PAC correctness without knowledge of instance parameters.

In addition, we introduce a relaxed formulation called the k -Pareto set identification problem, in which the learner returns a set $\hat{\mathcal{S}}$ of k elements satisfying $\hat{\mathcal{S}} \subseteq \mathcal{S}^*$ (or the full Pareto set if $k \geq |\mathcal{S}^*|$).

This relaxation handles practical situations where only a subset of the Pareto set needs to be identified. It is particularly meaningful in clinical studies, where one may only need to shortlist a few optimal strategies for further testing. We show that its complexity $H^k(\nu) \leq H(\nu)$ yields provably faster identification.

Empirically, our algorithms substantially outperform uniform and non-adaptive baselines, particularly when the number of arms is large. These results establish the first fully adaptive

algorithm for fixed-confidence PSI, bridging the methodological gap between classical LUCB-style exploration and PSI.

1.4.3 PSI with Linear Structure

In this contribution, we introduce and analyze *Pareto Set Identification (PSI) with linear structure*. This setting generalizes classical (unstructured) PSI by assuming that all arms share a low-dimensional parameterization: each arm $a \in [K]$ is associated with a known feature vector $x_a \in \mathbb{R}^h$, and its expected outcome vector is linear in an unknown parameter $\theta \in \mathbb{R}^{h \times d}$, that is

$$\mu_a = \theta^\top x_a, \quad \theta \in \mathbb{R}^{h \times d} \text{ and } x_a \in \mathbb{R}^h,$$

where each arm a has σ -subgaussian marginals. This structural assumption induces correlations between arms, allowing information gathered from one arm to benefit all others.

The goal remains to identify the Pareto-optimal set $\mathcal{S}^*(\mu)$ with probability at least $1 - \delta$, while minimizing the number of samples. In contrast to previous unstructured algorithms that sample arms uniformly or according to empirical dominance relations, we design a fully *structure-aware* sampling strategy guided by principles of optimal experimental design (Soare et al. 2014).

At the heart of our approach lies a connection between PSI estimation and G -optimal design. Let $V(w)$ denote the matrix associated with an allocation $w \in \Delta_K$, defined as

$$V(w) = \sum_{i=1}^K w_i x_i x_i^\top.$$

The key observation is that if arms are pulled according to a static allocation w , and an empirical least-squares estimate of θ is built, the uncertainty (covariance) on each arm's estimated mean $\hat{\mu}_a = \hat{\theta}^\top x_a$ is governed by $\|x_a\|_{V(w)^{-1}}^2 \cdot \Sigma$, where Σ is the (assumed common) covariance matrix of each arm's marginals. An allocation is thus informative if it minimizes the worst-case prediction variance across the set of relevant arms, that is,

$$\rho^* = \operatorname{argmin}_{w \in \Delta_K} \max_{a \in [K]} \|x_a\|_{V(w)^{-1}}^2.$$

This optimization corresponds exactly to the classical G -optimal design criterion (Fiez et al. 2019).

Adaptive algorithm and analysis. We design an elimination-style algorithm that estimates θ using a G -optimal design, and then estimate the empirical PSI gaps similarly to the technique used in Section 1.4.1. We show that the resulting algorithm enjoys instance-dependent guarantees: given $\delta \in (0, 1)$, its sample complexity is upper bounded by

$$\mathcal{O} \left(\inf_{w \in \Delta_K} \max_{k \in [K]} \frac{\|x_k\|_{V(w)^{-1}}^2}{\Delta_k^2} \log \frac{Kd}{\delta \Delta_{\min}} \cdot \log \frac{1}{\Delta_{\min}} \right)$$

with probability at least $1 - \delta$. Further refinement allows one to upper bound the sample complexity by

$$\log_2(2/\Delta_{(1)}) + \mathcal{O}\left(\sum_{i=2}^h \frac{\sigma^2}{\Delta_{(i)}^2} \log\left(\frac{Kd}{\delta} \log_2\left(\frac{2}{\Delta_{(i)}}\right)\right)\right),$$

where (\cdot) is a re-indexing such that $\Delta_{\min} := \Delta_{(1)} \leq \Delta_{(2)} \leq \dots \leq \Delta_{(K)}$. Thus, the sample complexity only depends on the h -smallest gaps, which is a major improvement when $h \ll K$.

In terms of fixed-budget guarantees, we prove that the probability of misidentification of the Pareto set for the proposed algorithm is upper-bounded by

$$\exp\left(-\frac{T}{c\sigma^2 H_{2,\text{lin}}(\nu) \lceil \log_2 h \rceil} + \log C(h, d, K)\right),$$

where $C(h, d, K) = 2d(K + h + \lceil \log_2 h \rceil)$, c is an absolute constant and

$$H_{2,\text{lin}}(\nu) := \max_{k \in [h]} k \Delta_{(k)}^{-2},$$

which is a major improvement over Eq. (1.7), again when $h \ll K$.

From an application perspective, this structured approach provides a natural way to share information across correlated treatments or dosages, *e.g.*, where key components of those treatments are correlated (the vaccine platform, dosages, schedule, or formulation, etc.).

1.4.4 Asymptotically Optimal Algorithms for Fixed-Confidence PSI

The first algorithms introduced for fixed-confidence PSI provided only near-optimal sample complexity guarantees, with non-negligible gaps between upper and lower bounds. In contrast, the Best Arm Identification (BAI) literature has established sharp asymptotic characterizations: information-theoretic lower bounds (Garivier & Kaufmann 2016) and matching algorithms such as TRACK-AND-STOP (Garivier & Kaufmann 2016) or Bayesian approaches (Russo 2016; Jourdan, Degenne & Kaufmann 2023) achieve the optimal expected sample complexity

$$\mathbb{E}_\nu[\tau_\delta] = T^*(\nu) \log(1/\delta) (1 + o(1)).$$

Reaching a similar level of optimality in PSI is significantly more challenging, since the learner must classify all arms simultaneously and the space of alternative hypotheses is combinatorial. Hence, unlike the single-objective case, no algorithm was known to achieve *asymptotic instance-optimality* for PSI.

Information-theoretic lower bound. Let $\mathcal{I} \subset \mathbb{R}^d$ denote the set of possible mean vectors, and $\mathcal{S}^*(\lambda)$ the Pareto set associated with the collection $\{\lambda_1, \dots, \lambda_K\}$. To guarantee

δ -correctness, an algorithm must statistically distinguish the true parameter μ from every alternative $\lambda \in \mathcal{I}^K$ whose Pareto set differs from $\mathcal{S}^*(\mu)$. Define the set of such alternatives as

$$\text{Alt}(\mathcal{S}^*(\mu)) = \{\lambda \in \mathcal{I}^K : \mathcal{S}^*(\lambda) \neq \mathcal{S}^*(\mu)\}.$$

Using information-divergence arguments analogous to those of [Garivier & Kaufmann 2016](#), we obtain the following lower bound on the expected stopping time of any δ -PAC PSI algorithm.

Theorem 1.4.1 (Fixed-confidence lower bound). *Let $\mathcal{D} := \{\mathcal{N}(\mu, \Sigma) : \mu \in \mathcal{I}\}$ denote the set of d -dimensional multivariate Gaussian arms with known covariance Σ and mean vectors in \mathcal{I} . Any algorithm that is δ -correct on all bandit instances in \mathcal{D}^K satisfies, for all ν parameterized by $\mu \in \mathcal{I}^K$,*

$$\mathbb{E}_\nu[\tau_\delta] \geq T^*(\nu) \log\left(\frac{1}{2.4\delta}\right),$$

where the characteristic time $T^*(\nu)$ is defined by

$$T^*(\nu)^{-1} := \sup_{w \in \Delta_K} \inf_{\lambda \in \text{Alt}(\mathcal{S}^*(\mu))} \left[\sum_{i=1}^K \frac{w_i}{2} \|\mu_i - \lambda_i\|_{\Sigma^{-1}}^2 \right].$$

An algorithm is said to be *asymptotically optimal* if

$$\lim_{\delta \rightarrow 0} \frac{\mathbb{E}_\nu[\tau_\delta]}{\log(1/\delta)} = T^*(\nu).$$

This bound characterizes the minimal information required to distinguish the true Pareto set from all incorrect ones. In BAI, algorithms achieving such tight guarantees typically rely on *generalized likelihood-ratio* (GLR) stopping rules, which admit closed-form expressions for exponential-family models. For PSI, this stopping rule generalizes to

$$\tau = \inf \left\{ t \geq 1 : \inf_{\lambda \in \text{Alt}(\widehat{S}_t)} \sum_{k=1}^K N_{t,k} \|\hat{\mu}_{t,k} - \lambda_k\|_{\Sigma^{-1}}^2 > \beta(t, \delta) \right\},$$

but the inner minimization over $\lambda \in \text{Alt}(\widehat{S}_t)$ has no known closed-form solution, making direct implementation intractable for PSI, due to its combinatorial and non-convex structure.

Posterior-resampling stopping rule. To circumvent this limitation, we propose a posterior-based alternative inspired by Bayesian sequential testing. At time t , let P_t denote the posterior distribution over θ given the observed history \mathcal{H}_{t-1} , and let $\theta_t^1, \theta_t^2, \dots$ be *i.i.d.* samples from P_t . Fix a recommendation $\widehat{S}_t = \mathcal{S}^*(\hat{\mu}_t)$ and a threshold $M(t, \delta)$. We define the *posterior geometric stopping rule* as

$$\tau = \inf \left\{ t \geq 1 : \forall m \leq M(t, \delta), \mathcal{S}^*(\theta_t^m) = \widehat{S}_t \right\}, \quad (1.8)$$

that is, we stop once $M(t, \delta)$ consecutive posterior samples yield the same optimal set as the current recommendation.

By proving novel anti-concentration results for Gaussian vectors, we calibrate the threshold $M(t, \delta)$ to ensure δ -correctness.

This rule is conceptually close to Chernoff’s GLR stopping: both trigger when the posterior (or likelihood) ratio between the true and alternative hypotheses becomes sufficiently large. However, unlike GLR-based tests, it does not require explicit evaluation of the posterior probabilities—only the ability to draw posterior samples, which is computationally straightforward for standard conjugate models. This makes it particularly attractive for high-dimensional PSI problems where $\text{Alt}(\mathcal{S}^*(\mu))$ has a complex combinatorial structure.

Main result. We combine this stopping rule with a randomized adaptive allocation strategy based on posterior samples. We prove that the resulting algorithm is asymptotically instance-optimal:

$$\mathbb{E}_\nu[\tau_\delta] = T^*(\nu) \log(1/\delta) (1 + o(1)),$$

and achieves optimal posterior contraction rates. This constitutes the first successful application of posterior sampling to multi-objective pure exploration. Beyond theory, the algorithm is computationally practical and statistically efficient, providing a randomized and interpretable design principle for early-phase clinical trials where decisions must be both adaptive and provably optimal.

1.4.5 PSI with Linear Feasibility Constraints

In many real-world applications, not all arms are admissible: a vaccine dose may need to satisfy toxicity thresholds, a hardware configuration must respect power or latency limits, or a treatment must achieve a minimum efficacy level before being considered. Such constraints are often linear in nature, motivating the definition of the feasible set

$$P = \{x \in \mathbb{R}^d : Ax \leq b\},$$

with $A \in \mathbb{R}^{m \times d}$ and $b \in \mathbb{R}^m$. Given this set, the learner must identify the subset of arms that are both feasible and Pareto-optimal among feasible ones:

$$\mathcal{S}_{\text{feas}}^*(\mu) = \left\{ i \in [K] : \mu_i \in P \text{ and } \nexists j \in [K], \mu_j \in P, \mu_i \prec \mu_j \right\}.$$

Compared to unconstrained PSI, this setting introduces a dual challenge: each arm must be simultaneously classified by both *feasibility* and *dominance*. The sample complexity is thus driven by two quantities: the Pareto gap Δ_i , characterizing the uncertainty about optimality, and the feasibility margin

$$\eta_i := \begin{cases} \text{dist}(\mu_k, P^c), & \text{if } \mu_k \in P, \\ \text{dist}(\mu_k, \partial P), & \text{otherwise.} \end{cases},$$

which measures how close an arm’s mean lies to the constraint boundary.

We propose the *explainable constrained Adaptive Pareto Exploration* (e-CAPE) algorithm, which generalizes the APE algorithm to constrained environments.

The algorithm recommends a partition $R_\tau = (O_\tau, S_\tau, I_\tau)$ such that

$$i) O_\tau = \mathcal{S}_{\text{feas}}(\mu), \quad ii) S_\tau \subseteq \text{SubOpt}(\mu) \text{ and } I_\tau \subseteq \{k : \mu_k \notin P\} \quad (1.9)$$

holds at stopping time τ , where (i) $\mathcal{S}_{\text{feas}}(\mu)$ is the Pareto set of $\{\mu_k : \mu_k \in P, k \in [K]\}$, $\text{SubOpt}(\mu) := \{i \in [K] \mid \exists j \text{ such that } \mu_i \prec \mu_j \text{ and } \mu_j \in P\}$ be the set of arms that are dominated by a feasible arm.

Since some arms can be both infeasible and dominated by a feasible arm, such arms could be correctly classified either in S_τ or I_τ . Therefore, multiple correct ways to partition $[K]$ may exist, so multiple correct partitions (O_τ, S_τ, I_τ) . Given the polyhedron P , we define the set of valid answers for (S_τ, I_τ) as $\mathcal{M}(P, \mu)$.

The learner maintains empirical estimates of both feasibility and Pareto gaps, and dynamically allocates samples according to their estimated classification difficulty. At each iteration, arms near the Pareto-feasibility boundary are targeted. This mechanism achieves a near-optimal balance between exploration of uncertain arms and exploitation of structural information encoded by P .

We proved that e-cAPE correctly identifies a partition in the fixed-confidence setting, with sample complexity characterized by a constrained instance complexity term

$$\mathbb{E}_\nu[\tau_\delta] \lesssim \mathcal{O}(C_{\mathcal{M}}^*(\nu) \log(KdC_{\mathcal{M}}^*(\nu)/\delta)), \quad (\text{for small } \delta),$$

where

$$C(\nu, S, I) := \sum_{i \in \mathcal{S}_{\text{feas}}} \frac{1}{\min(\Delta_i^2(\mathcal{S}_{\text{feas}} \cup S), \eta_i^2)} + \sum_{i \in S} \frac{1}{\Delta_i^2(\mathcal{S}_{\text{feas}} \cup S)} + \sum_{i \in I} \frac{1}{\eta_i^2}, \quad (1.10)$$

and, for $\mathcal{X} \subseteq [K]$, $\Delta_i^2(\mathcal{X})$ denotes the PSI gaps computed in the $|\mathcal{X}|$ -armed bandit with means $\{\mu_k : k \in \mathcal{X}\}$. $C_{\mathcal{M}}^*$ is the complexity of the *easiest to verify* partition:

$$C_{\mathcal{M}}^*(\nu) := \min_{(O, S, I) \in \mathcal{M}(P, \mu)} C(\nu, S, I).$$

We further establish an information-theoretic lower bound of the same order, proving that the dependence on $C_{\mathcal{M}}^*$ is unavoidable.

Finally, experiments on synthetic and realistic multi-objective datasets (Mark C et al. 2013) confirm that e-cAPE significantly outperforms uniform and unconstrained sampling, especially when feasible arms lie near the boundary of P . In practice, such constrained exploration directly models the setting of early-stage vaccine development, where identifying acceptable trade-offs between immunogenicity and safety is more realistic than seeking a single optimal dose.

1.5 Outline of the Dissertation

The dissertation is organized into five technical chapters plus a conclusion chapter. Each technical chapter addresses a distinct aspect of Pareto set identification in multi-objective bandits. The progression follows a natural path, from gap-based and finite-sample analyses to asymptotically optimal, constrained, and structured extensions.

- **Chapter 2** introduces the *fixed-budget* formulation of PSI, relevant when the sampling horizon is predetermined. It presents the *Empirical Gap Elimination* (EGE) framework, extending Successive Rejects and Sequential Halving to multi-objective settings, and establishes the first information-theoretic lower bound for PSI under fixed budget.
- **Chapter 3** addresses PSI in the *fixed-confidence* regime. It develops adaptive elimination algorithms with near-optimal instance-dependent sample complexity and discusses relaxations such as k -Pareto set identification, where the learner outputs a manageable representative subset of the Pareto-optimal arms.
- **Chapter 4** extends PSI to *multi-output linear models*, where arm means depend linearly on known features. Exploiting linear structure, we design sample-efficient algorithms for structured action spaces, connecting PSI estimation with optimal experimental design.
- **Chapter 5** studies *Bayesian algorithms for PSI*, based on posterior sampling and geometric stopping rules. These algorithms achieve asymptotic instance-optimality while remaining computationally efficient, offering a Bayesian randomized algorithm for PSI.
- **Chapter 6** introduces *constrained PSI*, where feasible arms must satisfy linear inequalities (e.g., safety or resource constraints). We develop adaptive algorithms that jointly assess feasibility and dominance and derive complexity measures that couple feasibility margins with Pareto gaps.
- **Chapter 7** concludes the dissertation, outlining future research directions, including large-scale PSI, preference-driven objectives, and extensions to reinforcement learning and other sequential decision frameworks.

Each chapter is largely self-contained, but contributes to a unified view of multi-objective pure exploration, bridging theoretical foundations, algorithmic design, and experiments toward applications to adaptive clinical settings such as vaccine trials.

1.6 Publications

The chapters of this dissertation are primarily based on peer-reviewed papers I co-authored during this doctoral work. Each contribution has been published in international machine learning conferences, reflecting different facets of the Pareto set identification problem.

Conference Papers

- **C. Kone**, E. Kaufmann, L. Richert.
Adaptive Algorithms for Relaxed Pareto Set Identification.
Proceedings of the Thirty-Seventh Conference on Neural Information Processing Systems (NeurIPS 2023).

- **C. Kone**, E. Kaufmann, L. Richert.
Bandit Pareto Set Identification: The Fixed-Budget Setting.
Proceedings of the 27th International Conference on Artificial Intelligence and Statistics (AISTATS 2024).
- **C. Kone**, E. Kaufmann, L. Richert.
Bandit Pareto Set Identification in a Multi-Output Linear Model.
Proceedings of the 28th International Conference on Artificial Intelligence and Statistics (AISTATS 2025).
- **C. Kone**, M. Jourdan, E. Kaufmann.
Pareto Set Identification with Posterior Sampling.
Proceedings of the 28th International Conference on Artificial Intelligence and Statistics (AISTATS 2025).
- **C. Kone**, E. Kaufmann, L. Richert.
Constrained Pareto Set Identification with Bandit Feedback.
Proceedings of the Forty-Second International Conference on Machine Learning (ICML 2025).

Together, these publications form the core of the dissertation:

- Chapter 2—*Fixed-Budget PSI*—is based on our AISTATS 2024 paper introducing the Empirical Gap Elimination framework, the first theoretical treatment of PSI in this regime.
- Chapter 3—*Adaptive Algorithms for Fixed-Confidence PSI*—draws on our NeurIPS 2023 work proposing LUCB-type adaptive elimination strategies and relaxed objectives such as $\text{PSI-}k$.
- Chapter 4—*Multi-Output Linear Models for PSI*—follows our AISTATS 2025 paper extending PSI to structured linear models via optimal experimental design.
- Chapter 5—*Bayesian PSI*—is based on our AISTATS 2025 paper on posterior-sampling-based algorithms achieving asymptotic optimality.
- Chapter 6—*Constrained PSI*—builds on our ICML 2025 paper introducing feasibility-aware Pareto identification with linear constraints.

Additional Publications

- **C. Kone**, K. Jamieson.
Optimal Posterior Sampling for Policy Identification in Tabular Markov Decision Processes.
Proceedings of the 29th International Conference on Artificial Intelligence and Statistics (AISTATS 2026).

Although not directly focused on PSI, this work builds on the posterior sampling principles developed in Chapter 5 to design computationally efficient algorithms for policy identification in reinforcement learning.

Chapter 2

Pareto Set Identification: The Fixed-Budget Setting

This chapter addresses the problem of *Pareto Set Identification* in the *fixed-budget* setting. Our motivation stems from applications such as clinical trials, where the number of participants is pre-specified and the exploration budget is thus limited. The goal is then to identify the Pareto set with minimal error.

We propose *Empirical Gap Elimination* (EGE), a family of algorithms that combine (i) a principled estimation of the difficulty of classifying each arm as Pareto-optimal or not, with (ii) elimination-style strategies tailored to the fixed-budget regime. We provide a theoretical analysis showing that two concrete instances, EGE-SR and EGE-SH, achieve an error probability that decreases exponentially with the budget, with an exponent supported by an information-theoretic lower bound.

This chapter is based on joint work with Émilie Kaufmann and Laura Richert, published in the proceedings of *AISTATS 2024*.

2.1	Introduction	30
2.2	Algorithmic contributions	33
2.2.1	Empirical Gap Elimination	33
2.2.2	Particular instances	35
2.2.3	Beyond exact PSI	36
2.3	Main theoretical results	37
2.3.1	Upper bound on error probability	37
2.3.2	Lower bound on the error probability	39
2.3.3	Sketch of proof of Theorem 2.3.1	40
2.4	Numerical study and discussion	45
2.5	Additional proofs	47
2.5.1	Analysis of EGE	48
2.5.2	Analysis of EGE-SR- k	53
2.5.3	Simplifying the sub-optimality gaps	58

2.1 Introduction

As discussed in Chapter 1, the multi-armed bandit problem has been predominantly studied as a single-objective stochastic optimization problem, mainly for reward maximization, (Lattimore & Szepesvari 2020) or for identifying the arm with the largest expected value (best arm identification, (Even-Dar et al. 2002; Audibert & Bubeck 2010)). However, many practical problems involve multiple, possibly conflicting objectives. This motivates the study of multi-objective pure exploration, where the goal is to identify the set of arms whose mean performance vectors are Pareto optimal.

Previous work has mainly focused on the *fixed-confidence* formulation of Pareto Set Identification (PSI) (Zuluaga, Sergent, et al. 2013; Auer et al. 2016; Zuluaga, Krause, et al. 2016; Ararat & Tekin 2023): given a risk $\delta \in (0, 1)$, the agent stops at a random time τ and recommends a set \hat{S}_τ such that $\mathbb{P}(\hat{S}_\tau \neq \mathcal{S}^*) \leq \delta$, while minimizing its stopping time τ .

In contrast, in the dual *fixed-budget* formulation, the total number of samples T is fixed in advance and the goal is to minimize the misidentification probability. This unexplored setting is particularly relevant in applications such as adaptive clinical trials, where the number of participants is known in advance due to ethical and/or financial constraints. In such contexts, the challenge is to allocate samples efficiently across competing treatment arms to maximize the probability of correctly identifying the most promising vaccine or drug combinations. Beyond clinical trials, similar trade-offs arise in multi-objective design problems such as optimizing hardware configurations for both latency and energy efficiency (Zuluaga, Sergent, et al. 2013), or in A/B/n testing for recommender systems involving multiple engagement metrics (Mehrotra et al. 2020).

In this chapter, we introduce the first algorithms for Pareto Set Identification in the fixed-budget setting. Our approach extends the principles of Successive Rejects (Audibert & Bubeck 2010) and Sequential Halving (Karnin et al. 2013); two canonical methods for best-arm identification, to the multi-objective case. We propose the *Empirical Gap Elimination* (EGE) framework, which relies on a novel arm-classification rule based on *empirical gaps*.

To handle practical situations where only a subset of the Pareto set needs to be identified, we further introduce the *k-relaxation* of PSI, denoted PSI- k . This relaxation requires the learner to return any subset of k Pareto-optimal arms, thereby interpolating between identifying a single good arm ($k = 1$) and recovering the full Pareto set ($k = K$). Such a formulation is particularly meaningful in clinical studies, where one may only need to shortlist a few good-performing strategies for further testing.

We show that the proposed algorithms, EGE-SR and EGE-SH, achieve instance-dependent exponential convergence rates matching known lower bounds up to logarithmic factors (Section 2.2).

Related work. Building on the context introduced in Chapter 1, where we reviewed the theoretical and practical motivations for pure exploration in multi-objective bandits, this chapter focuses on the fixed-budget formulation of Pareto Set Identification (PSI). While several works have addressed PSI in the fixed-confidence setting (Zuluaga, Sergent, et al.

2013; Auer et al. 2016; Zuluaga, Krause, et al. 2016; Ararat & Tekin 2023), no prior work provides theoretical guarantees for PSI when the sampling budget is fixed, despite its practical relevance in resource-limited domains such as clinical trials or recommender systems. In the single-objective case, classical algorithms like Successive Rejects and Sequential Halving achieve exponential decay of the misidentification probability with the budget (Audibert & Bubeck 2010; Karnin et al. 2013), yet no single strategy attains a universally optimal rate across all instances (Degenne 2023).

A common approach to extend fixed-confidence methods to the fixed-budget setting is to tune their risk parameter δ according to the available budget T and the instance hardness $H(\nu)$, as in the UGapEb algorithm (Gabillon et al. 2012).

Indeed, a fixed-confidence algorithm for PSI with stopping time τ typically takes a risk parameter $\delta \in (0, 1)$ as input and satisfies guarantees of the form

$$\mathbb{P}\left(\widehat{S}_\tau = \mathcal{S}^* \text{ and } \tau \leq C H(\nu) \log\left(d \frac{H(\nu)}{\delta}\right)\right) \geq 1 - \delta,$$

where \widehat{S}_τ denotes the estimated Pareto set and $C > 0$ is a universal constant. Given a fixed sampling budget T , one may attempt to tune δ as a function of T and the instance hardness $H(\nu)$ by choosing δ_T such that $C H(\nu) \log(H(\nu)/\delta_T) \leq T$. Under this calibration, the corresponding error probability behaves approximately as

$$\delta_T \approx H(\nu) \exp\left(-\frac{T}{C H(\nu)}\right).$$

However, such oracle tuning presupposes knowledge of $H(\nu)$, which is rarely available in practice. In this chapter, we overcome this limitation by introducing the EGE framework for fixed-budget PSI. Our method adaptively estimates the classification difficulty of each arm and integrates these estimates into elimination-based allocation rules, leading to two practical algorithms: EGE-SR and EGE-SH.

Learning model. We consider a stochastic multi-armed bandit model with K arms, where each arm $a \in [K]$ is associated with an unknown probability distribution ν_a supported on \mathbb{R}^d , $d \geq 1$. The mean vector of arm a is denoted $\mu_a = \mathbb{E}_{X \sim \nu_a}[X] \in \mathbb{R}^d$, and we write $\mu = (\mu_1, \dots, \mu_K)$. Throughout, we assume subgaussian marginals: for every arm a and each coordinate $c \in [d]$, the random variable X_a^c (with $X_a \sim \nu_a$) is σ -subgaussian, that is,

$$\forall \lambda \in \mathbb{R}, \quad \log \mathbb{E}[\exp(\lambda(X_a^c - \mathbb{E}[X_a^c]))] \leq \frac{\lambda^2 \sigma^2}{2},$$

for some common variance proxy $\sigma^2 > 0$.

The learner interacts with this environment over R rounds, sequentially selecting arms from an active set $\mathcal{A}_r \subseteq [K]$ and observing independent samples $(Z_{r,a})_{a \in \mathcal{A}_r}$, where each $Z_{r,a} \sim \nu_a$ satisfies $\mathbb{E}[Z_{r,a} | a] = \mu_a$. Let \mathbb{P}_ν and \mathbb{E}_ν denote respectively the probability law and expectation under the joint process $\{Z_{r,a}\}_{r,a}$, and define the filtration $\mathcal{H}_r = \sigma(\mathcal{A}_1, Z_1, \dots, \mathcal{A}_r, Z_r)$, representing the information available up to round r . A sampling strategy is said to be adaptive if each selection A_t is \mathcal{H}_{t-1} -measurable. Given a

total sampling budget T , the learner outputs an estimator \widehat{S}_T of the true Pareto-optimal set $\mathcal{S}^*(\nu)$, based on the full observation history \mathcal{H}_T . The performance of a strategy is measured by its probability of misidentification

$$e_T(\nu) := \mathbb{P}_\nu \left(\widehat{S}_T \neq \mathcal{S}^*(\nu) \right),$$

which quantifies the probability of returning an incorrect Pareto set under a fixed sampling budget.

Definition 2.1.1. Given two arms $i, j \in [K]$, i is weakly (Pareto) dominated by j (denoted by $\mu_i \leq \mu_j$) if for any $c \in \{1, \dots, d\}$, $\mu_i^c \leq \mu_j^c$. The arm i is (Pareto) dominated by j ($\mu_i \preceq \mu_j$ or $i \preceq j$) if i is weakly dominated by j and there exists $c \in \{1, \dots, d\}$ such that $\mu_i^c < \mu_j^c$. The arm i is strictly (Pareto) dominated by j ($\mu_i \prec \mu_j$ or $i \prec j$) if for any $c \in \{1, \dots, d\}$, $\mu_i^c < \mu_j^c$.

The Pareto set $\mathcal{S}^*(\nu)$ is defined as

$$\mathcal{S}^*(\nu) := \{i \in [K] \text{ s.t. } \nexists j \in [K] : \mu_i \prec \mu_j\},$$

and will be denoted by \mathcal{S}^* when ν is clear from the context. Any arm $a \in \mathcal{S}^*$ will be called (Pareto) *optimal* and an arm $a \notin \mathcal{S}^*$ is called *sub-optimal*.

To characterize the difficulty of the problem, we introduce below some quantities that allow us to measure how much an arm is dominated. For any arms i, j , we let

$$\begin{aligned} m(i, j) &:= \min_{c \leq d} [\mu_j^c - \mu_i^c], \\ M(i, j) &:= \max_{c \leq d} [\mu_i^c - \mu_j^c]. \end{aligned}$$

If $\mu_i \not\prec \mu_j$, then $M(i, j)$ is the smallest uniform increase of j that makes it dominate i . If $\mu_i \prec \mu_j$ then $m(i, j)$ is the smallest increase of any component of i which makes it non-dominated by j .

We recall the gaps introduced in Chapter 1 and proven (Auer et al. 2016) to characterize fixed-confidence PSI. For a sub-optimal arm $i \notin \mathcal{S}^*$,

$$\Delta_i := \Delta_i^* := \max_{j \in \mathcal{S}^*} m(i, j), \quad (2.1)$$

which is the smallest quantity that should be added component-wise to μ_i to make i appear Pareto optimal w.r.t. $\{\mu_k : k \in [K] \setminus \{i\}\}$. For an optimal arm $i \in \mathcal{S}^*$,

$$\Delta_i := \min \{\delta_i^+, \delta_i^-\} \quad (2.2)$$

where

$$\begin{aligned} \delta_i^+ &:= \min_{j \in \mathcal{S}^* \setminus \{i\}} \min [M(i, j), M(j, i)], \\ \delta_i^- &:= \min_{j \in [K] \setminus \mathcal{S}^*} [(M(j, i))_+ + \Delta_j], \end{aligned}$$

with the convention $\min_{\emptyset} = +\infty$. δ_i^+ accounts for how much i is close to dominate (or to be dominated by) another optimal arm, while δ_i^- translates, in a sense, the smallest “margin” from an optimal arm i to the sub-optimal arms. The difficulty (*i.e.*, near-optimal sample complexity) of fixed-confidence PSI is characterized by the complexity term

$$H(\nu) := \sum_{a=1}^K \frac{1}{\Delta_a^2}.$$

In this chapter, we show that it is also a relevant complexity measure for the fixed-budget setting. Our goal is to propose elimination algorithms that estimate the gaps of each arm online and sequentially discard and classify arms with larger gaps.

Toward this goal, we begin by rewriting and simplifying the gaps identified by [Auer et al. 2016](#), which, in their current form, directly depend on the Pareto set, making them apparently difficult to estimate online. Recalling that $(x)_+ := \max\{x, 0\}$, the following holds.

Lemma 2.1.2. *For any arm $k \in [K]$,*

$$\Delta_k = \begin{cases} \Delta_k^* = \max_{j \in [K]} m(k, j) & \text{if } k \notin \mathcal{S}^* \\ \delta_k^* & \text{if } k \in \mathcal{S}^* \end{cases},$$

where $\delta_k^* := \min_{j \neq k} [M(k, j) \wedge (M(j, k)_+ + (\Delta_j^*)_+)]$. In particular, for any arm k , $\Delta_k = \max\{\Delta_k^*, \delta_k^*\}$.

The proof of this result is given in Section 2.5. This rewriting removes the explicit dependency on the Pareto Set. Below, we propose a generic elimination algorithm that relies on a simple empirical estimation of these gaps.

2.2 Algorithmic contributions

We introduce the family of Empirical Gap Elimination (EGE) algorithms. They are generic algorithms that take as input a number of rounds, an elimination scheme (dictating the fraction of arms discarded at the end of each round), and a sampling scheme, dictating the fraction of the budget allocated to each round.

2.2.1 Empirical Gap Elimination

An Empirical Gap Elimination algorithm uses a round-based structure. The algorithm is parameterized by a number of rounds $R > 0$, an arm schedule vector $\lambda = (\lambda_1, \dots, \lambda_R, \lambda_{R+1}) \in [K]^{R+1}$ satisfying $\lambda_r > \lambda_{r+1}$ where λ_r indicates how many arms are active in round r , and an allocation vector $t = (t_1, \dots, t_R) \in [T]^R$, where t_r is the number of samples gathered

from each active arm in round r . Given the maximal budget T , these vectors should further satisfy the following:

$$\lambda_1 = K, \quad \lambda_{R+1} \in \{0, 1\} \quad \text{and} \quad (2.3)$$

$$\sum_{r=1}^R \lambda_r t_r \leq T, \quad (2.4)$$

ideally with an equality.

In each round r , the algorithm maintains a set of active arms, denoted by \mathcal{A}_r , that is of size λ_r , and collects t_r new samples from each arm in \mathcal{A}_r . The total number of samples gathered from each arm in \mathcal{A}_r in the first r rounds¹ is therefore

$$n_r = \sum_{s=1}^r t_s.$$

For each $a \in \mathcal{A}_r$, the empirical estimate of μ_a at the end of round r is denoted by $\hat{\mu}_{r,a} := \frac{1}{n_r} \sum_{s=1}^{n_r} Z_{a,s}$ where $Z_{a,s}$ denote the s -th observation drawn *i.i.d.* from distribution ν_a . These estimates are used to carefully decide which arm to explore in the next round, based on an appropriate notion of *empirical gap*. We introduce the empirical quantities

$$\begin{aligned} m(i, j; r) &:= \min_{c \leq d} [\hat{\mu}_{r,j}^c - \hat{\mu}_{r,i}^c], \\ M(i, j; r) &:= \max_{c \leq d} [\hat{\mu}_{r,i}^c - \hat{\mu}_{r,j}^c] \end{aligned}$$

and the empirical Pareto set at the end of round r

$$\begin{aligned} S_r &:= \{i \in \mathcal{A}_r : \nexists j \in \mathcal{A}_r \text{ such that } \hat{\mu}_{r,i} \prec \hat{\mu}_{r,j}\}, \\ &= \{i \in \mathcal{A}_r : \forall j \in \mathcal{A}_r \setminus \{i\}, M(i, j; r) > 0\}. \end{aligned}$$

Finally, we define for any arm $i \in \mathcal{A}_r$ the empirical gaps

$$\begin{aligned} \hat{\Delta}_{i,r}^* &:= \max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j; r), \\ \hat{\delta}_{i,r}^* &:= \min_{j \in \mathcal{A}_r \setminus \{i\}} [M(i, j; r) \wedge (M(j, i; r)_+ + (\hat{\Delta}_{i,r}^*)_+)], \end{aligned}$$

and for any $i \in \mathcal{A}_r$

$$\hat{\Delta}_{i,r} := \begin{cases} \hat{\Delta}_{i,r}^* & \text{if } i \in \mathcal{A}_r \setminus S_r, \\ \hat{\delta}_{i,r}^* & \text{else.} \end{cases} \quad (2.5)$$

These gaps are empirical variants of the gaps stated in Lemma 2.1.2 evaluated on the active arms. In particular, it is simple to observe that

$$\hat{\Delta}_{i,r} = \max\{\hat{\delta}_{i,r}^*, \hat{\Delta}_{i,r}^*\},$$

Algorithm 2.1: EGE: Empirical Gap Elimination

Require: number of rounds R ; budget allocation t ; and elimination scheme λ

- 1 Initialize: Set of active arms $\mathcal{A}_1 \leftarrow \{1, \dots, K\}$; set of arms (so far) identified as Pareto-optimal $\mathcal{B}_1 \leftarrow \emptyset$; set of discarded arms $\mathcal{D}_1 \leftarrow \emptyset$;
- 2 **foreach** $r = 1, 2, \dots, R$ **do**
- 3 Collect t_r samples from each arm $a \in \mathcal{A}_r$;
- 4 Compute S_r the empirical Pareto set ;
- 5 Let \mathcal{A}_{r+1} be the set of λ_{r+1} arms in \mathcal{A}_r with the smallest empirical gaps $\widehat{\Delta}_{a,r}$;
 // ties broken in favor of arms in S_r
- 6 $\mathcal{B}_{r+1} \leftarrow \mathcal{B}_r \cup \{S_r \cap (\mathcal{A}_r \setminus \mathcal{A}_{r+1})\}$;
- 7 $\mathcal{D}_{r+1} \leftarrow \mathcal{D}_r \cup \{(\mathcal{A}_r \setminus \mathcal{A}_{r+1}) \setminus S_r\}$;
- 8 **return** $\widehat{S}_T = \mathcal{B}_{R+1} \cup \mathcal{A}_{R+1}$

which follows from the fact that when $i \in S_r$, an empirically Pareto-optimal arm $\widehat{\Delta}_{i,r}^* < 0$ and $\widehat{\delta}_{i,r}^* \geq 0$. On the other side, when $i \in \mathcal{A}_r \setminus S_r$, an empirical sub-optimal arm, we have $\widehat{\delta}_{i,r}^* < 0$ while $\widehat{\Delta}_{i,r}^* \geq 0$, as it follows from their definitions.

At the end of round r , the algorithm sorts the arms by increasing order of their empirical gaps: $\widehat{\Delta}_{(1),r} \leq \dots \leq \widehat{\Delta}_{(\lambda_r),r}$ and we define $\mathcal{A}_{r+1} = \{(1), \dots, (\lambda_{r+1})\}$. In case of ties, *i.e.*, if $\widehat{\Delta}_{(\lambda_{r+1}),r} = \widehat{\Delta}_{(\lambda_{r+1}+m),r}$ for some m , we first add arms in S_r to \mathcal{A}_{r+1} . We emphasize that this tie-breaking rule is crucial in our analysis. Arms in $\mathcal{A}_{r+1} \setminus \mathcal{A}_r$ are further classified as optimal (and added to the set \mathcal{B}_{r+1}) or sub-optimal (and added to the set \mathcal{D}_{r+1}) based on whether or not they belong to S_r . The output of the algorithm is the set $\mathcal{B}_{R+1} \cup \mathcal{A}_{R+1}$.

2.2.2 Particular instances

Arm elimination algorithms have been proposed for different unidimensional fixed-budget identification tasks (Audibert & Bubeck 2010; Bubeck, T. Wang, et al. 2013; Karnin et al. 2013) with different elimination rules, that could also be rewritten featuring some (simpler) gaps. In these works, two different types of arm schedules and sampling allocations have been mainly investigated:

- In Successive Rejects (SR), one arm is deactivated in each round, *i.e.*, $R = K - 1$ and $\lambda_r^{\text{SR}} = K - r + 1$ for all $r \leq K$. The sampling allocation proposed by Audibert & Bubeck 2010 satisfies $t_r^{\text{SR}} = n_r^{\text{SR}} - n_{r-1}^{\text{SR}}$ with $n_r = \left\lceil \frac{1}{\log(K)} \frac{T-K}{K+1-r} \right\rceil$ where $\overline{\log(K)} := 2^{-1} + \sum_{i=2}^K i^{-1}$ and $n_0 = 0$.

¹We emphasize that EGE algorithms do not discard samples between rounds

- In Sequential Halving (SH), one half of the active arms is de-activated at the end of each round, that is $R = \lceil \log_2(K) \rceil$ and for all $r \in \{1, \dots, \lceil \log_2(K) \rceil\}$, $\lambda_{r+1}^{\text{SH}} := \lceil \lambda_r^{\text{SH}}/2 \rceil$ (we easily verify that $\lambda_{R+1}^{\text{SH}} = 1$). The sampling allocation proposed by [Karnin et al. 2013](#) is uniformly spread across rounds, that is $t_r^{\text{SH}} := \left\lfloor \frac{T}{|\mathcal{A}_r| \lceil \log_2(K) \rceil} \right\rfloor$.

We refer to EGE-SR (resp. EGE-SH) as the instances of EGE using the same allocation as SR (resp. SH). Another instantiation of EGE may use the geometric allocation proposed by [Karpov & Zhang 2022](#), generalizing the SH-type of allocation.

For $d = 1$, the PSI problem coincides with BAI, and EGE-SR (resp. EGE-SH) coincides with SR (resp. SH²). Indeed, in that case $S_r = \{\hat{a}_r\}$ reduces to the empirical best arm, $\hat{\Delta}_{i,r} = \hat{\mu}_{\hat{a}_r,r} - \hat{\mu}_{i,r}$ for $i \neq \hat{a}_r$, $\hat{\Delta}_{\hat{a}_r,r} = \min_{i \in \mathcal{A}_r \setminus \{\hat{a}_r\}} \hat{\Delta}_{i,r}$. Then, our tie-breaking rule ensures that no arm is accepted before the last round, and at each round r , \mathcal{A}_{r+1} is defined as the λ_{r+1} arms in \mathcal{A}_r with the largest empirical means, and the final survival arm is recommended as optimal.

2.2.3 Beyond exact PSI

We propose an extension of EGE-SR to the PSI- k problem. In PSI- k , the learner aims to return a subset $\hat{\mathcal{S}}_T \subseteq \mathcal{S}^*$ of size k , or the full Pareto set if $|\mathcal{S}^*| < k$. In the fixed-budget setting, the performance of an algorithm is measured through the k -relaxed expected loss:

$$e_{T,k}(\nu) := \mathbb{E}_\nu[\mathcal{L}(\hat{\mathcal{S}}_T, k)],$$

where

$$\mathcal{L}(\hat{\mathcal{S}}, k) := \begin{cases} \mathbb{1}(\hat{\mathcal{S}} \subset \mathcal{S}^*), & \text{if } |\hat{\mathcal{S}}| = k, \\ \mathbb{1}(\hat{\mathcal{S}} = \mathcal{S}^*), & \text{otherwise.} \end{cases}$$

This loss penalizes the inclusion of any sub-optimal arm while tolerating partial identification of the Pareto set up to size k .

The EGE-SR algorithm can be naturally extended to PSI- k by adjusting its stopping condition to

$$\tau := \inf\{r : |\mathcal{B}_{r+1}| = k\} \wedge (K - 1),$$

where \mathcal{B}_r denotes the set of arms declared optimal after round r . Let N_τ be the total number of samples drawn at termination. In Section 2.5.2, we show that when the budget is large,

$$\mathbb{E}_\nu[\tau] \lesssim q \quad \text{and} \quad \mathbb{E}_\nu[N_\tau] \lesssim N_q,$$

where $q := K - |\mathcal{S}^*| + k$. Intuitively, this means that in the worst case, the $(K - |\mathcal{S}^*|)$ sub-optimal arms are eliminated before k optimal ones are confirmed, as some optimal arms may be needed to correctly classify borderline sub-optimal arms. Notably, q can be significantly smaller than $(K - 1)$; for instance, when all arms are Pareto-optimal ($[K] = \mathcal{S}^*$), we have $q = k$, yielding a substantial reduction in sample complexity.

²Besides the fact that the original SH algorithm for BAI discards samples collected in previous rounds to compute the empirical means.

Algorithm 2.2: EGE-SR- k : Any k -sized subset of the Pareto Set

Require: Parameter $k \in [K]$; number of rounds R ; budget allocation t ; and elimination scheme λ

- 1 Initialize: Set of active arms $\mathcal{A}_1 \leftarrow \{1, \dots, K\}$; set of arms (so far) identified as Pareto-optimal $\mathcal{B}_1 \leftarrow \emptyset$; set of discarded arms $\mathcal{D}_1 \leftarrow \emptyset$;
- 2 **foreach** $r = 1, 2, \dots, K - 1$ **do**
- 3 Collect t_r samples from each arm $a \in \mathcal{A}_r$;
- 4 Compute S_r the empirical Pareto set
- 5 Choose $i_r \in \operatorname{argmax}_{i \in \mathcal{A}_r} \widehat{\Delta}_{i,r}$;
 // ties broken in favor of arms in S_r
- 6 $\mathcal{A}_{r+1} \leftarrow \mathcal{A}_r \setminus \{i_r\}$;
- 7 $\mathcal{B}_{r+1} \leftarrow \mathcal{B}_r \cup \{S_r \cap \{i_r\}\}$;
- 8 $\mathcal{D}_{r+1} \leftarrow \mathcal{D}_r \cup \{\{i_r\} \cap (\mathcal{A}_r \setminus S_r)\}$;
- 9 **if** $|\mathcal{B}_{r+1}| = k$ **then**
- 10 **break and return** $\widehat{S}_T \leftarrow \mathcal{B}_{r+1}$;

// the identified Pareto set contains k or fewer arms

- 11 **return** $\widehat{S}_T = \mathcal{B}_K \cup \mathcal{A}_K$

2.3 Main theoretical results

We present the main guarantees of EGE algorithms.

2.3.1 Upper bound on error probability

Exact PSI. We first propose an analysis of EGE for a generic number of rounds R , arm schedule λ , and sampling allocation t satisfying (2.3) and (2.4). It features the quantity

$$\widetilde{T}^{R,t,\lambda}(\nu) := \min_{r \in [R]} \left[\sum_{s=1}^r t_s \right] \cdot \Delta_{(\lambda_{r+1}+1)},$$

in which the dependency in ν is captured in the gaps.

Theorem 2.3.1. *Let a budget $T \geq K$ be given and ν be a bandit with marginally σ -subgaussian arms. Then Empirical Gap Elimination with number of rounds R , budget allocation t , and elimination scheme λ satisfies*

$$e_T^{EGE}(\nu) \leq 2(K-1) |\mathcal{S}^*| R d \cdot e^{-\frac{\widetilde{T}^{R,t,\lambda}(\nu)}{144\sigma^2}}.$$

This result shows that the probability of failure of EGE decreases exponentially fast with $\tilde{T}^{R,t,\lambda}(\nu)$.

Introducing

$$H_2(\nu) := \max_{k \in [K]} k \Delta_{(k)}^{-2},$$

where $k \mapsto (k)$ is a permutation such that $\Delta_{(1)} \leq \dots \leq \Delta_{(K)}$, we obtain the following.

Corollary 2.3.2. *Let $T \geq K$ and ν be a bandit with σ -subgaussian marginals. Then EGE-SR satisfies*

$$e_T^{\text{SR}}(\nu) \leq 2(K-1)^2 |\mathcal{S}^*| d \cdot e^{-\frac{T-K}{144\sigma^2 H_2(\nu) \lceil \log(K) \rceil}},$$

and for EGE-SH, $e_T^{\text{SH}}(\nu)$ is upper-bounded by

$$2(K-1) \lceil \log_2(K) \rceil |\mathcal{S}^*| d \cdot e^{-\frac{T}{288\sigma^2 H_2(\nu) \lceil \log_2(K) \rceil}}.$$

Proof. For EGE-SR, we have $R = K - 1$, $\lambda_r^{\text{SR}} = K + 1 - r$ and $t_r^{\text{SR}} = n_r^{\text{SR}} - n_{r-1}^{\text{SR}}$ where $n_r^{\text{SR}} = \left\lceil \frac{1}{\log(K)} \frac{T-K}{K+1-r} \right\rceil$, which yields

$$\begin{aligned} \tilde{T}^{R,t^{\text{SR}},\lambda^{\text{SR}}}(\nu) &:= \min_{r \in [K-1]} n_r^{\text{SR}} \Delta_{(\lambda_r^{\text{SR}}+1)}^2, \\ &\geq \min_{r \in [K-1]} \frac{\Delta_{(K+1-r)}^2}{\log(K)} \frac{T-K}{K+1-r}, \\ &= \frac{T-K}{\log(K)} \frac{1}{\max_{r \in \{2, \dots, K\}} r \Delta_{(r)}^{-2}} \\ &= \frac{T-K}{\log(K)} \frac{1}{H_2(\nu)}. \end{aligned}$$

The analysis of EGE-SH is proven similarly in Section 2.5. □

The complexity measure $H_2(\nu)$ featured in our error exponent satisfies $H_2(\nu) \leq H(\nu) \leq H_2(\nu) \cdot \log(2K)$, following similar steps as in [Audibert & Bubeck 2010](#). For BAI ($d = 1$), we essentially recover the existing guarantees for SR and SH, whose error bounds also feature $H_2(\nu)$, up to constant factors inside the exponential and an extra multiplicative K factor for Sequential Halving. Still, to our knowledge, this is the first analysis of the variant of SH that does not discard samples between rounds, which often performs (much) better in practice.

In the general case ($d \geq 1$), we note that the bounds obtained for EGE-SR and EGE-SH are hard to compare: the latter has an improved polynomial dependence ($K \log_2(K)$ instead of K^2) but a worse constant inside the exponential. As we shall see in the experiments, both algorithms have actually pretty close performance (like in BAI, see [Karnin et al. 2013](#)). Moreover, they both outperform a simple baseline using Uniform Allocation (UA) and recommending the Pareto set of the empirical means. We note that Theorem 2.3.1 yields an upper bound on the error probability of this strategy (by choosing $R = 1$, $t_1 = T/K$ and $\lambda_2 = 0$, for which $\tilde{T}^{1,t,\lambda}(\nu) = n_1 \Delta_{(\lambda_2+1)}^2 = T/(K \Delta_{(1)}^{-2})$), which we add to our summary in Table 2.1. This bound can be much worse than that for EGE-SR/SH when the gaps are distinct.

Table 2.1: Upper bounds on $e_T(\nu)$ for different algorithms (up to constants). APE-FB is an "oracle" algorithm introduced in [Kone, Kaufmann, et al. 2024](#) by converting a fixed-confidence PSI algorithm. APE-FB requires the parameter $H(\nu)$ to run.

Algorithm	Error Probability
EGE-SR	$K^2 \mathcal{S}^* d \cdot e^{-T/(H_2(\nu) \log K)}$
EGE-SH	$K \log(K) \mathcal{S}^* d \cdot e^{-T/(2H_2(\nu) \log K)}$
APE-FB*	$K \log(T) d \cdot e^{-T/H(\nu)}$
UA	$K \mathcal{S}^* d \cdot e^{-T/(K\Delta_{(1)}^{-2})}$

PSI- k relaxation. To introduce the upper-bound on the expected loss $e_{T,k}(\nu)$, we define $\omega_{(k)}$ to be the k -th largest gap among the optimal arms: $\omega_{(k)} := \max_{i \in \mathcal{S}^*}^k \Delta_i$ with $\omega_{(k)} = 0$ if $|\mathcal{S}^*| < k$. Our bound features the complexity measure $H_2^{(k)}(\nu) := \max_{i \in [K]} i(\Delta_{(i)}^{(k)})^{-2}$ with the k -relaxed gaps

$$\Delta_i^{(k)} := \begin{cases} \max\{\Delta_i, \omega_{(k)}\} & \text{if } i \in \mathcal{S}^* \\ \Delta_i & \text{else.} \end{cases} \quad (2.6)$$

Theorem 2.3.3. *Let $k \in [K]$. EGE-SR- k satisfies*

$$e_{T,k}(\nu) \leq 2(K-1)^2 |\mathcal{S}^*| d \cdot e^{-\frac{T-K}{144\sigma^2 H_2^{(k)}(\nu) \log(K)}}.$$

This result is particularly insightful when there are many optimal arms and some of them are easy to identify as such (large gaps). Indeed when $|\mathcal{S}^*| \approx K$ and $k \ll |\mathcal{S}^*|$, $H_2^{(k)}(\nu)$ can be an order of magnitude smaller than $H_2(\nu)$. We also note that when $k > |\mathcal{S}^*|$ (and PSI- k reduces to PSI), $H_2^{(k)}(\nu) = H_2(\nu)$ and we recover the result of [Theorem 2.3.2](#).

The theorem below bounds the expected stopping time and the number of samples used at stopping.

Theorem 2.3.4. *Fix $k < |\mathcal{S}^*|$ and let $q := K - |\mathcal{S}^*| + k$. Then*

$$\begin{aligned} \mathbb{E}[\tau] &\leq q + 2(K-1)|\mathcal{S}^*|(K-q-1)qde^{-\frac{T-K}{144\sigma^2 H_2^{(k)}(\nu) \log(K)}} \quad \text{and} \\ \mathbb{E}[N_\tau] &\leq N_q + 2(K-1)|\mathcal{S}^*|(K-q-1)qdT e^{-\frac{T-K}{144\sigma^2 H_2^{(k)}(\nu) \log(K)}}. \end{aligned}$$

This result suggests that for this relaxed problem, the algorithm might not need to use the whole budget, in particular when T is large. For example, in the case where $[K] = \mathcal{S}^*$, so $q = k$ and we roughly use N_k samples which can be way smaller than N_{K-1} , still, the probability of error improves upon the exact PSI setting.

2.3.2 Lower bound on the error probability

We present a lower bound for a class of instances. We define \mathcal{I}^K to be the set of means $\mu := (\mu_1, \dots, \mu_K)$ such that each sub-optimal arm i is only dominated by a single arm,

denoted by i^* (that has to belong to \mathcal{S}^*) and that for each optimal arm j there exists a unique sub-optimal arm which is dominated by j , denoted by \underline{j} . We further assume that optimal arms are not too close to arms that they do not dominate: for any sub-optimal arm i and optimal arm j such that $\mu_i \not\prec \mu_j$,

$$M(i, j) \geq 3 \max(\Delta_i, \Delta_j).$$

Let $\nu := (\nu_1, \dots, \nu_K)$ be an instance whose means $\mu \in \mathcal{I}^K$ and such that $\nu_i \sim \mathcal{N}(\mu_i, \sigma^2 I)$. For every $i \in [K]$ we define the alternative instance $\nu^{(i)} := (\nu_1, \dots, \nu_i^{(i)}, \dots, \nu_K)$ in which only the mean of arm i is modified to:

$$\mu_i^{(i)} := \begin{cases} \mu_i - 2\Delta_i e_{c_i} & \text{if } i \in \mathcal{S}^*(\nu), \\ \mu_i + 2\Delta_i e_{c_i} & \text{else,} \end{cases} \quad (2.7)$$

where e_1, \dots, e_d denotes the canonical basis of \mathbb{R}^d and $c_i := \operatorname{argmin}_c [\mu_{i^*}^c - \mu_i^c]$, with $\nu^{(0)} := \nu$. The following result is proven in [Kone, Kaufmann, et al. 2024](#).

Theorem 2.3.5 (Theorem 2 of [Kone, Kaufmann, et al. 2024](#)). *Let $\mu := (\mu_1, \dots, \mu_K) \in \mathcal{I}^K$ and $\nu := (\nu_1, \dots, \nu_K)$ where $\nu_i \sim \mathcal{N}(\mu_i, \sigma^2 I)$. For any algorithm alg , there exists $i \in \{0, \dots, K\}$ such that*

$$e_T^{\text{alg}}(\nu^{(i)}) \geq \frac{1}{4} e^{-\frac{2T}{\sigma^2 H(\nu^{(i)})}}.$$

In particular, there exists some instances $\tilde{\nu} \in \mathcal{M}^K$ such that $e_T^{\text{alg}}(\tilde{\nu}) \geq \frac{1}{4} \exp(-\frac{2T}{\sigma^2 H(\tilde{\nu})})$. In such instances, the decay rate of EGE-SR and EGE-SH is optimal up to constants and $\log(K)$ factors, and that of APE-FB is optimal up to constant factors, when the complexity is known.

2.3.3 Sketch of proof of Theorem 2.3.1

We define for any arms i, j and round r the events

$$\begin{aligned} \xi_{i,j,r} &:= \left\{ \left\| (\hat{\mu}_{i,n_r} - \hat{\mu}_{j,n_r}) - (\mu_i - \mu_j) \right\|_\infty \leq \gamma \Delta_{(\lambda_{r+1}+1)} \right\} \\ \mathcal{E}_\gamma^1 &:= \bigcap_{r \in [R]} \bigcap_{i \in \mathcal{S}^*} \bigcap_{j \in [K]} \xi_{i,j,r}, \text{ for any } \gamma > 0. \end{aligned}$$

We shall prove that there exists some $\gamma > 0$ such that EGE does not make any error on the event \mathcal{E}_γ^1 . That is, no sub-optimal arm is added to $(\mathcal{B}_r)_{1 \leq r \leq R}$ and no optimal arm is added to $(\mathcal{D}_r)_{1 \leq r \leq R}$, in any round r , and the possibly remaining arm in \mathcal{A}_{R+1} is a Pareto-optimal arm.

To do so, an important step is to justify that any sub-optimal arm should be deactivated before the optimal arm that dominates it the most. More formally, for any sub-optimal arm i , we let

$$i^* \in \operatorname{argmax}_{j \in \mathcal{S}^*} m(i, j),$$

which by definition is such that $\Delta_i = m(i, i^*)$. For a sub-optimal arm i , we know that $i^* \in \mathcal{S}^*$ always exists. More importantly, i^* could be the only arm dominating i . Therefore,

it is crucial to ensure that i is no longer active before discarding i^* , otherwise i could appear as optimal *w.r.t.* the remaining active arms.

Defining

$$\mathcal{P}_r := \{ \forall i \notin \mathcal{S}_*, i \in \mathcal{A}_r \Rightarrow i^* \in \mathcal{A}_r \}, \quad (2.8)$$

we state the following concentration result. The first one controls the deviation of the empirical $M(i, j; r)$ and $m(i, j; r)$ from their actual values.

Lemma 2.3.6. *On the event \mathcal{E}_γ^1 , for all $r \in [R]$ and $i, j \in \mathcal{A}_r$, if $i \in \mathcal{S}^*$ or $j \in \mathcal{S}^*$ then,*

$$\begin{aligned} |M(i, j; r) - M(i, j)| &\leq \gamma \Delta_{(\lambda_{r+1}+1)} \text{ and} \\ |m(i, j; r) - m(i, j)| &\leq \gamma \Delta_{(\lambda_{r+1}+1)}. \end{aligned}$$

The second one builds on it to prove that the empirical gaps cannot be too far from their true gaps, if \mathcal{P}_r holds. It is complementary to Lemma 2.5.2 stated in the main paper.

Lemma 2.3.7. *Let $\gamma > 0$ and assume \mathcal{E}_γ^1 holds. At round $r \in [R]$, for any sub-optimal arm $i \in \mathcal{A}_r$, if $i^* \in \mathcal{A}_r$ and i^* does not empirically dominate i then $\Delta_i^* \leq \gamma \Delta_{(\lambda_{r+1}+1)}$.*

Lemma 2.3.6, Lemma 2.3.7, and Lemma 2.3.8 are proven at the end of the chapter.

Lemma 2.3.8. *Assume that \mathcal{E}_γ^1 holds and let $r \in [R]$ such that \mathcal{P}_r holds. Then for any $i \in \mathcal{A}_r$*

$$\widehat{\Delta}_{i,r} - \Delta_i \geq \begin{cases} -2\gamma \Delta_{(\lambda_{r+1}+1)} & \text{if } i \in \mathcal{S}^*, \\ -\gamma \Delta_{(\lambda_{r+1}+1)} & \text{else.} \end{cases}$$

This result then permits us to prove by induction that \mathcal{P}_r holds in any round r , when γ is small enough.

Lemma 2.3.9. *Let $\gamma < 1/6$. On the event \mathcal{E}_γ^1 , for any $r \in [R + 1]$, \mathcal{P}_r holds. In particular, for any sub-optimal arm i , i^* cannot be deactivated before i .*

Proof. We assume that the event \mathcal{E}_γ^1 holds and we prove the result by induction on r . \mathcal{P}_r trivially holds for $r = 1$ as all arms are active. Let $r \geq 1$ such that \mathcal{P}_r holds. We shall prove that for any $i \in \mathcal{A}_{r+1} \cap (\mathcal{S}^*)^c$, i^* cannot be deactivated at the end of round r . A first observation is that the empirical gap $\widehat{\Delta}_{i,r}$ of each arm i satisfies

$$\widehat{\Delta}_{i,r} = \max \left\{ \widehat{\Delta}_{i,r}^*, \widehat{\delta}_{i,r}^* \right\} \quad (2.9)$$

which follows from the fact that when $i \in \mathcal{S}_r$, $\widehat{\Delta}_{i,r}^* < 0 \leq \widehat{\delta}_{i,r}^*$ and when $i \in \mathcal{A}_r \setminus \mathcal{S}_r$, $\widehat{\delta}_{i,r}^* < 0 \leq \widehat{\Delta}_{i,r}^*$. Using Lemma 2.3.8 and the inductive hypothesis enables us to prove that

$$\forall i \in \mathcal{A}_r, \widehat{\Delta}_{i,r} \geq \Delta_i - 2\gamma \Delta_{(\lambda_{r+1}+1)}. \quad (2.10)$$

To prove that \mathcal{P}_{r+1} holds, we proceed by contradiction and assume that there exists a sub-optimal arm $i \in \mathcal{A}_r$ (and therefore $i^* \in \mathcal{A}_r$ by \mathcal{P}_r) such that $i \in \mathcal{A}_{r+1}$ and $i^* \in \mathcal{A}_r \setminus \mathcal{A}_{r+1}$.

A first observation is that as there are λ_{r+1} arms in \mathcal{A}_{r+1} and i^* is deactivated at the end of round r , there exists an arm $a_r \in \mathcal{A}_{r+1} \cup \{i^*\}$ such that $\Delta_{a_r} \geq \Delta_{(\lambda_{r+1}+1)}$ and $\widehat{\Delta}_{i^*} \geq \widehat{\Delta}_{a_r}$. We now consider two cases depending on whether i^* is empirically optimal or empirically sub-optimal.

Case 1: arm $i^* \notin S_r$ i.e., i^* is empirically sub-optimal. We have

$$\max_{j \in \mathcal{A}_r \setminus \{i^*\}} m(i^*, j; r) := \widehat{\Delta}_{i^*, r}^* = \widehat{\Delta}_{i^*, r} \geq \widehat{\Delta}_{a_r, r}$$

Next, using Lemma 2.3.6 (on the LHS) and Equation (2.10) (on the RHS) yields

$$\begin{aligned} \max_{j \in \mathcal{A}_r \setminus \{i^*\}} m(i^*, j) &\geq \Delta_{a_r} - 3\gamma \Delta_{(\lambda_{r+1}+1)}, \\ &\geq (1 - 3\gamma) \Delta_{(\lambda_{r+1}+1)}, \end{aligned}$$

where the last inequality follows since $\Delta_{a_r} \geq \Delta_{(\lambda_{r+1}+1)}$. As i^* is optimal, the LHS of the previous inequality is negative. So it follows that $0 \geq (1 - 3\gamma) \Delta_{(\lambda_{r+1}+1)}$ which yields a contradiction if $3\gamma < 1$.

Case 2: arm $i^* \in S_r$ i.e., i^* is empirically optimal.

We first prove that i^* does not empirically dominate i . Indeed, if i were dominated by i^* , we would have $i \notin S_r$, so $\widehat{\Delta}_{i, r} = \widehat{\Delta}_{i, r}^* > 0$ and since $i^* \in S_r$,

$$\widehat{\Delta}_{i^*, r} = \widehat{\delta}_{i^*, r}^* \leq (M(i, i^*; r))_+ + (\widehat{\Delta}_{i, r}^*)_+, \quad (2.11)$$

$$= 0 + \Delta_{i, r}. \quad (2.12)$$

Recalling that $i \in \mathcal{A}_{r+1}$ we also have $\widehat{\Delta}_{i^*, r} \geq \widehat{\Delta}_{i, r}$, which combined with Equation (2.12) yields

$$\widehat{\Delta}_{i^*, r} = \widehat{\Delta}_{i, r}. \quad (2.13)$$

However, the tie-breaking rule of EGE ensures that Equation (2.13) is not possible since $i^* \in S_r$ and it is deactivated while $i \in \mathcal{A}_{r+1} \cap (S_r)^c$ (in case of an equality in the gaps, empirically sub-optimal arms are removed). Therefore i is not empirically dominated by i^* . Hence, by Lemma 2.3.7

$$\Delta_i^* \leq \gamma \Delta_{(\lambda_{r+1}+1)}. \quad (2.14)$$

Moreover, since $a_r \in \mathcal{A}_{r+1} \cup \{i^*\}$, using the definition of $\widehat{\Delta}_{i^*, r}$ yields

$$M(i, i^*; r)_+ + (\widehat{\Delta}_{i, r}^*)_+ \geq \widehat{\Delta}_{a_r, r}$$

Using Lemma 2.3.6, Lemma 2.3.8 applied to i and Equation (2.10) applied to a_r , we obtain

$$\begin{aligned} M(i, i^*)_+ + \Delta_i^* &\geq \Delta_{a_r} - 5\gamma \Delta_{(\lambda_{r+1}+1)} \\ &\geq (1 - 5\gamma) \Delta_{(\lambda_{r+1}+1)} \end{aligned}$$

so

$$\Delta_i^* \geq (1 - 5\gamma)\Delta_{(\lambda_{r+1}+1)}. \quad (2.15)$$

Combining this inequality with Equation (2.14) yields a contradiction if $6\gamma < 1$. \square

As an immediate consequence of Lemma 2.3.9, if \mathcal{A}_{R+1} contains one arm (i.e., $\lambda_{R+1} = 1$) and \mathcal{E}_γ^1 holds, then it is an optimal arm.

Note that in BAI ($d = 1$), for any sub-optimal arm i , i^* is the unique best arm, and Lemma 2.3.9 is sufficient to ensure the correctness of EGE on the good event. However, in the general case, we must ensure that no optimal arm is rejected and no suboptimal arm is accepted. This is proven in the lemma below.

Lemma 2.3.10. *Let $\gamma < 1/6$. Then, on the event \mathcal{E}_γ^1 , $\widehat{S}_T = \mathcal{S}^*$.*

Proof. We show that on \mathcal{E}_γ^1 , every arm is well classified, so the recommended set is the true Pareto optimal set. First by Lemma 2.3.9 we know that on \mathcal{E}_γ^1 , the property

$$(i \in \mathcal{A}_r \cap (\mathcal{S}^*)^c \Rightarrow i^* \in \mathcal{A}_r)$$

holds for every $r \in [R]$. As in the proof of Lemma 2.3.9 (see Equation (2.10)), this enables us to prove in Lemma 2.3.8

$$\forall r \in [R], \forall i \in \mathcal{A}_r, \widehat{\Delta}_{i,r} \geq \Delta_i - 2\gamma\Delta_{(\lambda_{r+1}+1)} \quad (2.16)$$

We now establish that at the end of every round $r \in [R]$, no mis-classification can occur. That is for every arm $i \in \mathcal{A}_r \setminus \mathcal{A}_{r+1}$:

- a) if $i \notin S_r$ (that is, i is added to \mathcal{D}_r) then $i \notin \mathcal{S}^*$,
- b) if $i \in S_r$ (that is, i is added to \mathcal{B}_r) then $i \in \mathcal{S}^*$.

Let $i \in \mathcal{A}_r \setminus \mathcal{A}_{r+1}$. Since $\lambda_{r+1} = |\mathcal{A}_{r+1}|$ should remain active at the end of round r , if i is removed at the end of the round, there exists $a_r \in \mathcal{A}_{r+1} \cup \{i\}$ such that $\Delta_{a_r} \geq \Delta_{(\lambda_{r+1}+1)}$ and a_r has a smaller empirical gap than i .

Case 1: $i \notin S_r$ i.e., i is empirically sub-optimal. As $a_r \in \mathcal{A}_{r+1} \cup \{i\}$ has a smaller empirical gap than i , we have

$$\widehat{\Delta}_{i,r}^* \geq \widehat{\Delta}_{a_r,r},$$

Assume by contradiction $i \in \mathcal{S}^*$. Using Lemma 2.3.6 and Equation (2.16) yields

$$\max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j) \geq \Delta_{a_r} - 3\gamma\Delta_{(\lambda_{r+1}+1)} \geq (1 - 3\gamma)\Delta_{(\lambda_{r+1}+1)}, \quad (2.17)$$

When $3\gamma < 1$, the RHS of Equation 2.17 is positive, so this inequality yields

$$\max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j) > 0,$$

that is there exists $j \in \mathcal{A}_r$ such that

$$m(i, j) > 0$$

so there exists j such that $\mu_i \prec \mu_j$, which contradicts the assumption $i \in \mathcal{S}^*$. Therefore, $i \notin \mathcal{S}^*$: i is a sub-optimal arm.

Case 2: $i \in S_r$ *i.e.*, i is empirically optimal. Since a_r has a larger empirical gap than i and $i \in S_r$, we have

$$\min_{j \in \mathcal{A}_r \setminus \{i\}} M(i, j; r) \geq \widehat{\delta}_i^* \geq \widehat{\Delta}_{a_r, r} \quad (2.18)$$

Assume by contradiction that i is sub-optimal. By Lemma 2.3.9 (for $\gamma < 1/6$), $i^* \in \mathcal{A}_r$. Combining with (2.18) yields

$$M(i, i^*; r) \geq \widehat{\Delta}_{a_r, r}.$$

Further using Lemma 2.3.6 and Equation (2.16) yields

$$M(i, i^*) \geq (1 - 3\gamma)\Delta_{(\lambda_{r+1}+1)}.$$

Then by taking $6\gamma < 1$, the RHS of the inequality is positive so

$$M(i, i^*) > 0, \quad (2.19)$$

which is not possible as $\mu_i \prec \mu_{i^*}$. Therefore, $i \in \mathcal{S}^*$: it is an optimal arm.

This proves that on the event \mathcal{E}_γ^1 , with $\gamma < 1/6$ we have $\mathcal{B}_R \subseteq \mathcal{S}^*$ and $\mathcal{D}_R \subseteq (\mathcal{S}^*)^c$. Moreover, if there is a remaining active arm in \mathcal{A}_{R+1} (which happens if and only if $\lambda_{R+1} = 1$), this arm has to be optimal by Lemma 2.3.9. Therefore, on \mathcal{E}_γ^1 , the set recommended by EGE is $\widehat{\mathcal{S}}_T = \mathcal{B}_R \cup \mathcal{A}_{R+1} = \mathcal{S}^*$. \square

Thanks to Hoeffding's inequality, we have

$$\begin{aligned} \mathbb{P}((\mathcal{E}_\gamma^1)^c) &\leq \sum_{r=1}^R \sum_{i \in \mathcal{S}^*} \sum_{j \neq i} \sum_{c=1}^d \mathbb{P}(|(\widehat{\mu}_{i, n_r}^c - \widehat{\mu}_{j, n_r}^c) - (\mu_i^c - \mu_j^c)| > \gamma \Delta_{\lambda_{r+1}+1}), \\ &\leq 2|\mathcal{S}^*|(K-1)d \sum_r \exp\left(-\frac{\gamma^2 n_r \Delta_{\lambda_{r+1}+1}^2}{4\sigma^2}\right) \quad (\text{Hoeffding's inequality}), \\ &\leq 2|\mathcal{S}^*|(K-1)dR \exp\left(-\frac{\gamma^2 \widetilde{T}^{R,t,\lambda}(\nu)}{4\sigma^2}\right). \end{aligned}$$

Applying the above to $\gamma_x = 1/6 - x$ for $x \geq 0$, and letting $x \rightarrow 0$ yields

$$e_T^{\text{EGE}}(\nu) \leq 2|\mathcal{S}^*|(K-1)dR \exp\left(-\frac{\widetilde{T}^{R,t,\lambda}}{144\sigma^2}\right).$$

The analysis of EGE-SR- k is given at the end of the chapter.

2.4 Numerical study and discussion

Experimental setup. We evaluate our algorithms on both synthetic and real-world-inspired instances, comparing them to Uniform Allocation (UA) and the fixed-budget APE-FB algorithm (Kone, Kaufmann, et al. 2023). APE-FB depends on a parameter $a_\alpha = \alpha \frac{25}{36} \frac{T-K}{H(\nu)}$ with $\alpha \in 1/10, 1, 10$, for which guarantees hold when $\alpha \leq 1$ and are optimal at $\alpha = 1$. We also include a heuristic variant, APE-FB-ADAPT, which estimates $H(\nu)$ online. Each experiment is averaged over 4,000 random runs, reporting the \log_{10} of the mean misidentification rate. Following Audibert & Bubeck 2010; Karnin et al. 2013, we set the sampling budget to $T = H(\nu)$. The R -round implementation of EGE has time complexity $\mathcal{O}(RK^2d)$ and memory complexity $\mathcal{O}(K^2)$.

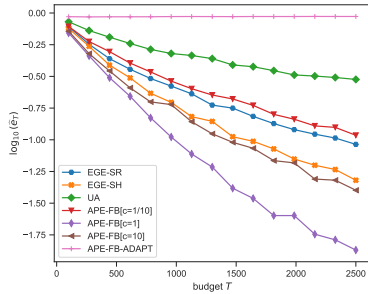


Figure 2.1: Application 1: COV-BOOST trial

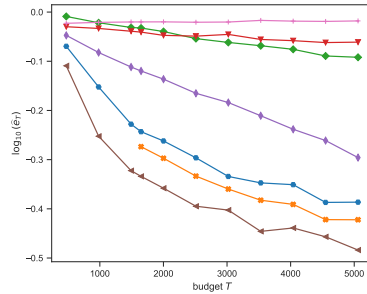


Figure 2.2: Application 2: Sorting Networks dataset.

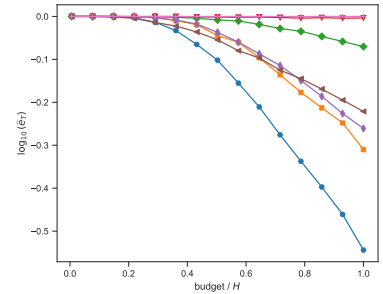


Figure 2.3: Arms on a convex Pareto set.

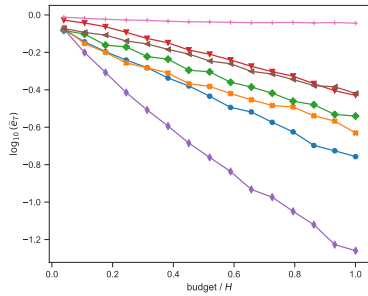


Figure 2.4: Each sub-optimal i is only dominated by i^* .

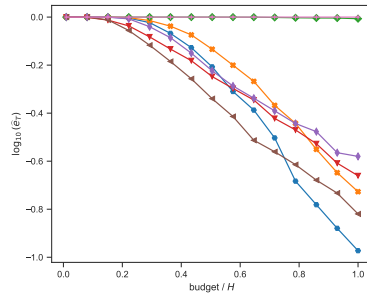


Figure 2.5: $K = 200$ arms on the unit circle.

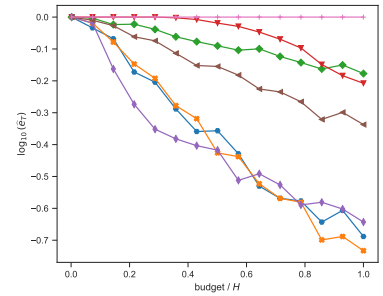


Figure 2.6: High dimension ($d = 10$) with 2 group of arms.

Datasets. Experiments are conducted on two real-world-inspired datasets. The first, COV-BOOST (Munro et al. 2021), is a phase II clinical trial on 2,883 participants evaluating the immunogenicity of $K = 20$ COVID-19 vaccination strategies, each defined by a triplet of vaccines used as first, second, and third doses. Using the three reported immunogenicity indicators—two antibody-based and one cellular—we construct a $d = 3$ Gaussian bandit under the log-normal assumption of Munro et al. 2021. The resulting model, reused in

subsequent chapters, simulates a realistic multi-objective pure-exploration task; full preprocessing details and statistics are given in Appendix 7. The second dataset, SNW (Zuluaga, Milder, et al. 2012), contains 206 sorting network architectures evaluated by two conflicting objectives: FPGA area (resource usage) and throughput (processing rate). Measurements are stochastic due to circuit randomness, and the data are modeled as a $d = 2$, $K = 206$ Gaussian bandit representing the trade-off between hardware cost and performance, with a fixed sampling budget of $T = 5,000$.

Synthetic instances. We test all algorithms on four synthetic instances covering various Pareto-set geometries, dimensions, and noise structures. The first instance (*Convex front*) has $K = 60$ arms in $d = 2$, with the first 20 arms placed on a smooth convex frontier $\mu_i = (x_i^2, 1/(4x_i^2))^\top$ for $x_i \in [0.55, 0.95]$, and the remaining 40 sampled from $\{(x, y) \in [0.1, 0.8]^2 : xy \leq 1/5\}$, forming a dense dominated region. The second instance (*Paired dominance*) includes $K = 10$ arms with $|\mathcal{S}^*| = 2$, where each sub-optimal arm is dominated by a unique optimal one. We set $\mu_1 = (0.4, 0.75)^\top$, $\mu_2 = (0.75, 0.4)^\top$, and define $\mu_{2i+1} = (0.45 + 0.2^i, 0.35 - 0.2^i)^\top$, $\mu_{2i+2} = (0.10 + 0.20^i, 0.70 - 0.20^i)^\top$ for $i = 1, \dots, 4$. The third instance (*Circular set*) has $K = 200$ arms in $d = 2$, uniformly distributed on a unit circle with Gaussian noise $\sigma = 1/4$, where $\mu_i = (\cos \beta_i, \sin \beta_i)^\top$ for β_i uniformly sampled in two disjoint arcs. Finally, the fourth instance (*High-dimensional clusters*) explores $d = 10$ with $K = 50$ arms, the first 30 drawn uniformly from $[0.2, 0.45]^{10}$ and the remaining 20 from $[0.55, 0.75]^{10}$, forming two well-separated clusters.

Results. In all experiments, the uniform allocation (UA) baseline is largely outperformed by both EGE-SR and EGE-SH (with no clear ordering between the two algorithms). This is particularly the case when there are many arms and for complex Pareto sets. By estimating the hardness to classify each arm, EGE eventually allocates more samples to arms that are difficult to classify, leading to a smaller error probability. APE-FB is a good competitor to our EGE algorithms; however, it is not robust to the hyper-parameter a , which requires the knowledge of $H(\nu)$ to be properly tuned. Our proposed heuristic that estimates the complexity online fails dramatically.

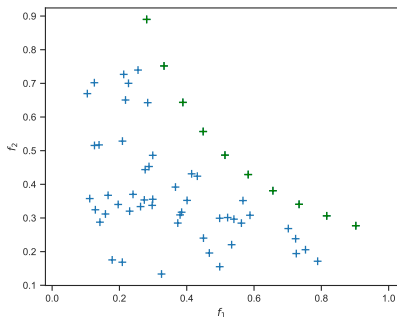


Figure 2.7: Synthetic Experiment 1: Group of arms on a convex Pareto set.

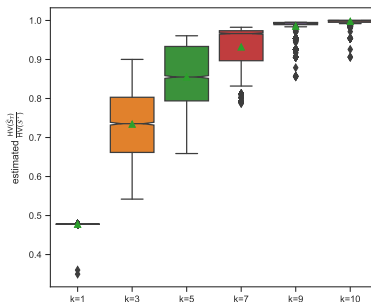


Figure 2.8: Hyper-volume fraction of the returned set on Exp.1 (convex Pareto set).

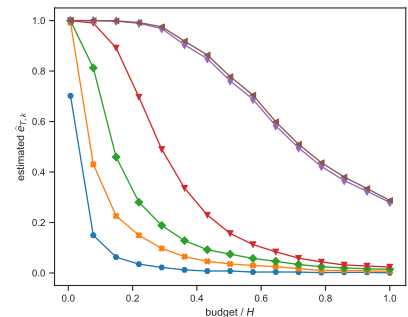


Figure 2.9: Estimated PSI- k loss for different values of k on Exp.1 (convex Pareto set).

Experiments on PSI- k . We further evaluate EGE-SR- k the convex Pareto set instance for different values of k , reporting the average \log_{10} loss over 4,000 independent trials. The PSI- k loss is defined as

$$\mathcal{L}(\widehat{S}_T, k) = \begin{cases} \mathbb{1}(\widehat{S}_T \subset \mathcal{S}^*), & \text{if } |\widehat{S}_T| = k, \\ \mathbb{1}(\widehat{S}_T = \mathcal{S}^*), & \text{otherwise.} \end{cases}$$

The empirical results (Fig. 2.9) confirm the theoretical prediction that the misidentification probability decreases exponentially fast with the sampling budget T . Smaller values of k yield smaller losses, consistent with the monotonic increase of $H_2(\nu)^{(k)}$ with k . Moreover, for $k = |\mathcal{S}^*|$, the loss can be smaller than for $k > |\mathcal{S}^*|$. This reflects the structural difference between PSI- k , which stops once k optimal arms are identified, and PSI, which continues until all arms are classified, potentially introducing additional errors in later stages.

Hyper-volume analysis. To assess the representativeness of the returned sets, we use the hyper-volume (HV) metric (Daulton et al. 2020), which measures the region dominated by a set of points. Given a reference point \mathbf{r} and a set $\mathcal{X} \subset \mathbb{R}^d$ the hyper-volume of S is

$$\text{HV}(\mathcal{X}) = \lambda \left(\bigcup_{\mathbf{x} \in S} [\mathbf{r}, \mathbf{x}] \right),$$

where λ denotes the Lebesgue measure on \mathbb{R}^d and $[\mathbf{r}, \mathbf{x}]$ is the line segment connecting \mathbf{r} and \mathbf{x} . To simplify notation we let for a discrete set $S \subset [K]$, $\text{HV}(S) := \text{HV}(\{\mu_i : i \in S\})$. Since for any set S , $\text{HV}(S)$ equals the hyper-volume of its Pareto front, we evaluate the normalized coverage

$$\alpha_{\text{HV}} = \frac{\text{HV}(\widehat{S}_T)}{\text{HV}(\mathcal{S}^*)}.$$

Higher values of α_{HV} indicate broader coverage of the Pareto-optimal region. In Fig. 2.8, even for $k = 1$, the identified arm already covers nearly half of the area dominated by the true Pareto set. This suggests that the most influential arms, those contributing most to the hyper-volume, are typically detected first by EGE-SR- k , confirming the efficiency of its allocation and discarding strategy.

Discussion. This chapter introduced the first algorithms for Pareto Set Identification in the fixed-budget setting. The proposed EGE framework, combined with sequential rejection or halving strategies, achieves an exponential decay of the misidentification probability with the budget, and its convergence exponent is provably tight for specific bandit instances. Empirical evaluations confirm that EGE-SR and EGE-SH not only outperform uniform allocation but also remain competitive with oracle algorithms requiring prior knowledge of the instance complexity $H(\nu)$.

2.5 Additional proofs

We present some proof elements for this chapter; we may refer to Kone, Kaufmann, et al. 2024 for complementary proofs.

2.5.1 Analysis of EGE

We recall for any arms i, j and round r the events

$$\begin{aligned}\xi_{i,j,r} &:= \left\{ \left\| (\hat{\mu}_{i,n_r} - \hat{\mu}_{j,n_r}) - (\mu_i - \mu_j) \right\|_\infty \leq \gamma \Delta_{(\lambda_{r+1}+1)} \right\} \\ \mathcal{E}_\gamma^1 &:= \bigcap_{r \in [R]} \bigcap_{i \in \mathcal{S}^*} \bigcap_{j \in [K]} \xi_{i,j,r}, \text{ for any } \gamma > 0.\end{aligned}$$

For each sub-optimal arm i , we define i^* to be an arbitrary element in $\operatorname{argmax}_{k \in \mathcal{S}^*} m(i, k)$, which by definition yields $\Delta_i^* = m(i, i^*)$. Introducing the property

$$\mathcal{P}_r := \{ \forall i \notin \mathcal{S}^*, i \in \mathcal{A}_r \Rightarrow i^* \in \mathcal{A}_r \},$$

the first step of the proof consists in proving that \mathcal{P}_r holds for all r .

Corollary 2.3.2. *Let $T \geq K$ and ν be a bandit with σ -subgaussian marginals. Then EGE-SR satisfies*

$$e_T^{\text{SR}}(\nu) \leq 2(K-1)^2 |\mathcal{S}^*| d \cdot e^{-\frac{T-K}{144\sigma^2 H_2(\nu) \lceil \log_2(K) \rceil}},$$

and for EGE-SH, $e_T^{\text{SH}}(\nu)$ is upper-bounded by

$$2(K-1) \lceil \log_2(K) \rceil |\mathcal{S}^*| d \cdot e^{-\frac{T}{288\sigma^2 H_2(\nu) \lceil \log_2(K) \rceil}}.$$

Proof. The case of EGE-SR has been proven in the main. For EGE-SH we have $R = \lceil \log_2(K) \rceil$, $n_r^{\text{SH}} := \sum_{s=1}^r t_s^{\text{SH}} \geq \frac{T}{\lambda_r^{\text{SH}} \lceil \log_2(K) \rceil}$ and $\lambda_r^{\text{SH}} \leq 2(\lambda_{r+1}^{\text{SH}} + 1)$. Then

$$\begin{aligned}\tilde{T}^{R, t^{\text{SH}}, \lambda^{\text{SH}}}(\nu) &:= \min_{r \in \{1, \dots, \lceil \log_2(K) \rceil\}} n_r^{\text{SH}} \Delta_{(\lambda_{r+1}^{\text{SH}}+1)}^2, \\ &\geq \min_{r \in \{1, \dots, \lceil \log_2(K) \rceil\}} \frac{\Delta_{(\lambda_{r+1}^{\text{SH}}+1)}^2}{\lambda_{r+1}^{\text{SH}} + 1} \times \frac{T}{2 \lceil \log_2(K) \rceil}, \\ &= \frac{T / (2 \lceil \log_2(K) \rceil)}{\max_{r \in \lceil \log_2(K) \rceil} (\lambda_{r+1}^{\text{SH}} + 1) \Delta_{(\lambda_{r+1}^{\text{SH}}+1)}^{-2}} \\ &\geq \frac{T}{2H_2(\nu) \lceil \log_2(K) \rceil}.\end{aligned}$$

□

Lemma 4.5.1. *On the event \mathcal{E}_γ^1 , for all $r \in [R]$ and $i, j \in \mathcal{A}_r$, if $i \in \mathcal{S}^*$ or $j \in \mathcal{S}^*$ then,*

$$\begin{aligned}|M(i, j; r) - M(i, j)| &\leq \gamma \Delta_{(\lambda_{r+1}+1)} \text{ and} \\ |m(i, j; r) - m(i, j)| &\leq \gamma \Delta_{(\lambda_{r+1}+1)}.\end{aligned}$$

Proof. Assume $i \in \mathcal{S}^*$ or $j \in \mathcal{S}^*$. We have

$$\begin{aligned} |M(i, j; r) - M(i, j)| &= \left| \max_c [\hat{\mu}_{i, n_r}^c - \hat{\mu}_{j, n_r}^c] - \max_c [\mu_i^c - \mu_j^c] \right|, \\ &\stackrel{(a)}{\leq} \max_c |(\hat{\mu}_{i, n_r}^c - \hat{\mu}_{j, n_r}^c) - (\mu_i^c - \mu_j^c)|, \\ &= \|(\hat{\mu}_{i, n_r} - \hat{\mu}_{j, n_r}) - (\mu_i - \mu_j)\|_\infty, \\ &\stackrel{(b)}{\leq} \gamma \Delta_{(\lambda_{r+1}+1)}. \end{aligned}$$

where (a) follows by reverse triangle inequality and (b) holds on the event \mathcal{E}_γ^1 . The second part of the lemma follows from

$$|m(i, j; r) - m(i, j)| = | -M(i, j; r) + M(i, j) |,$$

as $M(i, j) = -m(i, j)$ and $M(i, j; r) = -m(i, j; r)$. \square

Lemma 2.3.7. *Let $\gamma > 0$ and assume \mathcal{E}_γ^1 holds. At round $r \in [R]$, for any sub-optimal arm $i \in \mathcal{A}_r$, if $i^* \in \mathcal{A}_r$ and i^* does not empirically dominate i then $\Delta_i^* \leq \gamma \Delta_{(\lambda_{r+1}+1)}$.*

Proof. We have by definition

$$\begin{aligned} \hat{\mu}_{i, n_r} \not\leq \hat{\mu}_{i^*, n_r} &\implies \exists c : \hat{\mu}_{i, n_r}^c > \hat{\mu}_{i^*, n_r}^c \\ &\implies \exists c : (\hat{\mu}_{i, n_r}^c - \mu_i^c) - (\hat{\mu}_{i^*, n_r}^c - \mu_{i^*}^c) > \mu_{i^*}^c - \mu_i^c \geq m(i, i^*), \end{aligned}$$

so

$$\|(\hat{\mu}_{i, n_r} - \hat{\mu}_{i^*, n_r}) - (\mu_i - \mu_{i^*})\|_\infty \geq \Delta_i^*,$$

which, on the event \mathcal{E}_γ^1 is only possible if

$$\Delta_i^* \leq \gamma \Delta_{(\lambda_{r+1}+1)}.$$

\square

We recall the definition of the property

$$\mathcal{P}_r = \{ \forall i \notin \mathcal{S}^*, i \in \mathcal{A}_r \Rightarrow i^* \in \mathcal{A}_r \}.$$

and start by proving a first intermediate result towards proving Lemma 2.3.8.

Lemma 2.5.1. *Assume that \mathcal{E}_1^c holds and let $r \in [R]$ such that \mathcal{P}_r holds. Then for any $i \in \mathcal{A}_r$,*

$$(\hat{\Delta}_{i, r}^*)_+ - (\Delta_i^*)_+ \leq 2\gamma \Delta_{(\lambda_{r+1}+1)} \quad \text{and} \quad (\hat{\Delta}_{i, r}^*)_+ - (\Delta_i^*)_+ \geq -\gamma \Delta_{(\lambda_{r+1}+1)}.$$

Proof. We define the max of an empty set to be $-\infty$. We first analyze the case $i \in \mathcal{S}^*$. When $i \in \mathcal{S}^*$, we have

$$(\Delta_i^*)_+ = 0 = \left(\max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j) \right)_+,$$

which yields

$$\begin{aligned}
 |(\widehat{\Delta}_{i,r}^*)_+ - (\Delta_i^*)_+| &= \left| \left(\max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j; r) \right)_+ - \left(\max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j) \right)_+ \right|, \\
 &\leq \left| \left(\max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j; r) \right) - \left(\max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j) \right) \right| \quad (\text{since } |x_+ - y_+| \leq |x - y|), \\
 &\leq \max_{j \in \mathcal{A}_r \setminus \{i\}} |m(i, j; r) - m(i, j)|, \\
 &\leq \gamma \Delta_{(\lambda_{r+1}+1)},
 \end{aligned}$$

where the last inequality follows from Lemma 2.3.6. We now assume that i is a sub-optimal arm. Since \mathcal{P}_r holds at round r , $i^* \in \mathcal{A}_r$ and

$$\Delta_i^* = \max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j).$$

Let $\hat{i} \in \operatorname{argmax}_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j; r)$ then

$$(\widehat{\Delta}_{i,r}^*)_+ - (\Delta_i^*)_+ = (m(i, \hat{i}; r))_+ - (m(i, i^*))_+, \quad (2.20)$$

$$\geq (m(i, i^*; r))_+ - (m(i, i^*))_+ \quad (\text{since } i^* \in \mathcal{A}_r). \quad (2.21)$$

We further note that

$$\begin{aligned}
 |(m(i, i^*; r))_+ - (m(i, i^*))_+| &\leq |m(i, i^*; r) - m(i, i^*)|, \\
 &\leq \gamma \Delta_{(\lambda_{r+1}+1)},
 \end{aligned}$$

which follows from Lemma 2.3.6. Combining the latter inequality with (2.21) yields

$$(\widehat{\Delta}_{i,r}^*)_+ - (\Delta_i^*)_+ \geq -\gamma \Delta_{(\lambda_{r+1}+1)}. \quad (2.22)$$

We also have

$$\begin{aligned}
 (\widehat{\Delta}_{i,r}^*)_+ - (\Delta_i^*)_+ &\leq (m(i, \hat{i}; r))_+ - (m(i, \hat{i}))_+, \\
 &\leq |m(i, \hat{i}; r) - m(i, \hat{i})|, \\
 &= |M(i, \hat{i}) - M(i, \hat{i}; r)|, \\
 &\leq \|(\mu_i - \mu_{\hat{i}}) - (\hat{\mu}_{i, n_r} - \hat{\mu}_{\hat{i}, n_r})\|_\infty, \\
 &= \|((\mu_i - \mu_{i^*}) - (\hat{\mu}_{i, n_r} - \hat{\mu}_{i^*, n_r})) + ((\mu_{i^*} - \mu_{\hat{i}}) - (\hat{\mu}_{i^*, n_r} - \hat{\mu}_{\hat{i}, n_r}))\|_\infty, \\
 &\leq 2\gamma \Delta_{(\lambda_{r+1}+1)},
 \end{aligned}$$

where the last inequality follows from the triangle inequality and Lemma 2.3.6. This proves the lemma. \square

We then prove Lemma 2.5.2, which can be viewed as a ‘‘symmetric’’ version of Lemma 2.3.8.

Lemma 2.5.2. *Assume that \mathcal{E}_γ^1 holds. Let $r \in [R]$ such that \mathcal{P}_r holds. Then, for any sub-optimal arm $i \in \mathcal{A}_r$,*

$$\left| \widehat{\Delta}_{i,r}^* - \Delta_i^* \right| \leq 2\gamma\Delta_{(\lambda_{r+1}+1)}$$

and for any optimal arm $i \in \mathcal{A}_r$,

$$\widehat{\delta}_{i,r}^* \geq \Delta_i - 2\gamma\Delta_{(\lambda_{r+1}+1)}.$$

Proof. Let $i \in \mathcal{A}_r \cap (\mathcal{S}^*)^c$. By assumption, $i^* \in \mathcal{A}_r$, so

$$\Delta_i := \Delta_i^* = \max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j),$$

then we have

$$\begin{aligned} |\widehat{\Delta}_{i,r}^* - \Delta_i| &= \left| \left(\max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j; r) \right) - \left(\max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j) \right) \right|, \\ &\stackrel{(a)}{\leq} \max_{j \in \mathcal{A}_r \setminus \{i\}} |m(i, j; r) - m(i, j)|, \\ &\leq \max_{j \in \mathcal{A}_r \setminus \{i\}} |M(i, j) - M(i, j; r)|, \\ &\leq \max_{j \in \mathcal{A}_r \setminus \{i\}} \|(\mu_i - \mu_j) - (\hat{\mu}_{i,n_r} - \hat{\mu}_{j,n_r})\|_\infty \\ &= \max_{j \in \mathcal{A}_r \setminus \{i\}} \|((\mu_i - \mu_{i^*}) - (\hat{\mu}_{i,n_r} - \hat{\mu}_{i^*,n_r})) + ((\mu_{i^*} - \mu_j) - (\hat{\mu}_{i^*,n_r} - \hat{\mu}_{j,n_r}))\|_\infty, \\ &\stackrel{(b)}{\leq} 2\gamma\Delta_{(\lambda_{r+1}+1)}, \end{aligned}$$

where (a) follows by reverse triangle inequality and (b) follows by triangle inequality and Lemma 2.3.6. So we have proved the first statement of the lemma. Before proving the second statement, recall that

$$\widehat{\delta}_{i,r}^* = \min_{j \in \mathcal{A}_r \setminus \{i\}} [M(i, j; r) \wedge (M(j, i; r)_+ + (\widehat{\Delta}_{j,r}^*)_+)].$$

Let $i \in \mathcal{A}_r \cap \mathcal{S}^*$. For any $j \in \mathcal{A}_r$ we have

$$\begin{aligned} |M(j, i; r)_+ - M(j, i)_+| &\leq |M(j, i; r) - M(j, i)|, \\ &\leq \gamma\Delta_{(\lambda_{r+1}+1)} \quad (\text{Lemma 2.3.6}), \end{aligned}$$

and similarly

$$|M(i, j; r) - M(i, j)| \leq \gamma\Delta_{(\lambda_{r+1}+1)}$$

which, combined, yields

$$M(j, i; r)_+ \geq M(j, i)_+ - \gamma\Delta_{(\lambda_{r+1}+1)} \quad \text{and} \quad (2.23)$$

$$M(i, j; r) \geq M(i, j) - \gamma\Delta_{(\lambda_{r+1}+1)}. \quad (2.24)$$

Letting $j \in \mathcal{A}_r$, we have by Lemma 2.5.1

$$(\widehat{\Delta}_{j,r}^*)_+ - (\Delta_j^*)_+ \geq -\gamma\Delta_{(\lambda_{r+1}+1)} \quad (2.25)$$

Combined, (2.25) and (2.23) yields

$$M(j, i; r)_+ + (\widehat{\Delta}_{j,r}^*)_+ \geq M(j, i)_+ + (\Delta_j^*)_+ - 2\gamma\Delta_{(\lambda_{r+1}+1)}$$

for any $j \in \mathcal{A}_r$. Which combined with (2.24) yields

$$[M(i, j; r) \wedge (M(j, i; r)_+ + (\widehat{\Delta}_{j,r}^*)_+)] \geq [M(i, j) \wedge (M(j, i)_+ + (\Delta_j^*)_+)] - 2\gamma\Delta_{(\lambda_{r+1}+1)}$$

for any $j \in \mathcal{A}_r$. Therefore,

$$\begin{aligned} \widehat{\delta}_{i,r}^* &\geq \left(\min_{j \in \mathcal{A}_r \setminus \{i\}} [M(i, j) \wedge (M(j, i)_+ + (\Delta_j^*)_+)] \right) - 2\gamma\Delta_{(\lambda_{r+1}+1)}, \\ &\geq \left(\min_{j \in [K] \setminus \{i\}} [M(i, j) \wedge (M(j, i)_+ + (\Delta_j^*)_+)] \right) - 2\gamma\Delta_{(\lambda_{r+1}+1)}, \\ &= \delta_i^* - 2\gamma\Delta_{(\lambda_{r+1}+1)}. \end{aligned}$$

□

Finally, we prove Lemma 2.3.8.

Lemma 2.3.8. *Assume that \mathcal{E}_γ^1 holds and let $r \in [R]$ such that \mathcal{P}_r holds. Then for any $i \in \mathcal{A}_r$*

$$\widehat{\Delta}_{i,r} - \Delta_i \geq \begin{cases} -2\gamma\Delta_{(\lambda_{r+1}+1)} & \text{if } i \in \mathcal{S}^*, \\ -\gamma\Delta_{(\lambda_{r+1}+1)} & \text{else.} \end{cases}$$

Proof of Lemma 2.3.8. Let $i \in \mathcal{A}_r \cap (\mathcal{S}^*)^c$. Since $i \in (\mathcal{S}^*)^c$, $\Delta_i = \Delta_i^*$. If $i \notin S_r$ then $\widehat{\Delta}_{i,r} = \widehat{\Delta}_{i,r}^*$. If $i \in S_r$ then $\widehat{\Delta}_{i,r} = \widehat{\delta}_{i,r}^*$ and

$$\widehat{\Delta}_{i,r}^* \leq 0 \leq \widehat{\delta}_{i,r}^*$$

which follows since $i \in S_r$. Therefore in both cases

$$\begin{aligned} \widehat{\Delta}_{i,r} - \Delta_i &\geq \widehat{\Delta}_{i,r}^* - \Delta_i^* \\ &= \max_{j \in \mathcal{A}_r} m(i, j; r) - m(i, i_*) \\ &\geq m(i, i_*; r) - m(i, i_*) \\ &\geq -\gamma\Delta_{(\lambda_{r+1}+1)}, \end{aligned}$$

where the last inequality uses Lemma 2.3.6.

Similarly, let $i \in \mathcal{A}_r \cap \mathcal{S}^*$. We have $\Delta_i = \delta_i^*$ (Lemma 2.1.2) and $\widehat{\Delta}_{i,r} = \widehat{\delta}_{i,r}^*$ if $i \in S_r$. If $i \in \mathcal{A}_r \setminus S_r$ (empirically sub-optimal) then

$$\widehat{\Delta}_{i,r} = \widehat{\Delta}_{i,r}^* \geq 0$$

and there exists $j \in \mathcal{A}_r$ such that $\widehat{\mu}_{i,n_r} \prec \widehat{\mu}_{j,n_r}$ so for this arm j , $M(i, j; r) < 0$. Therefore $\widehat{\delta}_{i,r}^* < 0$ and

$$\widehat{\delta}_{i,r}^* < 0 \leq \widehat{\Delta}_{i,r} = \widehat{\Delta}_{i,r}^*.$$

In all cases, we have

$$\widehat{\Delta}_{i,r} - \Delta_i \geq \widehat{\delta}_{i,r}^* - \Delta_i \geq -2\gamma\Delta_{(\lambda_{r+1}+1)},$$

where the last inequality follows from Lemma 2.5.2. □

2.5.2 Analysis of EGE-SR- k

In this section, we analyze EGE-SR- k and bound its stopping time. We define for any $\gamma > 0$ and $t \in [K - 1]$:

$$\mathcal{E}_\gamma^2(t) := \bigcap_{r \in [t]} \bigcap_{i \in \mathcal{S}^*} \bigcap_{j \in [K]} \left\{ \left\| (\hat{\mu}_{i,n_r} - \hat{\mu}_{j,n_r}) - (\mu_i - \mu_j) \right\|_\infty \leq \gamma \Delta_{(K+1-r)}^{(k)} \right\}$$

and in particular

$$\mathcal{E}_\gamma^2 := \mathcal{E}_\gamma^2(K - 1).$$

We define also $\tilde{k} = \min \{k, |\mathcal{S}^*|\}$.

For any round r , let $\alpha(r) = |\mathcal{B}_r|$ denote the number of arms so far identified as optimal at the beginning of round r . We denote these arms by $a_1, \dots, a_{\alpha(r)}$. We say that the algorithm has made an error before round r if $\mathcal{B}_r \cap (\mathcal{S}^*)^c \neq \emptyset$ or $\mathcal{D}_r \cap \mathcal{S}^* \neq \emptyset$. We will show that on $\mathcal{E}_\gamma^2(t)$, the algorithm does not make any error until the end of round τ_k^t such that $\alpha(\tau_k^t + 1) = \tilde{k}$.

Lemma 2.5.3. *Let $\gamma < 1/6$, $t \in [R]$ and $k \leq |\mathcal{S}^*|$. Let $\tau_k^t := \min \{r \in [t] : \alpha(r + 1) = \tilde{k}\} \wedge t$. On the event $\mathcal{E}_\gamma^2(t)$, the algorithm makes no error until the end round τ_k^t and for any $r \leq \tau_k^t + 1$, if $j \in \mathcal{A}_r$ is a sub-optimal arm then $j^* \in \mathcal{A}_r$. In particular, on the event $\mathcal{E}_\gamma^2(t)$, $\{a_1, \dots, a_{\alpha(\tau_k^t + 1)}\} \subset \mathcal{S}^*$.*

Said otherwise, this lemma states that the first arms that will be declared as optimal by EGE-SR- k will be actually optimal if $\mathcal{E}_\gamma^2(t)$ holds for γ small enough.

Proof. In the sequel, we assume that $\mathcal{E}_\gamma^2(t)$ holds. We prove the correctness by induction on the round r . Let $r \in [t]$ and let $\mathcal{B}_r = \{a_1, \dots, a_{\alpha(r)}\}$ denote the arms so far identified as optimal. Let \mathcal{P}_r be the property “for any sub-optimal arm $i \in \mathcal{A}_r$, $i^* \in \mathcal{A}_r$ and no error occurred so far”.

\mathcal{P}_r trivially holds for $r = 1$ and $\alpha(1) = 0$. We now assume that it holds until the beginning of round r and that $\alpha(r) < \tilde{k}$. We will show that the arm i_r de-activated at the end of round r is well classified and that for any sub-optimal arm $j \in \mathcal{A}_{r+1}$, $j^* \in \mathcal{A}_{r+1}$.

Lemma 2.5.4. *If \mathcal{P}_r holds at round r and $\alpha(r) < \tilde{k}$, there exists $a \in \mathcal{A}_r$ such that $\Delta_a^{(k)} = \Delta_a$ and $\Delta_a^{(k)} \geq \Delta_{(K+1-r)}^{(k)}$.*

Proof. At the beginning of round r , it remains $K - r + 1$ active arms so there exists $a \in \mathcal{A}_r$ such that $\Delta_a^{(k)} \geq \Delta_{(K+1-r)}^{(k)}$. If a is sub-optimal then $\Delta_a^{(k)} = \Delta_a$. Otherwise, if a is an optimal arm, since $\alpha(r) < \tilde{k}$ and no error has occurred so far (by assumption) then there exists one of the optimal arms $a' \in \mathcal{A}_r$ (one of those with the \tilde{k} largest gaps) such that $\Delta_{a'}^{(k)} = \Delta_{a'}$ and $\Delta_{(K+1-r)}^{(k)} \leq \Delta_a^{(k)} \leq \Delta_{a'}^{(k)}$.

□

Lemma 2.3.6 still holds for the event $\mathcal{E}_\gamma^2(t)$ with the modified gaps introduced earlier. We state the following lemma, which is similar to Lemma 2.5.2.

Lemma 2.5.5. *Assume that the event $\mathcal{E}_\gamma^2(t)$ holds. Let $r \in [t]$ and assume that for any sub-optimal arm $j \in \mathcal{A}_r$, $j^* \in \mathcal{A}_r$. Then, for any sub-optimal arm $i \in \mathcal{A}_r$,*

$$|\widehat{\Delta}_{i,r}^* - \Delta_i^*| \leq 2\gamma\Delta_{(K+1-r)}^{(k)}$$

and for any optimal arm $i \in \mathcal{A}_r$,

$$\widehat{\delta}_{i,r}^* \geq \Delta_i - 2\gamma\Delta_{(K+1-r)}^{(k)}.$$

As already noted in the proof of Lemma 2.3.9 and Lemma 2.3.8, it is simple to see that the empirical gap $\widehat{\Delta}_{i,r}$ of each arm i satisfies

$$\widehat{\Delta}_{i,r} = \max \left\{ \widehat{\Delta}_{i,r}^*, \widehat{\delta}_{i,r}^* \right\}. \quad (2.26)$$

It follows from the fact that when $i \in S_r$, $\widehat{\Delta}_{i,r}^* < 0 \leq \widehat{\delta}_{i,r}^*$ and when $i \in \mathcal{A}_r \setminus S_r$, $\widehat{\delta}_{i,r}^* < 0 \leq \widehat{\Delta}_{i,r}^*$.

Using Lemma 2.5.5 and the inductive hypothesis \mathcal{P}_r enables us to prove that

$$\forall i \in \mathcal{A}_r, \widehat{\Delta}_{i,r} \geq \Delta_i - 2\gamma\Delta_{(K+1-r)}^{(k)} \quad (2.27)$$

We split the proof in three steps to prove that \mathcal{P}_{r+1} holds.

Step 1: i_r is well classified. Let i_r be empirically sub-optimal ($i_r \notin S_r$) and assume $i_r \in \mathcal{S}^*$.

Since i_r is removed at the end of round r , by the inductive assumption and using Equation 2.27 and Lemma 2.3.6 (adapted) we have

$$\max_{j \in \mathcal{A}_r \setminus \{i_r\}} m(i_r, j) \geq \Delta_i - 3\gamma\Delta_{(K+1-r)}^{(k)}, \quad \forall i \in \mathcal{A}_r. \quad (2.28)$$

By Lemma 2.5.4, there exists $a \in \mathcal{A}_r$ such that $\Delta_a = \Delta_a^{(k)} \geq \Delta_{(K+1-r)}^{(k)}$. Applying (2.28) to this arm a yields

$$\max_{j \in \mathcal{A}_r \setminus \{i_r\}} m(i_r, j) \geq (1 - 3\gamma)\Delta_{(K+1-r)}^{(k)}. \quad (2.29)$$

Recalling that $\gamma < 1/6$ we see that the LHS of (2.29) is positive, so there exists j such that

$$\mu_{i_r} \prec \mu_j$$

which contradicts the assumption $i_r \in \mathcal{S}^*$.

Now, if i_r is empirically optimal, i.e., $i_r \in S_r$ assume i_r is a sub-optimal arm that is $i_r \notin \mathcal{S}^*$. By the hypothesis \mathcal{P}_r , $i_r^* \in \mathcal{A}_r$ and by definition of $\widehat{\Delta}_{i_r,r}$ and since i_r is removed, we have (Lemma 2.3.6 and Equation 2.27)

$$M(i_r, i_r^*) \geq \Delta_i - 3\gamma\Delta_{(K+1-r)}^{(k)} \quad \forall i \in \mathcal{A}_r. \quad (2.30)$$

Using Lemma 2.5.4 as before, there exists $a \in \mathcal{A}_r$ such that $\Delta_{(K+1-r)}^{(k)} \leq \Delta_a^{(k)} = \Delta_a$. So the LHS of (2.30) is positive for $\gamma < 1/6$. That is

$$M(i_r, i_r^*) > 0,$$

which is impossible as $\mu_{i_r} \prec \mu_{i_r^*}$. This concludes the proof of this part: i_r is well classified.

Step 2: If $i_r \in \mathcal{S}^*$ is an optimal arm, then no active arm is “optimally” dominated by i_r :

$$\forall i \in \mathcal{A}_r \cap (\mathcal{S}^*)^c, i^* \neq i_r$$

To prove this, note that we can deduce from “Step 1” that $i_r \in S_r$. By contradiction, let $i \in \mathcal{A}_r \cap (\mathcal{S}^*)^c$ such that $i^* = i_r$. We claim that i is not dominated by i_r . Indeed, if i were dominated by i_r , we would have $i \notin S_r$ (i.e., empirically sub-optimal) so

$$\widehat{\Delta}_{i,r} = \widehat{\Delta}_{i,r}^* > 0,$$

and since $i_r \in S_r$, by definition of $\widehat{\delta}_{i_r,r}^*$ and noting that $\widehat{\Delta}_{i_r,r} \geq \widehat{\Delta}_{i,r}$, we would have

$$\begin{aligned} \widehat{\Delta}_{i_r,r} &\leq (M(i, i_r; r))_+ + (\widehat{\Delta}_i^*) \\ &= 0 + \widehat{\Delta}_{i,r}, \end{aligned}$$

that is $\widehat{\Delta}_{i_r,r} \leq \widehat{\Delta}_{i,r}$. However, our tie-breaking ensures that this inequality is impossible: since $i_r \in S_r, i \notin S_r$ and i_r is deactivated, we have $\widehat{\Delta}_{i_r,r} > \widehat{\Delta}_{i,r}$. Therefore, i is not empirically dominated by i_r . Moreover, since $i_r = i^*$ (by assumption), Lemma 2.3.7 (it trivially holds with the modified gaps) yields

$$\Delta_i^* \leq \gamma \Delta_{(K+1-r)}^{(k)}. \quad (2.31)$$

Recalling that i_r is deactivated, we have for any arm $j \in \mathcal{A}_r$,

$$(M(i, i_r; r))_+ + (\widehat{\Delta}_{i_r,r}^*)_+ \geq \widehat{\Delta}_{j,r} \quad (2.32)$$

On the other side, there exists $a \in \mathcal{A}_r$ such that $\Delta_{(K+1-r)}^{(k)} \leq \Delta_a^{(k)} = \Delta_a$ (Lemma 2.5.4). Applying Equation (2.27) and Lemma 2.3.6 and taking $j = a$ in Equation 2.32 yields

$$(M(i, i_r))_+ + (\Delta_i^*)_+ \geq (1 - 5\gamma) \Delta_{(K+1-r)}^{(k)}$$

that is

$$\Delta_i^* \geq (1 - 5\gamma) \Delta_{(K+1-r)}^{(k)}. \quad (2.33)$$

Both (2.31) and (2.33) cannot hold when $\gamma < 1/6$. So for any $i \in \mathcal{A}_r$

$$i^* \neq i_r.$$

Which completes the proof of “Step 2”. Combining, “Step 1&2” proves that \mathcal{P}_{r+1} holds.

Step 3: Conclusion

We have proved that if \mathcal{P}_r holds and $\alpha(r) < \tilde{k}$ then \mathcal{P}_{r+1} holds. Note that if $i_r \in S_r$ then $\alpha(r+1) = \alpha(r) + 1$ otherwise $\alpha(r+1) = \alpha(r)$. Therefore, no error occurs until the end of round $\min\{r', t\}$ where r' is such that $\alpha(r'+1) = \tilde{k}$. As a consequence $\mathcal{B}_{\tau_k^t+1} = \{a_1, \dots, a_{\alpha(\tau_k^t+1)}\} \subset \mathcal{S}^*$. And by the proof of "Step 2", if $j \in \mathcal{A}_{\tau_p^t}$ is a sub-optimal arm, then $j^* \in \mathcal{A}_{\tau_p^t+1}$.

□

Theorem 2.3.3. *Let $k \in [K]$. EGE-SR- k satisfies*

$$e_{T,k}(\nu) \leq 2(K-1)^2 |\mathcal{S}^*| d \cdot e^{-\frac{T-K}{144\sigma^2 H_2^{(k)}(\nu) \overline{\log(K)}}}.$$

Proof. The proof is a direct consequence of Lemma 2.5.3 and Hoeffding's inequality. Indeed, Lemma 2.5.3 yields that EGE-SR- k is correct on the event \mathcal{E}_γ^2 for $\gamma < 1/6$. Therefore, for any $0 < \gamma < 1/6$

$$e_{T,k}(\nu) \leq \mathbb{P}((\mathcal{E}_\gamma^2)^c).$$

Letting $0 < \gamma < 1/6$ fixed, by union bound and Hoeffding's inequality, it follows that

$$\mathbb{P}((\mathcal{E}_\gamma^2)^c) \leq 2(K-1)^2 |\mathcal{S}^*| d \exp\left(-\gamma^2 \frac{T-K}{4\sigma^2 H_2^{(k)}(\nu) \overline{\log(K)}}\right)$$

therefore,

$$e_{T,k}(\nu) \leq 2(K-1)^2 |\mathcal{S}^*| d \exp\left(-\gamma^2 \frac{T-K}{4\sigma^2 H_2^{(k)}(\nu) \overline{\log(K)}}\right)$$

and taking the limit to $1/6$ (as it holds for any $0 < \gamma < 1/6$) yields the expected bound. □

The theorem below bounds the expected stopping time and the number of samples used at stopping.

Theorem 2.5.6. *Fix $k < |\mathcal{S}^*|$ and let $q := K - |\mathcal{S}^*| + k$. Then*

$$\begin{aligned} \mathbb{E}[\tau] &\leq q + 2(K-1) |\mathcal{S}^*| (K-q-1) q d \exp\left(-\frac{T-K}{144\sigma^2 H_2^{(k)}(\nu) \overline{\log(K)}}\right) \text{ and} \\ \mathbb{E}[N_\tau] &\leq N_q + 2(K-1) |\mathcal{S}^*| (K-q-1) q d T \exp\left(-\frac{T-K}{144\sigma^2 H_2^{(k)}(\nu) \overline{\log(K)}}\right). \end{aligned}$$

This result suggests that for this relaxed problem, the algorithm might not need to use the whole budget, in particular when T is large. For example, consider a setting $[K] = \mathcal{S}^*$ then $q = k$ and we roughly use N_k samples which can be way smaller than N_{K-1} .

Proof. The idea is to show that if the algorithm has not stopped after round q then some high probability event must not hold. Let $\gamma < 1/6$ be fixed. We have

$$\begin{aligned}\mathbb{E}[\tau] &\leq q + \mathbb{E}[\tau \mathbb{1}(\tau > q)], \\ &\leq q + \sum_{s=q+1}^{K-1} \mathbb{P}(\tau \geq s).\end{aligned}$$

We claim that for any $s > q$,

$$\{\tau \geq s\} \subset (\mathcal{E}_\gamma^2(q))^c.$$

Indeed,

$$\begin{aligned}\tau \geq s &\implies \tau > q \\ &\implies \alpha(q+1) < k,\end{aligned}$$

However, by Lemma 2.5.3, if $\mathcal{E}_\gamma^2(q)$ holds and $\alpha(q+1) < k$ then no error has occurred until the end of round q , therefore $\mathcal{D}_q = (\mathcal{S}^*)^c$ and $|\mathcal{B}_q| = k$, which is not possible as $\alpha(q+1) < k$. So $\tau \geq s \implies (\mathcal{E}_\gamma^2(q))^c$ does not hold). Then

$$\mathbb{E}[\tau] \leq q + \sum_{s=q+1}^{K-1} \mathbb{P}((\mathcal{E}_\gamma^2(q))^c), \quad (2.34)$$

$$\leq q + \sum_{s=q+1}^{K-1} 2(K-1)|\mathcal{S}^*|qd \exp\left(-\gamma^2 \frac{T-K}{4\sigma^2 H_2^{(k)}(\nu) \overline{\log}(K)}\right), \quad (2.35)$$

$$\leq q + 2(K-1)|\mathcal{S}^*|(K-q-1)qd \exp\left(-\gamma^2 \frac{T-K}{4\sigma^2 H_2^{(k)}(\nu) \overline{\log}(K)}\right). \quad (2.36)$$

Similarly, we have

$$\mathbb{E}[N_\tau] \leq N_q + \mathbb{E}[(N_\tau \mathbb{1}(\tau > q))].$$

Next,

$$\begin{aligned}\mathbb{E}[N_\tau] &\leq N_q + \sum_{s=q+1}^{K-1} N_s \mathbb{P}(\tau \geq s) \\ &\leq N_q + \sum_{s=q+1}^{K-1} N_s \mathbb{P}((\mathcal{E}_\gamma^2(q))^c) \\ &\leq N_q + 2(K-1)|\mathcal{S}^*|qd \exp\left(-\gamma^2 \frac{T-K}{4\sigma^2 H_2^{(k)}(\nu) \overline{\log}(K)}\right) \sum_{s=q+1}^{K-1} N_s.\end{aligned}$$

Simple algebra yields for any $r \in [K - 1]$,

$$\begin{aligned}
 N_r &= (K - r)n_r + \sum_{s=1}^r n_s, \\
 &\leq \frac{T - K}{\overline{\log(K)}} + (K - r) + r + \frac{T - K}{\overline{\log(K)}} \sum_{s=1}^r \frac{1}{K + 1 - s} \\
 &= \frac{T - K}{\overline{\log(K)}} \left(1 + \sum_{s=K+1-r}^K s^{-1} \right) + K \\
 &\leq T.
 \end{aligned}$$

Therefore,

$$\mathbb{E}[N_\tau] \leq N_q + 2(K - 1)|\mathcal{S}^*|(K - q - 1)qdT \exp \left(-\gamma^2 \frac{T - K}{4\sigma^2 H_2^{(k)}(\nu) \overline{\log(K)}} \right) \quad (2.37)$$

As (2.36) and (2.37) hold for any $\gamma < 1/6$, taking the limit for a sequence $\gamma \rightarrow 1/6$ yields the expected constants in the exponent, which completes the proof. \square

2.5.3 Simplifying the sub-optimality gaps

In this section, we prove the lemma that allows us to remove the explicit dependency on \mathcal{S}^* in the expression of the sub-optimality gaps. Towards proving Lemma 2.1.2, we prove the result below.

Lemma 2.5.7. *For any sub-optimal arm a , there exists a Pareto optimal arm a^* such that $\mu_a \prec \mu_{a^*}$ and $\Delta_a = m(a, a^*) > 0$. Moreover, for any $i \in [K] \setminus \mathcal{S}^*$, $j \in \mathcal{S}^*$*

- i) $\max_{k \in \mathcal{S}^*} m(i, k) = \max_{k \in [K]} m(i, k)$,
- ii) *If $i \in \operatorname{argmin}_{k \in [K] \setminus \{j\}} M(j, k)$ then j is the unique arm such that $\mu_i \prec \mu_j$.*

Proof. Assume there are $p < n$ dominated arms. Without loss of generality, we may assume they are μ_1, \dots, μ_p . Let $i_1 \leq p$. Suppose that no Pareto-optimal arm dominates μ_{i_1} . Since μ_{i_1} is not optimal, by the latter assumption, there exists $i_2 \leq p$ such that $\mu_{i_1} \prec \mu_{i_2}$. If μ_{i_2} is dominated by a Pareto optimal arm, this arm also dominates μ_{i_1} (strict dominance is transitive) which contradicts the initial assumption. If not, there exists $i_3 \leq p$ such that $\mu_{i_1} \prec \mu_{i_2} \prec \mu_{i_3}$. Again we can use the same reasoning as before for i_3 . In any case we should stop in at most p steps, otherwise we would have $\mu_{i_1} \prec \mu_{i_2} \prec \dots \prec \mu_{i_p}$ and μ_{i_p} should be dominated by a Pareto-optimal arm, otherwise it would be itself Pareto-optimal, which is not the case. Therefore, for any $a \in [K] \setminus \mathcal{S}^*$, there exists $a^* \in \mathcal{S}^*$ such that $a^* \prec a$ and $\Delta_a = m(a, a^*) > 0$.

Letting i be a sub-optimal arm, since for any $a \in [K] \setminus \mathcal{S}^*$, there exists $a^* \in \mathcal{S}^*$ such that $a \prec a^*$, it follows that

$$\forall c \in [d], \mu_a^c - \mu_i^c < \mu_{a^*}^c - \mu_i^c,$$

which leads to $m(i, a) \leq m(i, a^*)$, so

$$\max_{j \in [K]} m(i, j) = \max_{j \in \mathcal{S}^*} m(i, j) > 0,$$

which completes the proof of the first point i). For the second point, let $q \in [K] \setminus \mathcal{S}^*$ and q' such that $q \prec q'$ and

$$q \in \operatorname{argmin}_{a \in [K] \setminus \{j\}} M(j, a).$$

By direct algebra, since $q \prec q'$, we have

$$M(j, q') < M(j, q),$$

which is impossible if $q' \neq j$ (because q belongs to the argmin). Therefore, if

$$q \in \operatorname{argmin}_{a \in [K] \setminus \{j\}} M(j, a)$$

is a sub-optimal arm, then j is the only arm such that $q \prec j$ (i.e., $\mu_q \prec \mu_j$). □

We now prove the following result, which follows from the lemma above.

Lemma 2.5.8. *For any Pareto optimal arm i , $\Delta_i \leq \min_{j \neq i} M(i, j)$.*

Proof of Lemma 2.5.8. If $\operatorname{argmin}_{j \neq i} M(i, j) \subset \mathcal{S}^*$, then the lemma follows from the definition of the gap of an optimal arm recalled in Section 4. If $\min_{j \neq i} M(i, j) = M(i, a)$, $a \notin \mathcal{S}^*$, then, from Lemma 2.5.7, i is the unique arm which dominates a so $\Delta_a = m(a, i)$ and using the definition of the gap of an optimal arm,

$$\begin{aligned} \Delta_i &\leq M(a, i)_+ + \Delta_a, \\ &= 0 + m(a, i) \leq M(i, a), \end{aligned}$$

where we have used the fact that $m(p, q) \leq M(q, p)$ for any pair of arms p, q (which follows from the definition). Therefore, for an optimal arm i , we always have

$$\Delta_i \leq \min_{j \neq i} M(i, j).$$

□

We can now finish by proving Lemma 2.1.2.

Proof of Lemma 2.1.2. For sub-optimal arms, the result follows from Lemma 2.5.7. It remains to prove the equality for sub-optimal arms. By definition, for an optimal arm i , we have

$$\Delta_i = \min \{ \delta_i^+, \delta_i^- \},$$

where

$$\delta_i^+ := \min_{j \in \mathcal{S}^* \setminus \{i\}} \min \{ M(i, j), M(j, i) \} \text{ and } \delta_i^- := \min_{j \in [K] \setminus \mathcal{S}^*} [(M(j, i))_+ + \Delta_j^*].$$

For any optimal arm i , $\Delta_i^* \leq 0$ (by direct calculation), so introducing

$$\delta_i^{-'} := \min \left\{ \delta_i^-, \min_{j \in \mathcal{S}^* \setminus \{i\}} M(j, i) \right\},$$

we have

$$\delta_i^{-'} = \min_{j \neq i} [M(j, i)_+ + (\Delta_j^*)_+]$$

Then, if

$$\min_{j \in \mathcal{S}^* \setminus \{i\}} M(i, j) = \min_{j \neq i} M(i, j), \quad (2.38)$$

holds, the result simply follows as for any optimal arm i ,

$$\begin{aligned} \Delta_i = \min \{ \delta_i^+, \delta_i^- \} &= \min \left\{ \min_{j \in \mathcal{S}^* \setminus \{i\}} M(i, j), \delta_i^{-'} \right\}, \\ &= \min(\min_{j \neq i} M(i, j), \delta_i^{-'}), \\ &= \min_{j \neq i} [M(i, j) \wedge (M(j, i)_+ + (\Delta_j^*)_+)], \\ &= \delta_i^*. \end{aligned}$$

In the sequel, assume (2.38) does not hold, that is assume

$$\min_{j \in \mathcal{S}^* \setminus \{i\}} M(i, j) > \min_{j \neq i} M(i, j).$$

From Lemma 2.5.7 we know that in this case, there exists a sub-optimal arm k such that i is the unique arm dominating k and

$$\Delta_k^* = m(k, i) \quad \text{and} \quad \min_{j \neq i} M(i, j) = M(i, k).$$

Thus, we have

$$\min_{j \neq i} M(i, j) = M(i, k), \quad (2.39)$$

$$\geq m(k, i) = \Delta_k^*, \quad (2.40)$$

$$\geq \delta_i^{-'} \quad (\text{since } i \text{ dominates } k, M(k, i)_+ = 0). \quad (2.41)$$

On the other side, as $\mu_k \prec \mu_i$, using the definition of Δ_i in particular δ_i^- directly yields

$$\Delta_i \leq \Delta_k^* = m(k, i) \quad (2.42)$$

$$\leq M(i, k) = \min_{j \neq i} M(i, j) \quad (2.43)$$

$$< \min_{j \in \mathcal{S}^* \setminus \{i\}} M(i, j). \quad (2.44)$$

We recall that

$$\Delta_i = \min \left\{ \delta_i^{-'}, \min_{j \in \mathcal{S}^* \setminus \{i\}} M(i, j) \right\},$$

which combined with (2.44) yields

$$\Delta_i = \delta_i^{-'}$$

and further combining with (2.41) yields

$$\begin{aligned} \Delta_i = \delta_i^{-'} &= \min \left\{ \min_{j \neq i} M(i, j), \delta_i^{-'} \right\} \\ &= \delta_i^*, \end{aligned}$$

which concludes the proof. □

Chapter 3

Adaptive Algorithms for Fixed-Confidence Pareto Set Identification

In this chapter, we study the problem of fixed-confidence identification of the Pareto set in multi-objective multi-armed bandits. Since the sample complexity required to exactly identify the Pareto set can be prohibitive, previous work has proposed relaxations that tolerate the inclusion of a few near-optimal arms. Here, we explore alternative and complementary relaxations that instead allow the learner to identify a relevant *subset* of the Pareto set.

We introduce a unified adaptive sampling strategy, called Adaptive Pareto Exploration (APE), which can be combined with different stopping rules to address multiple relaxations of the Pareto set identification problem within a single framework. We provide a theoretical analysis of the resulting algorithms, establishing instance-dependent sample complexity guarantees and quantifying the gains achieved when identifying at most k Pareto-optimal arms.

This chapter is based on joint work with Émilie Kaufmann and Laura Richert, published in the Proceedings of *NeurIPS 2023*.

3.1	Introduction	64
3.2	Adaptive Pareto Exploration	66
3.2.1	Generic algorithm(s)	67
3.2.2	Our instantiation	71
3.3	Main theoretical results	73
3.3.1	Sketch of proof of Theorem 3.3.1	74
3.4	Numerical study and discussion	76
3.5	Additional proofs	80
3.5.1	Probability of the good event	80
3.5.2	Sample complexity	81

3.1 Introduction

As discussed in Section 1.2 of the introductory chapter, fixed-confidence pure exploration problems aim to make statistically reliable decisions while explicitly controlling the probability of error through a prescribed risk δ . In this chapter, we revisit this setting for *Pareto Set Identification* (PSI) in multi-objective bandits, focusing on the design and analysis of *fully adaptive* algorithms that allocate samples where they are most informative.

Unlike in fixed-budget exploration—studied in Chapter 2—the total number of samples in this setting is not fixed in advance: sampling continues until the confidence in the identified set reaches the desired level $1 - \delta$. This formulation directly aligns with regulatory and ethical requirements in clinical research, where one must control the probability of an incorrect recommendation while minimizing the expected number of enrolled participants.

Beyond clinical trials, active Pareto set identification is also relevant in domains such as hardware and software design (Zuluaga, Sergent, et al. 2013), or multi-objective recommender systems and A/B/n testing, where one jointly optimizes engagement, diversity, and platform objectives (Mehrotra et al. 2020).

Relaxed identification tasks. Before introducing our relaxations, we recall the notion of Pareto dominance. Given two mean vectors $\mu_i, \mu_j \in \mathbb{R}^d$, we say that arm i is (strictly) dominated by arm j , denoted $\mu_i \prec \mu_j$, if $\mu_i^c < \mu_j^c$ for all coordinates $c \in \{1, \dots, d\}$. Intuitively, j outperforms i on every objective. For any $\varepsilon \geq 0$, we define the ε -Pareto set as

$$\mathcal{S}_\varepsilon^* := \{ i \in [K] : \nexists j \in [K] \text{ such that } \mu_i + \varepsilon \mathbf{1} \prec \mu_j \},$$

that is, the set of arms that remain non-dominated when all objectives are relaxed by an additive margin ε . In particular, \mathcal{S}_0^* coincides with the exact Pareto set \mathcal{S}^* .

Exact identification of \mathcal{S}^* can be prohibitively costly when the Pareto frontier is large or when many arms are nearly optimal. We therefore consider several relaxed formulations that capture realistic decision-making constraints:

- ε_1 -PSI (Auer et al. 2016): the returned set \widehat{S} must contain all Pareto-optimal arms and may include arms that become Pareto-optimal after an additive ε_1 relaxation, i.e., $\mathcal{S}^* \subseteq \widehat{S} \subseteq \mathcal{S}_{\varepsilon_1}^*$;
- $(\varepsilon_1, \varepsilon_2)$ -PSI: generalizing the ε -cover of Zuluaga, Krause, et al. 2016, this criterion allows the learner to ignore arms that are nearly equivalent (within ε_2) to others already selected, which is particularly meaningful when several treatments exhibit indistinguishable efficacy;
- ε_1 -PSI- k : introduced in Chapter 2 for the fixed-budget case, this version targets the identification of at most k promising candidates—such as a limited number of vaccine strategies that can proceed to later phases—balancing statistical confidence with practical feasibility.

These relaxations define distinct correctness criteria but share the same fixed-confidence goal: ensuring δ -correct identification while minimizing the expected number of samples.

Related work. The most closely related study is that of [Auer et al. 2016](#), who introduced the ε_1 -PSI relaxation and provided instance-dependent sample complexity bounds under subgaussian noise. In parallel, [Zuluaga, Sargent, et al. 2013](#); [Zuluaga, Krause, et al. 2016](#) proposed structured variants of PSI in which the mean vectors are smooth functions of arm descriptors modeled via Gaussian processes, leading to worst-case guarantees. Their notion of an ε -cover of the Pareto set is conceptually related to our $(\varepsilon_1, \varepsilon_2)$ -PSI formulation. Both Auer’s and Zuluaga’s algorithms rely on *uniform exploration* over an active set combined with accept/reject mechanisms: an approach that can be sample-inefficient when the number of arms is large, as already observed in single-objective Best Arm Identification ([Kaufmann & Kalyanakrishnan 2013](#); [Kaufmann, Cappé, et al. 2016](#)).

Our objective is to develop a fully adaptive exploration strategy that dynamically reallocates samples to the most ambiguous arms and objectives: those whose dominance relations remain uncertain.

Contributions. We introduce *Adaptive Pareto Exploration* (APE), a unified, adaptive framework for fixed-confidence PSI encompassing the three relaxations above. APE extends LUCB- and UGap-style principles from single-objective bandits ([Gabillon et al. 2012](#); [Kalyanakrishnan et al. 2012](#)) to the multi-objective setting by maintaining confidence intervals on pairwise dominance gaps between arms. Depending on the chosen stopping and recommendation rules, APE can address ε_1 -PSI, $(\varepsilon_1, \varepsilon_2)$ -PSI, or ε_1 -PSI- k tasks, with instance-dependent guarantees (Section 3.2). Our analysis further quantifies how relaxing the identification goal—through ε_1 , ε_2 , or k —provably reduces the required sample complexity. Empirical studies (Section 3.4) on synthetic and vaccine-inspired data corroborate these theoretical findings, showing that adaptive exploration yields substantial efficiency gains over uniform baselines.

Learning model. We formalize the *Pareto Set Identification* (PSI) problem and its relaxed variants. Let $K, d \in \mathbb{N}^*$. Each arm $a \in [K] := \{1, \dots, K\}$ is associated with an unknown distribution ν_a over \mathbb{R}^d with mean vector $\mu_a = (\mu_a^1, \dots, \mu_a^d)^\top$. We denote the bandit instance by $\nu := (\nu_1, \dots, \nu_K)$ and its mean matrix by $\mu := (\mu_1, \dots, \mu_K)$. At each round $t = 1, 2, \dots$, the learner selects an arm $A_t \in [K]$ based on past observations and observes an independent draw $Z_t \sim \nu_{A_t}$ with $\mathbb{E}[Z_t | A_t] = \mu_{A_t}$. Each coordinate of Z_t is assumed to be *1-subgaussian* conditionally on A_t . We recall that a random variable X is σ -subgaussian if and only if

$$\forall \lambda \in \mathbb{R}, \log \mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}] \leq \frac{\lambda^2 \sigma^2}{2}.$$

This standard assumption ensures concentration of empirical means around their expectations and covers both Gaussian and bounded observation models.

We denote by \mathbb{P}_ν the law of the stochastic process $(Z_t)_{t \geq 1}$ and by \mathbb{E}_ν the associated expectation. Let $\mathcal{H}_t := \sigma(A_1, Z_1, \dots, A_t, Z_t)$ be the σ -algebra encoding the history up to round t . The algorithm stops sampling at some random time τ , and finally outputs a recommendation \hat{S}_τ . The goal is to make a correct guess with high probability while using as few samples as possible.

Depending on the application, different notions of correctness can be considered, controlled by parameters $\varepsilon_1, \varepsilon_2 \geq 0$ and $k \in [K]$:

Definition 3.1.1. A set $\widehat{S} \subseteq [K]$ is correct for ε_1 -PSI if

$$\mathcal{S}^* \subseteq \widehat{S} \subseteq \mathcal{S}_{\varepsilon_1}^*,$$

where $\mathcal{S}_{\varepsilon_1}^*$ denotes the set of ε_1 -Pareto-optimal arms.

Definition 3.1.2. A set $\widehat{S} \subseteq [K]$ is an $(\varepsilon_1, \varepsilon_2)$ -cover of the Pareto set if $\widehat{S} \subseteq \mathcal{S}_{\varepsilon_1}^*$ and, for any $i \notin \widehat{S}$, either $i \notin \mathcal{S}^*$ or there exists $j \in \widehat{S}$ such that $\mu_i \prec \mu_j + \mathbf{1}\varepsilon_2$. This generalizes the ε -accurate set of [Zuluaga, Krause, et al. 2016](#) (recovered for $\varepsilon_1 = \varepsilon_2 = \varepsilon$) and relaxes the requirement of returning all near-optimal arms.

Definition 3.1.3. A set $\widehat{S} \subseteq [K]$ is correct for ε_1 -PSI- k if (i) $|\widehat{S}| = k$ and $\widehat{S} \subseteq \mathcal{S}_{\varepsilon_1}^*$, or (ii) $|\widehat{S}| < k$ and $\mathcal{S}^* \subseteq \widehat{S} \subseteq \mathcal{S}_{\varepsilon_1}^*$.

Given an objective (one of the above tasks) and a risk parameter $\delta \in (0, 1)$, an algorithm is said to be δ -correct if, with probability at least $1 - \delta$, its recommendation \widehat{S}_τ satisfies the corresponding correctness criterion. The number of samples τ required to reach this guarantee defines the algorithm's sample complexity.

We now introduce two quantities that play a central role in characterizing the (Pareto) optimality and sub-optimality relations between arms. For any arm pair $i, j \in [K]$, define

$$m(i, j) := \min_{1 \leq c \leq d} [\mu_j^c - \mu_i^c], \quad \text{and} \quad M(i, j) := \max_{1 \leq c \leq d} [\mu_i^c - \mu_j^c].$$

As explained in Chapter 1, these quantities capture the smallest and largest coordinate-wise performance gaps between the two arms.

3.2 Adaptive Pareto Exploration

We describe in this section our sampling rule, Adaptive Pareto Exploration and present three stopping and recommendation rules with which it can be combined to solve each of the proposed relaxations.

Let $N_{t,k} := \sum_{s=1}^{t-1} \mathbb{1}(A_s = k)$ be the number of times arm k has been pulled up to round t and $\hat{\mu}_{t,k} := N_{t,k}^{-1} \sum_{s=1}^{N_{t,k}} Z_{k,s}$ be the empirical mean of this arm at time t , where $Z_{k,s}$ denotes the s -th observation drawn *i.i.d.* from ν_k . For any arm pair $i, j \in [K]$, we let

$$m(i, j; t) := \min_{c \leq d} [\hat{\mu}_{t,j}^c - \hat{\mu}_{t,i}^c] \quad \text{and} \quad M(i, j; t) := \max_{c \leq d} [\hat{\mu}_{t,i}^c - \hat{\mu}_{t,j}^c].$$

The empirical Pareto set is defined as

$$\begin{aligned} S(t) &:= \{i \in [K] : \nexists j \in [K] : \hat{\mu}_{t,i} \prec \hat{\mu}_{t,j}\}, \\ &= \{i \in [K] : \forall j \in [K] \setminus \{i\}, M(i, j; t) > 0\}. \end{aligned}$$

3.2.1 Generic algorithm(s)

Adaptive Pareto Exploration relies on a *lower/upper confidence bound* approach, similar to single-objective BAI algorithms like UGap (Gabillon et al. 2012), LUCB (Kalyanakrishnan et al. 2012), and LUCB++ (Simchowitz et al. 2017). These three algorithms identify, in each round, two contentious arms. A leader arm b_t , a current guess for the optimal arm (defined as the empirical best arm or the arm with the smallest upper bound on its sub-optimality gap), and c_t : a contender of b_t ; the arm that is the most likely to outperform b_t (in all three algorithms, it is the arm with the largest upper confidence bound in $[K] \setminus \{b_t\}$). Then either both arms are pulled (LUCB, LUCB++) or the least explored among b_t and c_t is pulled (UGap).

The originality of our sampling rule lies in how to appropriately define b_t and c_t for the multi-objective setting. To define those, we suppose that there exists confidence intervals $[L_{i,j}^c(t, \delta), U_{i,j}^c(t, \delta)]$ on the difference of expected values for each pair of arms (i, j) and each objective $c \in \{1, 2, \dots, d\}$, such that introducing

$$\mathcal{E}_t := \bigcap_{i=1}^K \bigcap_{j \neq i} \bigcap_{c=1}^d \{L_{i,j}^c(t, \delta) \leq \mu_i^c - \mu_j^c \leq U_{i,j}^c(t, \delta)\} \quad \text{and} \quad \mathcal{E} = \bigcap_{t=1}^{\infty} \mathcal{E}_t, \quad (3.1)$$

we have $\mathbb{P}(\mathcal{E}) \geq 1 - \delta$. Concrete choices of these confidence intervals will be discussed in Section 3.2.2.

We further define

$$M^-(i, j; t) := \max_{c \leq d} L_{i,j}^c(t, \delta) \quad \text{and} \quad M^+(i, j; t) := \max_{c \leq d} U_{i,j}^c(t, \delta) \quad (3.2)$$

$$m^-(i, j; t) := -M^+(i, j; t) \quad \text{and} \quad m^+(i, j; t) := -M^-(i, j; t). \quad (3.3)$$

In the sequel, we simplify notation by omitting the dependence on δ and write $L_{i,j}^c(t)$ (resp. $U_{i,j}^c(t)$) in place of $L_{i,j}^c(t, \delta)$ (resp. $U_{i,j}^c(t, \delta)$).

Lemma 3.2.1. *For any round $t \geq 1$, if \mathcal{E}_t holds, then for any $i, j \in [K]$, $M^-(i, j; t) \leq M(i, j) \leq M^+(i, j; t)$ and $m^-(i, j; t) \leq m(i, j) \leq m^+(i, j; t)$.*

Noting that $\mathcal{S}_{\varepsilon_1}^* = \{i \in [K] : \forall j \neq i, M(i, j) + \varepsilon_1 > 0\}$, we define the following set of arms that are likely to be ε_1 -Pareto optimal:

$$\text{OPT}^{\varepsilon_1}(t) := \{i \in [K] : \forall j \in [K] \setminus \{i\}, M^-(i, j; t) + \varepsilon_1 > 0\}.$$

Sampling rule. In round t , Adaptive Pareto Exploration samples A_t , the least pulled arm among two candidate arms b_t and c_t given by

$$\begin{aligned} b_t &:= \operatorname{argmax}_{i \in [K] \setminus \text{OPT}^{\varepsilon_1}(t)} \min_{j \neq i} M^+(i, j; t), \\ c_t &:= \operatorname{argmin}_{j \neq b_t} M^-(b_t, j; t) \end{aligned}$$

The intuition for their definition is the following. Letting i be a fixed arm, note that $M(i, j) > 0$ for some j , if and only if there exists an objective c such that $\mu_i^c > \mu_j^c$, i.e., i is

not dominated by j . Moreover, the larger $M(i, j)$, the more i is non-dominated by j in the sense that there exists c such that $\mu_i^c \gg \mu_j^c$. Therefore, i is strictly optimal if and only if for all $j \neq i$, $M(i, j) > 0$ i.e., $\alpha_i := \min_{j \neq i} M(i, j) > 0$.

The larger α_i , the more optimal i looks, in the sense that for each arm $j \neq i$, there exists a component d_j for which i is way better than j . As the α_i are unknown, we define b_t as the maximizer of an optimistic estimate of the α_i 's. We further restrict the maximization to arms for which we are not already convinced that they are optimal (by Lemma 3.2.1, the arms in $\text{OPT}^{\varepsilon_1}(t)$ are (nearly) Pareto optimal on the event \mathcal{E}). Then, we note that for a fixed arm i , $M(i, j) < 0$ if and only if i is strictly dominated by j . And the smaller $M(i, j)$, the closer j is to dominating i (or largely dominates it): for any component c , $\mu_i^c - \mu_j^c$ is small (or negative). Thus, for a fixed arm i , $\text{argmin}_{j \neq i} M(i, j)$ can be seen as the arm that is the closest to dominating i (or the one that dominates it by the largest margin). By minimizing a lower confidence bound on the unknown quantity $M(b_t, j)$, our contender c_t can be interpreted as the arm that is the most likely to be (close to) dominating b_t . Gathering information on both b_t and c_t can be useful to check whether b_t can indeed be optimal.

Stopping and recommendation rule(s). Depending on the objective, Adaptive Pareto Exploration can be plugged in with different stopping rules, which are summarized in Table 3.1 with their associated recommendations. To define those, we introduce for all $i \in [K]$, $\varepsilon_1, \varepsilon_2 \geq 0$,

$$g_i^{\varepsilon_2}(t) := \max_{j \neq i} [m^-(i, j; t) + \varepsilon_2 \mathbb{1}(j \in \text{OPT}^{\varepsilon_1}(t))] \quad \text{and} \quad h_i^{\varepsilon_1}(t) := \min_{j \neq i} M^-(i, j; t) + \varepsilon_1.$$

and let $g_i(t) := g_i^0(t)$. Introducing

$$Z_1^{\varepsilon_1}(t) := \min_{i \in S(t)} h_i^{\varepsilon_1}(t), \quad \text{and} \quad Z_2^{\varepsilon_1}(t) := \min_{i \in S(t)^c} \max\{g_i(t), h_i^{\varepsilon_1}(t)\},$$

for ε_1 -PSI, our stopping rule is $\tau_{\varepsilon_1} := \inf \{t \geq K : Z_1^{\varepsilon_1}(t) > 0 \wedge Z_2^{\varepsilon_1}(t) > 0\}$ and the associated recommendation is $\mathcal{O}(\tau_{\varepsilon_1})$ where

$$\mathcal{O}(t) := S(t) \cup \{i \in S(t)^c : \nexists j \neq i : m^-(i, j; t) > 0\}$$

consists of the current empirical Pareto set plus some additional arms that have not yet been formally identified as sub-optimal. Those arms should be (ε_1) -Pareto optimal.

For $(\varepsilon_1, \varepsilon_2)$ -PSI we define a similar stopping rule $\tau_{\varepsilon_1, \varepsilon_2}$ where the stopping statistics are respectively replaced by

$$Z_1^{\varepsilon_1, \varepsilon_2}(t) := \min_{i \in S(t)} \max\{g_i^{\varepsilon_2}(t), h_i^{\varepsilon_1}(t)\} \quad \text{and} \quad Z_2^{\varepsilon_1, \varepsilon_2}(t) := \min_{i \in S(t)^c} \max\{g_i^{\varepsilon_2}(t), h_i^{\varepsilon_1}(t)\}$$

with the convention $\min_{\emptyset} = +\infty$, and the recommendation is $\text{OPT}^{\varepsilon_1}(\tau_{\varepsilon_1, \varepsilon_2})$.

To tackle the ε_1 -PSI- k relaxation, we propose to couple τ_{ε_1} with an additional stopping condition checking whether $\text{OPT}^{\varepsilon_1}(t)$ already contains k arms. That is, we stop at $\tau_{\varepsilon_1}^k := \min \{\tau_{\varepsilon_1}, \tau^k\}$ where $\tau^k := \inf \{t \geq K : |\text{OPT}^{\varepsilon_1}(t)| \geq k\}$ with associated recommendation $\text{OPT}^{\varepsilon_1}(\tau^k)$. Depending on the reason for stopping (τ_{ε_1} or τ^k), we follow the corresponding recommendation.

Table 3.1: Stopping conditions and associated recommendation

	Stopping Condition	Recommendation	Objective
τ_{ε_1}	$Z_1^{\varepsilon_1}(t) > 0 \wedge Z_2^{\varepsilon_1}(t) > 0$	$\mathcal{O}(\tau_{\varepsilon_1})$ or $\text{OPT}^{\varepsilon_1}(\tau_{\varepsilon_1})$	ε_1 -PSI
$\tau_{\varepsilon_1, \varepsilon_2}$	$Z_1^{\varepsilon_1, \varepsilon_2}(t) > 0 \wedge Z_2^{\varepsilon_1, \varepsilon_2}(t) > 0$	$\text{OPT}^{\varepsilon_1}(\tau_{\varepsilon_1, \varepsilon_2})$	$(\varepsilon_1, \varepsilon_2)$ -PSI
τ^k	$ \text{OPT}^{\varepsilon_1}(t) \geq k$	$\text{OPT}^{\varepsilon_1}(\tau^k)$	ε_1 -PSI- k

Lemma 3.2.2. *Assume \mathcal{E} holds. For ε_1 -PSI (resp. $(\varepsilon_1, \varepsilon_2)$ -PSI, ε_1 -PSI- k), Adaptive Pareto Exploration combined with the stopping rule τ_{ε_1} (resp. $\tau_{\varepsilon_1, \varepsilon_2}$, resp. $\tau_{\varepsilon_1}^k$) outputs a correct subset.*

Remark 3.1. We decoupled the presentation of the sampling rule from that of the “sequential testing” aspect (stopping and recommendation). We could even go further and observe that multiple tests could actually be run in parallel, for free. If we collect samples with APE (which only depends on ε_1), whenever one of the three stopping conditions given in Table 3.1 triggers, for any value of ε_2 or k , we can decide to stop and make the corresponding recommendation or continue and wait for another “more interesting” stopping condition to be satisfied. If \mathcal{E} holds, a recommendation made at any such time will be correct for the objective associated to the stopping criterion (third column in Table 3.1).

Proof of Lemma 3.2.2.

Step 1: correctness for ε_1 -PSI- k . We show the correctness of ε_1 -PSI- k (for any k), which includes ε_1 -PSI as a special case: ε_1 -PSI- K . Assuming \mathcal{E} holds, let $t = \tau_{\varepsilon_1}^k$ and $i \in \text{OPT}^{\varepsilon_1}(t)$. Since $i \in \text{OPT}^{\varepsilon_1}(t)$, for any $j \neq i$,

$$M(i, j) + \varepsilon_1 \stackrel{\mathcal{E}}{\geq} M^-(i, j; t) + \varepsilon_1 > 0,$$

that is $i \in \mathcal{S}_{\varepsilon_1}^*$. Therefore, on the event \mathcal{E} , $\text{OPT}^{\varepsilon_1}(t) \subseteq \mathcal{S}_{\varepsilon_1}^*$. Thus, if the stopping has occurred because $|\text{OPT}^{\varepsilon_1}(t)| \geq k$, since in this case $\text{OPT}^{\varepsilon_1}(t) \subseteq \mathcal{S}_{\varepsilon_1}^*$, all the recommended arms will be ε_1 -Pareto optimal. On the contrary, if $|\text{OPT}^{\varepsilon_1}(t)| < k$, then from the definition of $\tau_{\varepsilon_1}^k$ it holds that

$$Z_1^{\varepsilon_1}(t) > 0 \quad \text{and} \quad Z_2^{\varepsilon_1}(t) > 0.$$

We will show that $\mathcal{O}(t)$ recalled below:

$$\mathcal{O}(t) := S(t) \cup \{i \in S(t)^c : \nexists j \neq i : m^-(i, j; t) > 0\},$$

is an ε_1 -Pareto set (at $t = \tau_{\varepsilon_1}$) and it is a subset of $\text{OPT}^{\varepsilon_1}(t)$. We note that, for any $i \in \mathcal{O}(t)^c$, by the definition, and since $Z_2(t) > 0$,

$$\exists j \in [K] \text{ such that } m(i, j) \stackrel{\mathcal{E}}{\geq} m^-(i, j; t) > 0,$$

so i is a sub-optimal arm. Therefore,

$$\mathcal{S}^* \subseteq \mathcal{O}(t).$$

Moreover, for any $i \in \mathcal{O}(t) \cap S(t)$, since $Z_1^{\varepsilon_1}(t) > 0$ we have $h_i^{\varepsilon_1}(t) > 0$, so $i \in \text{OPT}^{\varepsilon_1}(t)$, and, as above

$$\min_{j \in [K] \setminus \{i\}} M(i, j) + \varepsilon_1 \stackrel{\mathcal{E}}{\geq} \min_{j \in [K] \setminus \{i\}} M^-(i, j; t) + \varepsilon_1 > 0. \quad (3.4)$$

If $i \in \mathcal{O}(t) \cap S(t)^c$, by definition of $\mathcal{O}(t)$, we have $g_i(t) < 0$. However, since $Z_2^{\varepsilon_1}(t) > 0$, $\max\{g_i(t), h_i^{\varepsilon_1}(t)\} > 0$ so we also have $h_i^{\varepsilon_1}(t) > 0$ and (3.4) applies. Thus, for any $i \in \mathcal{O}(t)$, again,

$$\min_{j \in [K] \setminus \{i\}} M(i, j) + \varepsilon_1 > 0,$$

that is $i \in \mathcal{S}_{\varepsilon_1}^*$, so $\mathcal{S}^* \subseteq \mathcal{O}(t) \subseteq \mathcal{S}_{\varepsilon_1}^*$ and $\mathcal{O}(t) \subseteq \text{OPT}^{\varepsilon_1}(t) \subseteq \mathcal{S}_{\varepsilon_1}^*$, which implies that $\text{OPT}^{\varepsilon_1}(t)$ is itself an (ε_1) -Pareto set.

Finally the correctness for ε_1 -PSI- k and ε_1 -PSI follows on \mathcal{E} .

Step 2: correctness for $(\varepsilon_1, \varepsilon_2)$ -PSI. Let $t = \tau_{\varepsilon_1, \varepsilon_2}$ and $i \in \text{OPT}^{\varepsilon_1}(t)$. Since $i \in \text{OPT}^{\varepsilon_1}(t)$, for any $j \neq i$, $M(i, j) + \varepsilon_1 \stackrel{\mathcal{E}}{\geq} M^-(i, j; t) + \varepsilon_1 > 0$ that is $i \in \mathcal{S}_{\varepsilon_1}^*$. Therefore, on the event \mathcal{E} , $\text{OPT}^{\varepsilon_1}(t) \subseteq \mathcal{S}_{\varepsilon_1}^*$. When the stopping time $\tau_{\varepsilon_1, \varepsilon_2}$ is reached, $Z_1^{\varepsilon_1, \varepsilon_2}(t) > 0$ and $Z_2^{\varepsilon_1, \varepsilon_2}(t) > 0$. Under this condition,

$$\text{OPT}^{\varepsilon_1}(t) \neq \emptyset.$$

Indeed, since $Z_1^{\varepsilon_1, \varepsilon_2}(t) > 0$ and $Z_2^{\varepsilon_1, \varepsilon_2}(t) > 0$, if $\text{OPT}^{\varepsilon_1}(t) = \emptyset$ then, by the stopping rule and since $\text{OPT}^{\varepsilon_1}(t) = \emptyset$, for any arm i , we would have $h_i^{\varepsilon_1}(t) < 0$ and $g_i^{\varepsilon_2}(t) > 0$. That is, for any arm $i \in [K]$,

$$\exists j \neq i \text{ such that } m(i, j) \stackrel{\mathcal{E}}{>} m^-(i, j, t) > 0,$$

so every arm would be strictly dominated, which is impossible since the Pareto set cannot be empty. Then, $\text{OPT}^{\varepsilon_1}(t) \neq \emptyset$ and for any $i \in \mathcal{O}(t)^c = \text{OPT}^{\varepsilon_1}(t)^c$, by the stopping rule it holds that $\max\{g_i^{\varepsilon_2}(t), h_i^{\varepsilon_1}(t)\} > 0$. Further, noting that for such arm $i \in \text{OPT}^{\varepsilon_1}(t)^c$, $h_i^{\varepsilon_1}(t) < 0$, we thus have $g_i^{\varepsilon_2}(t) > 0$, which, recalling that

$$g_i^{\varepsilon_2}(t) := \max_{j \neq i} [m^-(i, j; t) + \varepsilon_2 \mathbb{1}(j \in \text{OPT}^{\varepsilon_1}(t))],$$

implies

$$\exists j \neq i : m^-(i, j; t) + \varepsilon_2 \mathbb{1}(j \in \text{OPT}^{\varepsilon_1}(t)) > 0,$$

thus, on the event \mathcal{E} we have

$$m(i, j) + \varepsilon_2 \mathbb{1}(j \in \text{OPT}^{\varepsilon_1}(t)) > 0.$$

Therefore, for such arm i , either

- a) $\exists j \in [K]$ such that $m(i, j) > 0$ that is $\mu_i \prec \mu_j$ or
- b) $\exists j \in \text{OPT}^{\varepsilon_1}(t)$ such that $m(i, j) + \varepsilon_2 > 0$ that is $\mu_i \prec \mu_j + \varepsilon_2 \mathbf{1}$.

Combining these, $\text{OPT}^{\varepsilon_1}(t) \subseteq \mathcal{S}_{\varepsilon_1}^*$ and, for any $i \notin \text{OPT}^{\varepsilon_1}(t)$, either $i \notin \mathcal{S}_{\varepsilon_1}^*$ (i is ε_1 -sub-optimal) or there exists $j \in \text{OPT}^{\varepsilon_1}(t)$ such that $\mu_i \prec \mu_j + \varepsilon_2 \mathbf{1}$. Thus, $\text{OPT}^{\varepsilon_1}(t)$ is an $(\varepsilon_1, \varepsilon_2)$ -cover of the Pareto set and APE is correct for $(\varepsilon_1, \varepsilon_2)$ -PSI. \square

3.2.2 Our instantiation

We propose to instantiate the algorithms with confidence intervals on the difference between pair of arms. For any pair $i, j \in [K]$, we define a function $\beta_{i,j}$ such that for any $c \in [d]$, $U_{i,j}^c(t) = \hat{\mu}_{t,i}^c - \hat{\mu}_{t,j}^c + \beta_{i,j}(t, \delta)$ and $L_{i,j}^c(t) = \hat{\mu}_{t,i}^c - \hat{\mu}_{t,j}^c - \beta_{i,j}(t, \delta)$. We use the following confidence bonus from [Kaufmann & W.-M. Koolen 2021](#) for time-uniform concentration:

$$\beta_{i,j}(t, \delta) := 2 \sqrt{\left(C^g \left(\frac{\log(\frac{K_1}{\delta})}{2} \right) + \sum_{a \in \{i,j\}} \log(4 + \log(N_{t,a})) \right) \left(\sum_{a \in \{i,j\}} \frac{1}{N_{t,a}} \right)}, \quad (3.5)$$

where $K_1 := K(K-1)d/2$ and $C^g \approx x + \log(x)$ is a calibration function. As a result, we have the simple expressions $M^\pm(i, j; t) = M(i, j; t) \pm \beta_{i,j}(t, \delta)$ and $m^\pm(i, j; t) = m(i, j; t) \pm \beta_{i,j}(t, \delta)$. As an example, we state in [Algorithm 3.1](#) the pseudo-code of APE combined with the stopping rule suited for the k -relaxation of ε_1 -PSI, which we refer to as the ε_1 -APE- k algorithm.

Algorithm 3.1: APE: Adaptive Pareto Exploration

Require : (optional) ε_1 -PSI relaxation parameter $\varepsilon_1 \geq 0$; (optional) parameter $\varepsilon_2 \geq 0$; (optional) k -relaxation parameter $k \in [K]$

- 1 Initialize: sample each arm once, set $t = K$, $N_{K,i} = 1$ for any $i \in [K]$
- 2 **foreach** $t = K + 1, \dots$, **do**
 - 3 *// Empirical non-dominated set at time t*
 $S(t) \leftarrow \{i \in [K] : \forall j \in [K] \setminus \{i\}, M(i, j; t) > 0\}$
 - 4 *// Relaxed ε_1 -Pareto candidates*
 $OPT^{\varepsilon_1}(t) \leftarrow \{i \in [K] : \forall j \in [K] \setminus \{i\}, M(i, j; t) - \beta_{i,j}(t, \delta) + \varepsilon_1 > 0\}$
 - 5 **if** $|OPT^{\varepsilon_1}(t)| \geq k$ **then**
 - 6 **break** and output $OPT^{\varepsilon_1}(t)$
 - 7
 - 8 **if** $Z_1^{\varepsilon_1}(t) > 0 \wedge Z_2^{\varepsilon_1}(t) > 0$ **then**
 - 9 **break** and output
 $\mathcal{O}(t) = S(t) \cup \{i \in S(t)^c : \nexists j \neq i : m(i, j; t) - \beta_{i,j}(t, \delta) > 0\}$
 - 10 *// Leader / Most promising arm outside $OPT^{\varepsilon_1}(t)$*
 $b_t \leftarrow \operatorname{argmax}_{i \in [K] \setminus OPT^{\varepsilon_1}(t)} \min_{j \neq i} [M(i, j; t) + \beta_{i,j}(t, \delta)]$
 - 11 *// Challenger / Most critical opponent against b_t*
 $c_t \leftarrow \operatorname{argmin}_{j \neq b_t} [M(b_t, j; t) - \beta_{b_t,j}(t, \delta)]$
 - 12 *// Sample the less explored of $\{b_t, c_t\}$*
 Sample $A_t \leftarrow \operatorname{argmin}_{i \in \{b_t, c_t\}} N_{t,i}$; $t \leftarrow t + 1$; update statistics

Other instantiations are possible, such as the ones based on individual confidence bounds of the form $U_{i,j}(t) = U_i(t) - L_j(t)$ where $[L_i(t), U_i(t)]$ is a confidence interval on μ_i (Kalyanakrishnan et al. 2012; Auer et al. 2016) or refined KL confidence sequences (Kaufmann & Kalyanakrishnan 2013).

Special case: Best Arm Identification. When $d = 1$, APE reduces to sampling at each round t , the least sampled among

$$b_t := \operatorname{argmax}_i \left[\min_{j \neq i} U_{i,j}(t) \right], \quad \text{and} \quad c_t := \operatorname{argmin}_{j \neq b_t} L_{b_t,j}(t). \quad (3.6)$$

To compare against LUCB (Kalyanakrishnan et al. 2012) and UGap (Gabillon et al. 2012) which use confidence intervals on single arms, we would have $\beta_{i,j}(t, \delta) := \beta_i(t, \delta) + \beta_j(t, \delta)$, and $L_i(t) := \hat{\mu}_{t,i} - \beta_i(t, \delta)$ and $U_i(t) := \hat{\mu}_{t,i} + \beta_i(t, \delta)$: lower and upper confidence bounds on μ_i . Then (3.6) rewrites as

$$b_t := \operatorname{argmax}_i \left[U_i(t) - \max_{j \neq i} L_j(t) \right], \quad \text{and} \quad c_t := \operatorname{argmax}_{j \neq b_t} U_j(t).$$

This resembles the sampling rule of UGap, which defines

$$b_t^{\text{UGap}} := \operatorname{argmax}_i \left[L_i(t) - \max_{j \neq i} U_j(t) \right], \quad \text{and} \quad c_t^{\text{UGap}} := \operatorname{argmax}_{j \neq b_t} U_j(t),$$

and also pulls the least sampled so far. We further note that when $\varepsilon_1 = 0$, for any $i \in S(t)^c$, $g_i(t) > h_i^0(t)$. Indeed, by definition,

$$h_i^0(t) = \min_{j \neq i} [M(i, j; t) - \beta_{i,j}(t, \delta)] = \min_{j \neq i} [-m(i, j; t) - \beta_{i,j}(t, \delta)]$$

and since $i \in S(t)^c$, there exists i^* such that $m(i, i^*; t) > 0$ (i.e., $\hat{\mu}_{t,i} \prec \hat{\mu}_{t,i^*}$) and so

$$-m(i, i^*; t) - \beta_{i,i^*}(t, \delta) < m(i, i^*; t) - \beta_{i,i^*}(t, \delta).$$

Therefore,

$$\min_{j \neq i} [-m(i, j; t) - \beta_{i,j}(t, \delta)] := h_i^0(t) < \max_{j \neq i} [m(i, j; t) - \beta_{i,j}(t, \delta)] := g_i(t).$$

Thus for $\varepsilon_1 = 0$, $Z_2^0(t) = \min_{i \in S(t)^c} g_i(t)$. In particular, when $d = 1$, $\varepsilon_1 = 0$, the stopping time τ_0 can be simplified to $\tau_0 = \inf\{t \geq K : Z_1^0(t) > 0\}$. Letting \hat{a}_t denote the empirical best arm after t rounds, the stopping rule of APE (with the instantiation based on confidence intervals on pairs of arms) reduces to

$$\tau_0 = \inf \left\{ t \geq K : \forall i \neq \hat{a}_t, \frac{(\hat{\mu}_{\hat{a}_t,t} - \hat{\mu}_{t,i})^2}{2 \left(\frac{1}{N_{t,\hat{a}_t}} + \frac{1}{N_{t,i}} \right)} \geq 2C^g \left(\frac{\log(K_1/\delta)}{2} \right) + 2 \sum_{a \in \{\hat{a}_t, i\}} \log(4 + \log(N_{t,a})) \right\}$$

which is very close to a Generalized Likelihood Ratio (GLR) stopping rule assuming Gaussian rewards with variance 1 (which is also known to be correct for subgaussian rewards) Garivier & Kaufmann 2016; Kaufmann & W.-M. Koolen 2021.

3.3 Main theoretical results

We state our main theorem on the sample complexity of our algorithms and give a sketch of its proof. First, let us introduce some quantities that are needed to state the theorem.

As presented in Chapter 1, the sample complexity of ε_1 -Pareto set identification scales as a sum over the arms i of $1/(\Delta_i \vee \varepsilon_1)^2$ where Δ_i is called the sub-optimality gap of arm i and is defined as follows. For a sub-optimal arm $i \notin \mathcal{S}^*(\mu)$,

$$\Delta_i := \max_{j \in \mathcal{S}^*(\mu)} m(i, j),$$

which is the smallest quantity that should be added component-wise to μ_i to make i appear Pareto optimal w.r.t. $\{\mu_i : i \in [K]\}$. For a Pareto optimal arm $i \in \mathcal{S}^*(\mu)$, the definition is more involved:

$$\Delta_i := \begin{cases} \min_{j \in [K] \setminus \{i\}} \Delta_j & \text{if } \mathcal{S}^*(\mu) := \{i\} \\ \min\{\delta_i^+, \delta_i^-\} & \text{else,} \end{cases}$$

where

$$\delta_i^+ := \min_{j \in \mathcal{S}^*(\mu) \setminus \{i\}} \min\{M(i, j), M(j, i)\} \quad \text{and} \quad \delta_i^- := \min_{j \in [K] \setminus \mathcal{S}^*(\mu)} [(M(j, i))_+ + \Delta_j],$$

with $(x)_+ := \max\{x, 0\}$, $\forall x \in \mathbb{R}$. For $d = 1$, these gaps match the classical gaps in Best Arm Identification (Audibert & Bubeck 2010; Kaufmann, Cappé, et al. 2016).

We also introduce some additional notation required to express the contribution of the k -relaxation. Let $1 \leq k \leq K$. For any arm i , let $\omega_i = \min_{j \neq i} M(i, j)$ and define

$$\omega^k := \max_{i \in [K]}^k \omega_i, \quad \mathcal{S}^{*,k} := \operatorname{argmax}_{i \in [K]}^{1 \dots k} \omega_i,$$

with the k -th max and the first k argmax operators. Observe that $\omega^k > 0$ if and only if $|\mathcal{S}^*(\mu)| \geq k$.

Theorem 3.3.1. *Let $\delta \in (0, 1)$ be a risk parameter, $\varepsilon_1, \varepsilon_2 \geq 0$, $K \geq 2$, and $k \leq K$. Consider a bandit instance ν whose arm marginals are 1-subgaussian. Then, with probability at least $1 - \delta$, the APE algorithm returns a correct recommendation for either the ε_1 -PSI- k or the $(\varepsilon_1, \varepsilon_2)$ -PSI task (with the corresponding stopping rule), and it stops after at most*

$$\sum_{a \in [K]} \frac{88}{\tilde{\Delta}_a^2} \log \left(\frac{2K(K-1)d}{\delta} \log \left(\frac{12e}{\tilde{\Delta}_a} \right) \right)$$

samples, where for each $a \in [K]$,

$$\tilde{\Delta}_a := \begin{cases} \max\{\Delta_a, \varepsilon_1, \omega^k\}, & \text{for the } \varepsilon_1\text{-PSI-}k \text{ task,} \\ \max\{\Delta_a, \varepsilon_1\}, & \text{for the } (\varepsilon_1, \varepsilon_2)\text{-PSI task.} \end{cases}$$

Our theoretical results extend and refine existing bounds for approximate and relaxed Pareto set identification.

When $k = K$, the proposed algorithm provides a δ -correct solution for ε_1 -PSI, which corresponds to the ε_1 -PSI- K task. It improves upon [Auer et al. 2016](#) by constant factors and replaces their $\log \Delta^{-2}$ dependence with a milder $\log \log \Delta^{-1}$ term, while nearly matching their lower bound (Theorem 1.3.1).

Our strongest result concerns the ε_1 -PSI- k relaxation, which introduces a new notion of sub-optimality gap, quantifying how the relaxation parameter k reduces the effective sample complexity. In particular, for any arm $i \in \mathcal{S}^* \setminus \mathcal{S}^{*,k}$, we have $\max\{\Delta_i, \omega^k\} = \omega^k$, indicating that one should not pay more than the cost of identifying the k -th optimal arm. This observation parallels the findings of [Roy Chaudhuri & Kalyanakrishnan 2019](#) for the “top- k among the best- m ” problem. However, the authors have shown the relaxation only for the best m arms, while our result shows that even the sub-optimal arms should be sampled less.

For the general $(\varepsilon_1, \varepsilon_2)$ -PSI formulation, our upper bound does not make the dependence on ε_2 explicit, yet nearly matching lower bounds suggest that this dependence may vanish in some cases. In particular, when $d = 1$, $\varepsilon_1 = 0$, and $\varepsilon_2 > 0$ with a unique optimal arm a_* , $(0, \varepsilon_2)$ -PSI reduces to exact Best Arm Identification. Indeed, in this case, any $\widehat{S} \subseteq \mathcal{S}^* = \{a_*\}$ trivially satisfies the covering condition. The corresponding lower bounds ([Garivier & Kaufmann 2016](#); [Kaufmann, Cappé, et al. 2016](#); [Simchowitz et al. 2017](#)) are thus independent of ε_2 . This argument extends to the multi-objective setting when the Pareto set is a singleton, confirming that ε_2 has no impact on the intrinsic difficulty of such instances.

3.3.1 Sketch of proof of Theorem 3.3.1

Proof. Using Proposition 24 of [Kaufmann & W.-M. Koolen 2021](#) it is simple to prove that the choice of $\beta_{i,j}$ in (3.5) yields $\mathbb{P}(\mathcal{E}) \geq 1 - \delta$ for the good event \mathcal{E} defined in (3.1). Combining this result with Lemma 3.2.2 yields that APE is correct for each task with probability at least $1 - \delta$.

The idea of the remaining proof is to show that under the event \mathcal{E} , if APE has not stopped at the end of round t , then the selected arm A_t has not been explored enough. The first lemma showing this is specific to the stopping rule $\tau_{\varepsilon_1}^k$ used for ε_1 -PSI- k .

Lemma 3.3.2. *Let $\varepsilon_1 \geq 0$ and $k \in [K]$. If \mathcal{E}_t holds and $t < \tau_{\varepsilon_1}^k$ then $\omega^k \leq 2\beta_{A_t, A_t}(t, \delta)$.*

The following two lemmas are more general as they apply to different stopping rules.

Lemma 3.3.3. *Let $\varepsilon_1 \geq 0$. Let $\tau = \tau_{\varepsilon_1}^k$ for some $k \in [K]$ or $\tau = \tau_{\varepsilon_1, \varepsilon_2}$ for some $\varepsilon_2 \geq 0$. If \mathcal{E}_t holds and $t < \tau$ then $\Delta_{A_t} \leq 2\beta_{A_t, A_t}(t, \delta)$.*

Lemma 3.3.4. *Let $\varepsilon_1 \geq 0$ and τ be as in Lemma 3.3.3. If \mathcal{E}_t holds and $t < \tau$ then $\varepsilon_1 \leq 2\beta_{A_t, A_t}(t, \delta)$.*

The proofs of these three lemmas are given at the end of the chapter. They heavily rely on the specific definition of b_t and c_t . In particular, to prove Lemma 3.3.3 and 3.3.4, we first establish that when $t < \tau$, any arm $j \in [K]$ satisfies $m(b_t, j; t) \leq \beta_{b_t, j}(t, \delta)$.

The upper bound on the sample complexity follows, since using the lemmas above, we have that, if ε_1 -APE- k has not stopped during round t , i.e., $t < \tau_{\varepsilon_1}^k$ and the event \mathcal{E}_t holds, then

- $\omega^k \leq 2\beta_{A_t, A_t}(t, \delta)$ (Lemma 3.3.2),
- $\Delta_{A_t} \leq 2\beta_{A_t, A_t}(t, \delta)$ (Lemma 3.3.3),
- $\varepsilon_1 \leq 2\beta_{A_t, A_t}(t, \delta)$ (Lemma 3.3.4)

hold simultaneously. Introducing $\tilde{\Delta}_a := \max\{\omega^k, \varepsilon_1, \Delta_a\}$,

$$\begin{aligned} \tau_{\varepsilon_1}^k \mathbb{1}(\mathcal{E}) - 1 &\leq \sum_{t=1}^{\infty} \mathbb{1}(\mathcal{E}) \mathbb{1}(\tau_{\varepsilon_1}^k > t), \\ &\leq \sum_{t=1}^{\infty} \mathbb{1}(\max\{\omega^k, \varepsilon_1, \Delta_{A_t}\} \leq 2\beta_{A_t, A_t}(t, \delta)) \\ &\leq \sum_{a=1}^K \sum_{t=1}^{\infty} \mathbb{1}(\{A_t = a\} \wedge \{\tilde{\Delta}_a \leq 2\beta_{a, a}(t, \delta)\}) \\ &\leq \sum_{a=1}^K \inf\{n \geq 2 : \tilde{\Delta}_a > 2\beta^n\}, \end{aligned}$$

where β^n is the expression of $\beta_{i, j}(t, \delta)$ when $N_{t, i} = N_{t, j} = n$, that is

$$\beta^n = 2\sqrt{\left(C^g \left(\frac{\log\left(\frac{K+1}{\delta}\right)}{2}\right) + 2\log(4 + \log(n))\right) \frac{2}{n}}.$$

Then, an inversion result (see Lemma 19 in [Kone, Kaufmann, et al. 2023](#)) yields

$$\inf\{s \geq 2 : 2\beta^s < \tilde{\Delta}_a\} \leq \frac{88}{\tilde{\Delta}_a^2} \log\left(\frac{2K(K-1)d}{\delta} \log\left(\frac{12e}{\tilde{\Delta}_a}\right)\right).$$

Therefore,

$$\tau_{\varepsilon_1}^k \mathbb{1}(\mathcal{E}) \leq \sum_{a \in [K]} \frac{88}{\tilde{\Delta}_a^2} \log\left(\frac{2K(K-1)d}{\delta} \log\left(\frac{12e}{\tilde{\Delta}_a}\right)\right).$$

The same reasoning applies to $\tau_{\varepsilon_1, \varepsilon_2}$ as $\tau_{\varepsilon_1, \varepsilon_2} \leq \tau_{\varepsilon_1}^K$ holds with probability 1. □

LUCB-like instantiation. Theorem 3.3.1 bounds the sample complexity of APE on \mathcal{E} but as, for many algorithms in pure-exploration, we do not control what happens on \mathcal{E}^c . However, when APEs run with confidence bonuses similar to LUCB1 (Kalyanakrishnan et al. 2012), in the form

$$\beta_i(t, \delta) := \sqrt{\frac{2}{N_{t,i}} \log \left(\frac{5K dt^4}{2\delta} \right)}, \quad \text{and} \quad \forall (i, j) \in [K], \quad \beta_{i,j}(t, \delta) = \beta_i(t, \delta) + \beta_j(t, \delta), \quad (3.7)$$

we are able to tightly upper bound the expected stopping time. The following result is proven in Kone, Kaufmann, et al. 2023 following the lines of the proof of Theorem 3.3.1.

Theorem 3.3.5 (Expected sample complexity of APE (Kone, Kaufmann, et al. 2023, Theorem 5)). *Let $\varepsilon_1, \varepsilon_2 \geq 0$, $k \leq K$, and ν be a bandit instance with 1-subgaussian marginals. When APE is run with the parameters $(\beta_i)_i$ defined in (3.7), it is δ -correct for either ε_1 -PSI- k or $(\varepsilon_1, \varepsilon_2)$ -PSI. Moreover, its expected sample complexity satisfies*

$$\mathbb{E}_\nu[\tau] \leq \mathcal{O} \left(H(\nu) \log \frac{d H(\nu)}{\delta} \right),$$

where

$$H(\nu) := \sum_{a=1}^K \tilde{\Delta}_a^{-2}, \quad \text{with} \quad \tilde{\Delta}_a := \begin{cases} \max\{\Delta_a, \varepsilon_1, \omega^k\}, & \text{for the } \varepsilon_1\text{-PSI-}k \text{ task,} \\ \max\{\Delta_a, \varepsilon_1\}, & \text{for the } (\varepsilon_1, \varepsilon_2)\text{-PSI task.} \end{cases}$$

Kone, Kaufmann, et al. 2023 prove the following lower bound showing that in some scenarios, APE is optimal for ε_1 -PSI- k , up to $d \log(K)$ and constant multiplicative terms. We note that for ε_1 -PSI a lower bound featuring the gaps Δ_a and ε_1 was already derived by Auer et al. 2016.

Theorem 3.3.6 (Theorem 2 of Kone, Kaufmann, et al. 2023). *There exists a bandit instance ν with $|\mathcal{S}^*| = p \geq 3$ such that for $k \in \{p, p-1, p-2\}$, any δ -correct algorithm for 0-PSI- k verifies*

$$\mathbb{E}_\nu[\tau_\delta] \geq \frac{1}{d} \log \frac{1}{\delta} \sum_{a=1}^K \frac{1}{(\Delta_a^k)^2},$$

where $\Delta_a^k := \Delta_a + \omega^k$ and τ_δ is the stopping time of the algorithm.

3.4 Numerical study and discussion

Experimental setup. We compare *Adaptive Pareto Exploration* (APE) to the algorithm of Auer et al. 2016 (denoted PSI-Unif-Elim). For a fair comparison, both use our pairwise confidence bonuses $\beta_{i,j}(t, \delta)$.¹ Since anytime confidence bounds are known to be conservative, we set $K_1 = 1$ in (3.5) instead of its theoretical union-bound value. Unless stated otherwise, $\delta = 0.1$. Across all runs, the empirical error is well below δ .

¹Auer et al. 2016 already suggested this heuristic; we instantiate it for both methods.

Real-world-inspired evaluation (COV-BOOST). We reuse the COV-BOOST instance introduced in Chapter 2 and fully documented in Appendix 7. Briefly, the bandit has $K = 20$ arms (vaccine strategies) and $d = 3$ immunogenicity indicators derived from the processed trial data; we simulate a multivariate Gaussian model using the means and variances reported there. All preprocessing details, including the log-normal assumption and pooled variance estimates, are deferred to Appendix 7. We set $\varepsilon_1 = 0$ and compare PSI-Unif-Elim to APE- k for several values of k .

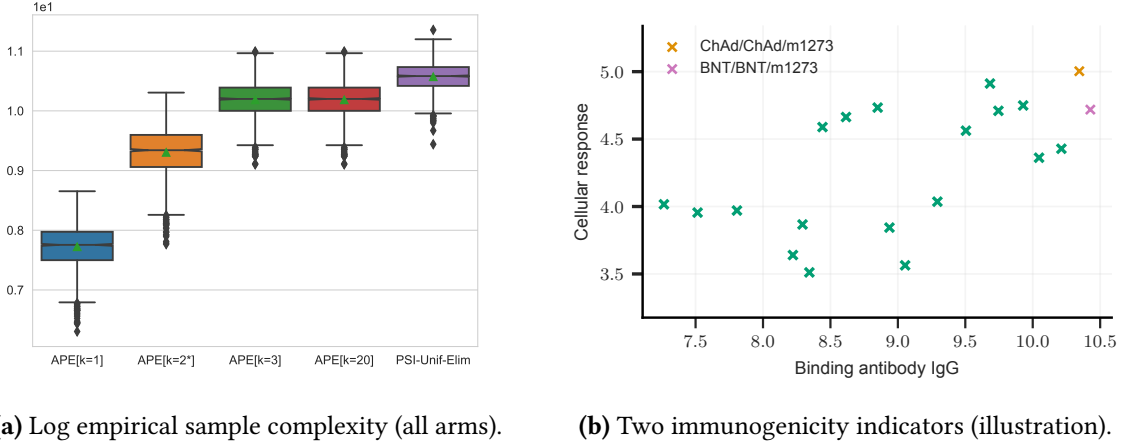


Figure 3.1: COV-BOOST: APE- k vs. PSI-Unif-Elim at $\delta = 0.1$ (2000 runs).

Since $|\mathcal{S}^*| = 2$, APE (with $k \geq 3$) outperforms PSI-Unif-Elim on exact PSI. The k -relaxation substantially reduces sample complexity: APE-2 stops as soon as two optimal arms are certified, whereas for $k = 3$ remaining arms must be ruled out as suboptimal, explaining the clear gap between $k = 2$ and $k = 3$.

Random Bernoulli instances. We generate 2000 multivariate Bernoulli bandits with $K = 5$ and $d \in \{2, 4, 8\}$ (independent marginals), set $\delta = 0.1$ and $\varepsilon_1 = 0.005$ (for runtime), and average the sample complexity over runs.

Table 3.2: Average sample complexity over 2000 random instances ($K = 5$). Average $|\mathcal{S}^*|$ is (2.30, 4.06, 4.93) for $d = 2, 4, 8$.

	ε_1 -APE-1	ε_1 -APE-2	ε_1 -APE-3	ε_1 -APE-4	ε_1 -APE-5	ε_1 -PSI-Unif-Elim
$d = 2$	811	39,530	109,020	145,777	150,844	190,625
$d = 4$	214	6,410	19,908	68,061	124,001	157,584
$d = 8$	119	204	405	1,448	20,336	27,270

APE (with $k = K$) uses 20–25% fewer samples than PSI-Unif-Elim, with larger gains as d increases—and consequently $|\mathcal{S}^*|$ as well. The k -relaxation drastically lowers budgets relative to exact PSI.

APE on $(\varepsilon_1, \varepsilon_2)$ -PSI. We now examine the empirical behavior of $(\varepsilon_1, \varepsilon_2)$ -APE for identifying an $(\varepsilon_1, \varepsilon_2)$ -cover. We set $\varepsilon_1 = 0$, $\delta = 0.01$, and vary $\varepsilon_2 \in \{0, 0.05, 0.1, 0.2, 0.25\}$. Results are averaged over 2000 independent runs with different random seeds on the same instance (described in Figure 3.2a). The environment is a multivariate Bernoulli bandit with $K = 5$ and $d = 2$ (independent marginals). In this toy example, $(\varepsilon_1, \varepsilon_2)$ -PSI is particularly meaningful since several arms are nearly equivalent: three Pareto-optimal vectors are manually chosen, while the remaining two are drawn uniformly at random.

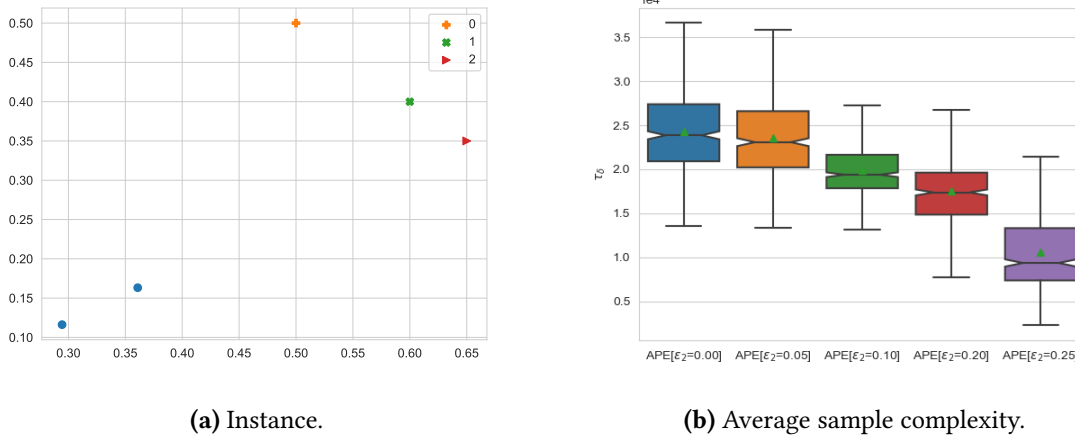


Figure 3.2: Toy instance with $\mathcal{S}^* = \{0, 1, 2\}$. The difference on the x- and y-axes is 0.1 between arms 0 and 1, and 0.05 between arms 1 and 2. Right: average sample complexity over 2000 runs.

As shown in Figure 3.2b, the sample complexity decreases as ε_2 increases. This trend is confirmed in Figure 3.3, which reports both the average sample complexity and the mean size of the returned cover for 50 evenly spaced values of $\varepsilon_2 \in [0, 0.5]$. Each major drop in Figure 3.3b corresponds to a value of ε_2 at which APE removes an (otherwise optimal) arm from the cover to save samples (see Figure 3.3a). Thus, a decrease in sample complexity roughly coincides with a reduction in the number of arms included in the final recommendation.

For large values of ε_2 (e.g., > 0.3), the sample complexity plateaus (Figure 3.3b). This behavior stems from the fact that, even when ε_2 is large, the algorithm must still certify at least one truly optimal arm (as reflected by the single-element covers in Figure 3.3a). Empirically, this limiting complexity is close to that of 0-APE-1 on the same instance (≈ 4073 samples).

Finally, Figure 3.4 shows the frequency with which each arm appears in the recommended set for three representative values of ε_2 . For $\varepsilon_2 = 0.15$, arm 0 is always recommended, while the others appear in about half of the runs. For $\varepsilon_2 = 0.4$, the algorithm almost always returns only arm 0—the one with the largest ω_i term, i.e., the easiest to identify as optimal.

Sampling allocations of APE vs. PSI-Unif-Elim. We now examine how the sampling allocations of APE and PSI-Unif-Elim differ on specific instances. In some configurations,

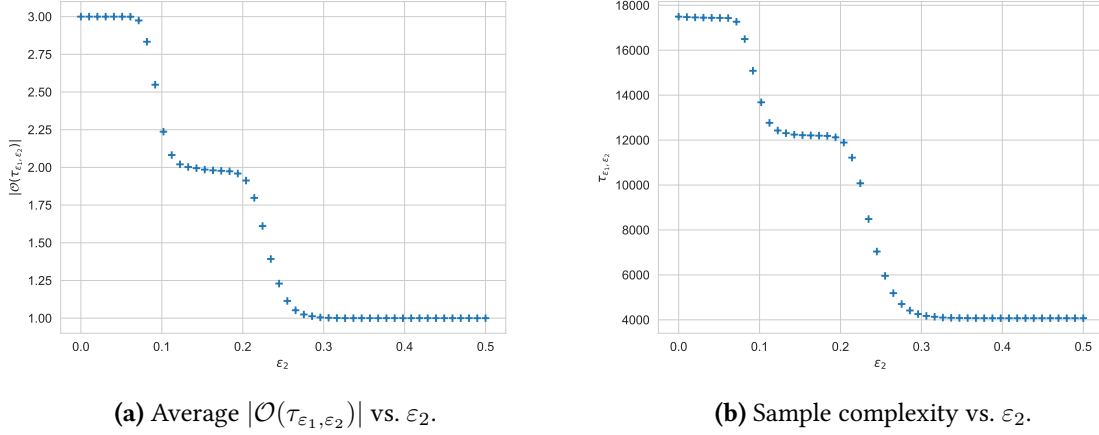


Figure 3.3: Average size of the returned cover (left) and corresponding sample complexity (right), both averaged over 2000 runs. The empirical error probability was negligible in all cases.

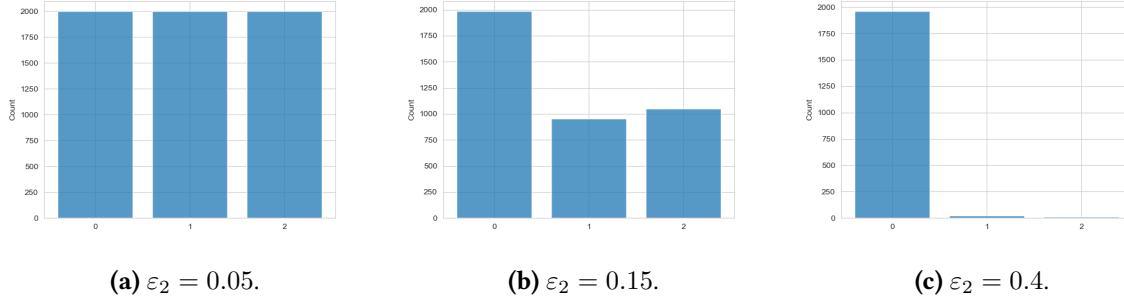


Figure 3.4: Frequency of each arm in the recommended set across 2000 runs for three representative values of ϵ_2 .

APE requires up to three times fewer samples than PSI-Unif-Elim. This improvement stems from PSI-Unif-Elim’s conservative strategy: it continues sampling arms already identified as (nearly) optimal until these arms are proven not to dominate any remaining active ones. For instance, in Figure 3.5a, arm 2 is easily recognized as optimal, yet it slightly dominates the suboptimal arm 1. Consequently, PSI-Unif-Elim keeps sampling arm 2 until arm 1 is eventually removed from the active set—likely when it is shown to be dominated by arm 3. Our adaptive sampling rule avoids this inefficiency by dynamically reallocating samples to the most uncertain comparisons instead of repeatedly pulling confirmed optimal arms.

As shown in Figure 3.5b, APE achieves roughly half the sample complexity of PSI-Unif-Elim on this instance. Table 3.3 details the average number of pulls taken by PSI-Unif-Elim divided by that of 0-APE- K for each arm. The primary source of inefficiency is evident: arm 2 is sampled nearly six times more often by PSI-Unif-Elim, highlighting its tendency to oversample dominant arms that are already known to be optimal.

Discussion. Across real-world-inspired and synthetic settings, APE consistently improves over PSI-Unif-Elim on exact PSI and delivers large additional savings under the k - and

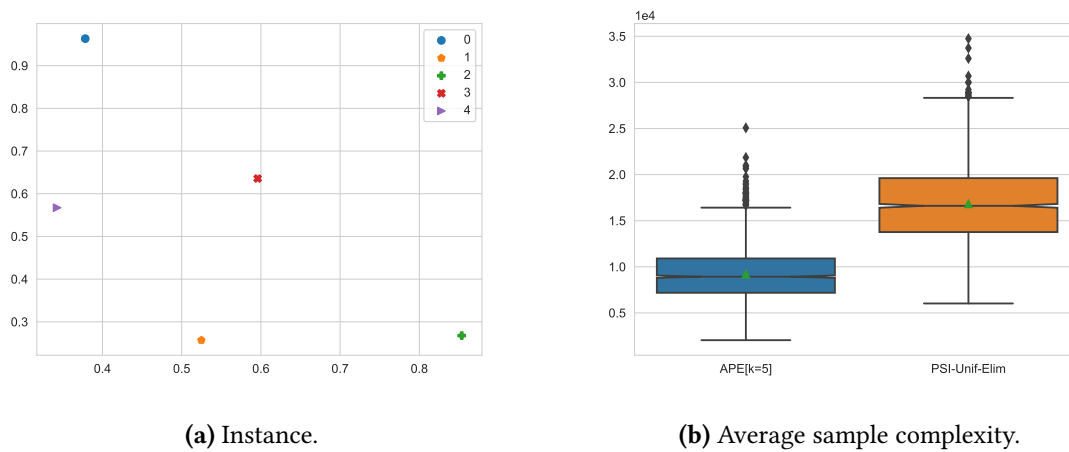


Figure 3.5: Illustrative instance with $\mathcal{S}^* = \{0, 2, 3\}$. Right: average sample complexity across repeated trials. APE avoids oversampling already-certified optimal arms.

Table 3.3: Average number of pulls under PSI-Unif-Elim divided by that under 0-APE- K for each arm. The largest gap arises from arm 2, which PSI-Unif-Elim oversamples to maintain dominance guarantees.

Arm	0	1	2	3	4
Average pull ratio (PSI-Unif-Elim / APE)	3.08	1.36	5.68	1.28	1.40

$(\varepsilon_1, \varepsilon_2)$ -relaxations. On average (Table 3.2), APE achieves about a 20% reduction in sample complexity for exact PSI, with substantially larger gains on specific instances where multiple Pareto-optimal arms dominate a suboptimal arm. As emphasized in Chapter 1 and Chapter 2, these relaxations align with practical decision constraints (e.g., limited advancement slots or near-duplicate treatments), making APE both statistically and operationally attractive.

3.5 Additional proofs

In this section, we gather results used in the correctness and sample-complexity analyses of APE, and provide their proofs. We recall the definition of the events

$$\mathcal{E}_t = \bigcap_{i=1}^K \bigcap_{j \neq i}^d \bigcap_{c=1}^d \{L_{i,j}^c(t) \leq \mu_i^c - \mu_j^c \leq U_{i,j}^c(t)\} \quad \text{and} \quad \mathcal{E} = \bigcap_{t=1}^{\infty} \mathcal{E}_t.$$

3.5.1 Probability of the good event

Proof of Lemma 3.2.1. This result follows from the definition of \mathcal{E}_t . Since

$$\mathcal{E}_t := \bigcap_{i=1}^K \bigcap_{j \neq i}^d \bigcap_{c=1}^d \{L_{i,j}^c(t) \leq \mu_i^c - \mu_j^c \leq U_{i,j}^c(t)\},$$

if \mathcal{E}_t holds, then for any i, j

$$M^-(i, j; t) := \max_c L_{i,j}^c(t) \leq M(i, j) := \max_c [\mu_i^c - \mu_j^c] \leq \max_c U_{i,j}^c(t) := M^+(i, j; t),$$

and the second point follows by noting that $m(i, j) = -M(i, j)$ and

$$m^+(i, j; t) := -M^-(i, j; t); m^-(i, j; t) := -M^+(i, j; t)$$

for any pair of arms. □

We note that when the algorithm uses a confidence bonus of the form $(\hat{\mu}_{t,i}^c - \hat{\mu}_{t,j}^c \pm \beta_{i,j}(t, \delta))$,

$$M^+(i, j; t) := \max_c U_{i,j}^c(t) = \max_c [\hat{\mu}_{t,i}^c - \hat{\mu}_{t,j}^c] + \beta_{i,j}(t, \delta) = M(i, j; t) + \beta_{i,j}(t, \delta),$$

$$M^-(i, j; t) := \max_c L_{i,j}^c(t) = \max_c [\hat{\mu}_{t,i}^c - \hat{\mu}_{t,j}^c] - \beta_{i,j}(t, \delta) = M(i, j; t) - \beta_{i,j}(t, \delta),$$

and the previous lemma implies that on \mathcal{E}_t ,

$$|M(i, j) - M(i, j; t)| \leq \beta_{i,j}(t, \delta) \quad \text{and} \quad |m(i, j) - m(i, j; t)| \leq \beta_{i,j}(t, \delta),$$

which is extensively used in our sample-complexity analysis.

3.5.2 Sample complexity

Lemma 3.3.2. *Let $\varepsilon_1 \geq 0$ and $k \in [K]$. If \mathcal{E}_t holds and $t < \tau_{\varepsilon_1}^k$ then $\omega^k \leq 2\beta_{A_t, A_t}(t, \delta)$.*

Proof. First, note that if $k > |\mathcal{S}^*|$, then the lemma holds trivially since $\omega^k < 0$. In the sequel, we assume \mathcal{E}_t holds and $k \leq |\mathcal{S}^*|$. If $t < \tau_{\varepsilon_1}^k$ then it holds that $|\text{OPT}^{\varepsilon_1}(t)| < k$. So $\mathcal{S}^{*,k} \cap \text{OPT}^{\varepsilon_1}(t)^c \neq \emptyset$. Let $i \in \mathcal{S}^{*,k} \cap \text{OPT}^{\varepsilon_1}(t)^c$, we have

$$\begin{aligned} \omega^k &\leq \omega_i = \min_{j \in [K] \setminus \{i\}} M(i, j), \\ &\stackrel{(a)}{\leq} \min_{j \in [K] \setminus \{i\}} M(i, j; t) + \beta_{i,j}(t, \delta), \\ &\stackrel{(b)}{\leq} \min_{j \in [K] \setminus \{b_t\}} M(b_t, j; t) + \beta_{b_t, j}(t, \delta), \\ &\leq M(b_t, c_t; t) + \beta_{b_t, c_t}(t, \delta), \\ &\stackrel{(c)}{\leq} 2\beta_{b_t, c_t}(t, \delta), \\ &\leq 2\beta_{A_t, A_t}(t, \delta), \end{aligned}$$

where (a) uses that \mathcal{E}_t holds and Lemma 3.2.1, (b) uses the definition of b_t and (c) follows from the definition of c_t and the fact that $b_t \notin \text{OPT}^{\varepsilon_1}(t)$, which yields $M(b_t, c_t; t) \leq \beta_{b_t, c_t}(t, \delta)$. The last inequality follows since A_t is the least sampled among b_t, c_t , and β is decreasing. □

Lemma 3.3.3. *Let $\varepsilon_1 \geq 0$. Let $\tau = \tau_{\varepsilon_1}^k$ for some $k \in [K]$ or $\tau = \tau_{\varepsilon_1, \varepsilon_2}$ for some $\varepsilon_2 \geq 0$. If \mathcal{E}_t holds and $t < \tau$ then $\Delta_{A_t} \leq 2\beta_{A_t, A_t}(t, \delta)$.*

Before proving Lemma 3.3.3, we recall the following lemma, proven in Chapter 2, which is used to derive an upper bound on the gap of an optimal arm.

Lemma 2.5.8. *For any Pareto optimal arm i , $\Delta_i \leq \min_{j \neq i} M(i, j)$.*

Proof of Lemma 3.3.3. Assume that \mathcal{E}_t holds. We consider four different cases depending on whether b_t and c_t are optimal or sub-optimal.

Case 1.1: b_t is a Pareto optimal arm. From the definition of the gap of an optimal arm and using Lemma 2.5.8 it follows $\Delta_{b_t} \leq M(b_t, c_t)$ which on \mathcal{E}_t and using Lemma 3.2.1 yields

$$\Delta_{b_t} + \varepsilon_1 \leq M(b_t, c_t; t) + \beta_{b_t, c_t}(t, \delta) + \varepsilon_1 \quad (3.8)$$

then, noting that there exists $j \in [K] \setminus \{b_t\}$ such that $M(b_t, j, t) + \varepsilon_1 \leq \beta_{b_t, j}(t, \delta)$, by definition of c_t , we have

$$M(b_t, c_t; t) + \varepsilon_1 \leq \beta_{b_t, c_t}(t, \delta), \quad (3.9)$$

therefore,

$$\Delta_{b_t} + \varepsilon_1 \leq 2\beta_{b_t, c_t}(t, \delta).$$

Case 1.2 b_t is a sub-optimal arm. By definition of c_t and using $M = -m$, we have

$$c_t \in \operatorname{argmax}_{j \in [K] \setminus \{b_t\}} m(b_t, j; t) + \beta_{b_t, j}(t, \delta), \quad (3.10)$$

then, from the definition of the gap of a sub-optimal arm and since \mathcal{E}_t holds, we know that there exists an arm b_t^* such that

$$\begin{aligned} \Delta_{b_t} = m(b_t, b_t^*) &\leq m(b_t, b_t^*; t) + \beta_{b_t, b_t^*}(t, \delta), \\ &\stackrel{(a)}{\leq} m(b_t, c_t; t) + \beta_{b_t, c_t}(t, \delta), \\ &\stackrel{(b)}{\leq} 2\beta_{b_t, c_t}(t, \delta). \end{aligned}$$

where (a) uses the definition of c_t and (b) uses Lemma 3.5.1.

Case 2.1 c_t is a Pareto optimal arm. If b_t is also an optimal arm, it follows that $\Delta_{c_t} \leq M(b_t, c_t)$ which on \mathcal{E}_t yields $\Delta_{c_t} \leq M(b_t, c_t; t) + \beta_{b_t, c_t}(t, \delta)$, then, similarly to case 1.1, we have $M(b_t, j; t) + \varepsilon_1 \leq \beta_{b_t, j}(t, \delta)$ so

$$\Delta_{c_t} + \varepsilon_1 \leq 2\beta_{b_t, c_t}(t, \delta).$$

Now, assume b_t is a sub-optimal arm. Then, by definition, $\Delta_{c_t} \leq M(b_t, c_t)_+ + \Delta_{b_t}$. Using a similar reasoning to case 1.2, it holds that $\Delta_{b_t} \leq m(b_t, c_t; t) + \beta_{b_t, c_t}(t, \delta)$, so

$$\begin{aligned}
 \Delta_{c_t} &\leq M(b_t, c_t)_+ + \Delta_{b_t}, \\
 &\leq (M(b_t, c_t; t) + \beta_{b_t, c_t}(t, \delta))_+ + m(b_t, c_t; t) + \beta_{b_t, c_t}(t, \delta), \\
 &= (-m(b_t, c_t; t) + \beta_{b_t, c_t}(t, \delta))_+ + m(b_t, c_t; t) + \beta_{b_t, c_t}(t, \delta), \\
 &\stackrel{(a)}{\leq} \max(2\beta_{b_t, c_t}(t, \delta), m(b_t, c_t; t) + \beta_{b_t, c_t}(t, \delta)) \\
 &\stackrel{(b)}{\leq} 2\beta_{b_t, c_t}(t, \delta).
 \end{aligned}$$

where (a) follows from $(x-y)_+ + (x+y) \leq \max(x+y, 2x)$ and (b) follows from $m(b_t, c_t; t) \leq \beta_{b_t, c_t}(t, \delta)$ (Lemma 3.5.1).

Case 2.2 c_t is a sub-optimal arm. We know that there exists an arm c_t^* such that $\Delta_{c_t} = m(c_t, c_t^*)$. If $c_t^* = b_t$ then, since $m(j, i) \leq M(i, j)$ (follows from the definition), we have

$$\begin{aligned}
 \Delta_{c_t} = m(c_t, c_t^*) &= m(c_t, b_t), \\
 &\leq M(b_t, c_t), \\
 &\stackrel{(a)}{\leq} M(b_t, c_t; t) + \beta_{b_t, c_t}(t, \delta), \\
 &\stackrel{(b)}{\leq} 2\beta_{b_t, c_t}(t, \delta),
 \end{aligned}$$

where (a) follows from \mathcal{E}_t and (b) has been already justified in the case 1.1. If $b_t \neq c_t^*$, then by definition of c_t , we have

$$m(b_t, c_t; t) + \beta_{b_t, c_t}(t, \delta) \geq m(b_t, c_t^*; t) + \beta_{b_t, c_t^*}(t, \delta),$$

which implies that there exists $c \in [d]$ such that

$$\hat{\mu}_{t, c_t}^c - \hat{\mu}_{t, b_t}^c + \beta_{b_t, c_t}(t, \delta) \geq \hat{\mu}_{t, c_t^*}^c - \hat{\mu}_{t, b_t}^c + \beta_{b_t, c_t^*}(t, \delta) \stackrel{\mathcal{E}_t}{\geq} \mu_{c_t^*}^c - \mu_{b_t}^c,$$

then recalling that $\beta_{i, j} = \beta_{j, i}$,

$$\mu_{c_t}^c - \mu_{b_t}^c + 2\beta_{b_t, c_t}(t, \delta) \stackrel{\mathcal{E}_t}{\geq} (\hat{\mu}_{t, c_t}^c - \hat{\mu}_{t, b_t}^c - \beta_{b_t, c_t}(t, \delta)) + 2\beta_{b_t, c_t}(t, \delta) \geq \mu_{c_t^*}^c - \mu_{b_t}^c.$$

Together, these imply that there exists $c \in [d]$ such that

$$\mu_{c_t^*}^c - \mu_{c_t}^c \leq 2\beta_{b_t, c_t}(t, \delta),$$

so

$$\Delta_{c_t} = \min_c [\mu_{c_t^*}^c - \mu_{c_t}^c] \leq 2\beta_{b_t, c_t}(t, \delta).$$

Combining all cases, we have proved that if $t < \min(\tau_{\varepsilon_1}^k, \tau_{\varepsilon_1, \varepsilon_2})$ then both

$$\Delta_{b_t} \leq 2\beta_{b_t, c_t}(t, \delta) \quad \text{and} \quad \Delta_{c_t} \leq 2\beta_{b_t, c_t}(t, \delta) \tag{3.11}$$

holds. Further noting that A_t is the least sampled among b_t, c_t and β is non-increasing, $\beta_{b_t, c_t}(t, \delta) \leq \beta_{A_t, A_t}(t, \delta)$, (3.11) yields

$$\Delta_{A_t} \leq 2\beta_{A_t, A_t}(t, \delta),$$

which completes the proof of Lemma 3.3.3. \square

The following lemma holds for each of the stopping times $\tau_{\varepsilon_1}, \tau_{\varepsilon_1, \varepsilon_2}$ and $\tau_{\varepsilon_1}^k$.

Lemma 3.3.4. *Let $\varepsilon_1 \geq 0$ and τ be as in Lemma 3.3.3. If \mathcal{E}_t holds and $t < \tau$ then $\varepsilon_1 \leq 2\beta_{A_t, A_t}(t, \delta)$.*

Proof. By Lemma 3.5.1, we have $m(b_t, c_t; t) \leq \beta_{b_t, c_t}(t, \delta)$ or equivalently

$$M(b_t, c_t; t) \geq -\beta_{b_t, c_t}(t, \delta). \quad (3.12)$$

Then, knowing that $b_t \notin \text{OPT}^{\varepsilon_1}(t)$, there exists an arm j such that $\varepsilon_1 + M(b_t, j; t) \leq \beta_{b_t, j}(t, \delta)$. Using further the definition of c_t , it follows that $\varepsilon_1 + M(b_t, c_t; t) \leq \beta_{b_t, c_t}(t, \delta)$. Combining this with inequality (3.12) and noting that A_t is the least sampled among b_t, c_t yields

$$\beta_{A_t, A_t}(t, \delta) \geq \beta_{b_t, c_t}(t, \delta) \geq \varepsilon_1/2.$$

\square

Lemma 3.5.1. *Let $\varepsilon_1 \geq 0, \varepsilon_2 \geq 0$ and $k \in [K]$. If $\tau = \tau_{\varepsilon_1}^k$ or $\tau = \tau_{\varepsilon_1, \varepsilon_2}$, the following holds. If $t < \tau$ then for any $j \in [K]$, $m(b_t, j; t) \leq \beta_{b_t, j}(t, \delta)$.*

Proof. The proof is split into two steps.

Step 1: If $t < \tau_{\varepsilon_1}^k$ then for any $j \in [K]$, $m(b_t, j; t) \leq \beta_{b_t, j}(t, \delta)$.

First, note that $t < \min(\tau_{\varepsilon_1}^k, \tau_{\varepsilon_1})$ implies that $Z_1^{\varepsilon_1}(t) \leq 0$ or $Z_2^{\varepsilon_1}(t) \leq 0$. By definition of b_t and noting that $M(i, j; t) = -m(i, j; t)$, we have

$$b_t \in \underset{i \in \text{OPT}^{\varepsilon_1}(t)^c}{\text{argmin}} \max_{j \neq i} m(i, j; t) - \beta_{i, j}(t, \delta). \quad (3.13)$$

so that if there exists j such that $m(b_t, j; t) > \beta_{b_t, j}(t, \delta)$, then

$$\max_{j \neq b_t} m(b_t, j; t) - \beta_{b_t, j}(t, \delta) > 0,$$

therefore,

$$\forall i \in \text{OPT}^{\varepsilon_1}(t)^c, \max_{j \neq i} m(i, j; t) - \beta_{i, j}(t, \delta) > 0 \quad \text{i.e.,} \quad g_i(t) > 0. \quad (3.14)$$

Furthermore, for any $i \in \text{OPT}^{\varepsilon_1}(t)$, $h_i^{\varepsilon_1}(t) > 0$. Combining these, if there exists j such that $m(b_t, j; t) > \beta_{b_t, j}(t, \delta)$ then, $Z_1^{\varepsilon_1}(t) > 0$ and $Z_2^{\varepsilon_1}(t) > 0$.

Step 2: If $t < \tau_{\varepsilon_1, \varepsilon_2}$ then for any $j \in [K]$, $m(b_t, j; t) \leq \beta_{b_t, j}(t, \delta)$.

Recall that by definition $t < \tau_{\varepsilon_1, \varepsilon_2}$ implies that $Z_1^{\varepsilon_1, \varepsilon_2}(t) \leq 0$ or $Z_2^{\varepsilon_1, \varepsilon_2}(t) \leq 0$. Using (3.13), if there exists j such that $m(b_t, j; t) > \beta_{b_t, j}(t, \delta)$, then

$$\max_{j \neq b_t} m(b_t, j; t) - \beta_{b_t, j}(t, \delta) > 0.$$

Combining this with

$$g_i^{\varepsilon_2}(t) := \max_{j \in [K] \setminus \{i\}} [m(i, j; t) - \beta_{i, j}(t, \delta) + \varepsilon_2 \mathbb{1}(j \in \text{OPT}^{\varepsilon_1}(t))],$$

yields

$$\forall i \in \text{OPT}^{\varepsilon_1}(t)^c, 0 < \max_{j \neq i} [m(i, j; t) - \beta_{i, j}(t, \delta)] \leq g_i^{\varepsilon_2}(t). \quad (3.15)$$

Furthermore, since we have

$$\forall i \in \text{OPT}^{\varepsilon_1}(t), h_i^{\varepsilon_1}(t) > 0, \quad (3.16)$$

the initial assumption would yield that for any arm i , $\max\{h_i^{\varepsilon_1}(t), g_i^{\varepsilon_2}(t)\} > 0$, so $Z_1^{\varepsilon_1, \varepsilon_2}(t) > 0$ and $Z_2^{\varepsilon_1, \varepsilon_2}(t) > 0$. We conclude that if $\tau = \tau_{\varepsilon_1}^k$ or $\tau_{\varepsilon_1, \varepsilon_2}$, $t < \tau$ implies that for any $j \in [K]$, $m(b_t, j; t) \leq \beta_{b_t, j}(t, \delta)$. \square

Chapter 4

Linear Model for Pareto Set Identification

This chapter extends the study of *Pareto Set Identification* (PSI) to a structured *multi-output linear bandit* model, where each arm is represented by a feature vector in \mathbb{R}^h and its mean outcome in \mathbb{R}^d depends linearly on this feature via a shared, unknown parameter matrix $\theta \in \mathbb{R}^{h \times d}$.

Building upon the foundations established in Chapters 2 and 3, we introduce, to the best of our knowledge, the first *optimal-design-based algorithms* for structured PSI and provide a unified analysis in both the fixed-budget and fixed-confidence settings. Our results show that the statistical complexity of linear PSI depends only on the h arms with the smallest suboptimality gaps, leading to substantial reductions in sample complexity when the number of arms is large.

This chapter is based on joint work with Émilie Kaufmann and Laura Richert, published in the proceedings of *AISTATS 2025*.

4.1	Introduction	88
4.1.1	Complexity measures for Pareto Set Identification	90
4.1.2	Least-squares estimation and optimal designs	91
4.2	Algorithmic contribution	92
4.2.1	Optimal designs and gap estimation	92
4.2.2	Fixed-budget algorithm	96
4.3	Main theoretical results	96
4.3.1	Fixed-budget	96
4.3.2	Fixed-confidence	97
4.3.3	Sketch of proofs	99
4.4	Numerical study and discussion	100
4.5	Additional proofs	102
4.5.1	Empirical gaps: Proof of Proposition 4.3.6	103
4.5.2	Probability of error: Proof of Theorem 4.3.1	108
4.5.3	Sample complexity: Proof of Theorem 4.3.3	111
4.5.4	Lower bounds	116
4.5.5	Concentration lemmas	118

4.1 Introduction

In previous chapters, we studied *Pareto Set Identification* (PSI) in the unstructured setting (fixed-budget and fixed-confidence), where each arm is a distribution over \mathbb{R}^d and no structure is assumed between the arms. In this chapter, we incorporate additional structure through a *multi-output linear model*, assuming that each arm k is associated with a known feature vector $x_k \in \mathbb{R}^h$ and that its mean reward vector satisfies

$$\mu_k = \theta^\top x_k, \quad \theta \in \mathbb{R}^{h \times d}.$$

This setting, which we refer to as *(multi-output) linear PSI*, generalizes both classical linear bandits ($d = 1$) and multi-objective bandits (with one-hot features, $h = K$). The goal remains to identify the set of arms whose mean vectors are not uniformly dominated by any other, while minimizing the number of samples required to achieve a given confidence level. We bridge the settings of the two previous chapters by extending gap-based PSI algorithms to the linear, multi-objective case.

In many applications, arms exhibit observable structure: for instance, candidate vaccines may be characterized by design features such as antigen dose, adjuvant type, or administration schedule, and their immunogenicity responses can be measured along several endpoints (antibody titers, neutralization, T-cell activity, [Munro et al. 2021](#)). Such features often induce correlations across arms that can be exploited to accelerate learning. Incorporating this structure allows the learner to infer the expected responses of untested candidates and substantially reduces the sample complexity required for identifying promising treatments. The same principle applies to engineering problems—such as algorithm configuration or hardware design—where arms correspond to structured configurations with measurable descriptors ([Almer et al. 2011a](#)).

Related work. The multi-output linear model was first explored by [S. Lu et al. 2019](#) in the context of *Pareto regret minimization*. It can also be seen as a special case of the multi-output kernel regression model of [Zuluaga, Krause, et al. 2016](#) under a linear kernel, for which the authors proposed the ε -PAL algorithm. Parameterized correctly for linear fixed-confidence PSI, ε -PAL has a sample complexity of order $(h^2/\Delta_{\min}^2) \log(1/\delta)$.

[Kim et al. 2023](#) extended the algorithm of [Auer et al. 2016](#) using a robust estimator to jointly minimize Pareto regret and identify the Pareto set, achieving a bound of order $(h/\Delta_{\min}^2) \log(1/\delta)$.

When $d = 1$, linear PSI reduces to the classical *linear bandit* pure exploration problem, whose sample complexity and optimal design connections have been extensively studied ([Soare et al. 2014](#); [Tao et al. 2018](#); [Degenne, Ménard, et al. 2020](#)). [Soare et al. 2014](#) first explored the link with optimal experimental design ([Pukelsheim 2006](#)), showing that the minimal sample complexity can be expressed as an optimal (XY) design. They proposed the first algorithms with sample complexity scaling in $(h/\Delta_{\min}^2) \log(1/\delta)$ where Δ_{\min} is the smallest gap in the model. [Tao et al. 2018](#) proposed a novel estimator of the regression parameter and a G-optimal-design-based algorithm, with a sample complexity in $\sum_{i=1}^h \Delta_{(i)}^{-2} \log(1/\delta)$

where $\Delta_{(1)} \leq \dots \leq \Delta_{(h)}$ are the h smallest gaps, which improves upon the complexity of the unstructured setting when $K \gg h$. Some algorithms even match the minimal sample complexity either in the asymptotic regime $\delta \rightarrow 0$ (Degenne, Ménard, et al. 2020; Jedra & Proutiere 2020b) or within multiplicative factors (Fiez et al. 2019). Some adaptive algorithms, such as LinGapE (Xu et al. 2018), are also very effective in practice, but without provably improving over unstructured algorithms in all instances.

The fixed-budget setting has been studied by Azizi et al. 2022; Yang & Tan 2022, who propose algorithms based on Sequential Halving (Karnin et al. 2013) where in each round, the active arms are sampled according to a G-optimal design. The best guarantees are those obtained by Yang & Tan 2022 who show that a budget T of order $\log_2(h) \sum_{i=1}^h \Delta_{(i)}^{-2} \log(1/\delta)$ is sufficient to ensure an error probability at most δ .

Katz-Samuels, Jain, et al. 2020 propose an elimination algorithm for both fixed-confidence and fixed-budget settings, using optimal design to minimize the *Gaussian width*—a complexity measure that may better capture non-asymptotic error rates. Extending such notions to linear PSI is challenging due to the complex structure of the set of alternative models with a different Pareto set.

Contributions. We formalize the *linear Pareto Set Identification* (linear PSI) problem as an extension of PSI to structured arms with features in \mathbb{R}^h , allowing the learner to exploit shared information across arms through their linear structure. We propose GEGER (G-optimal Empirical Gap Estimation), the first PSI algorithm based on G-optimal experimental design, applicable to both fixed-confidence and fixed-budget regimes. At each round, arms are sampled according to a G-optimal design computed over the feature space, ensuring maximal information gain for estimating the relevant gaps between arms.

Our unified analysis establishes instance-dependent sample complexity bounds showing that the hardness of linear PSI scales with the h smallest suboptimality gaps: replacing the unstructured dependence on all K arms with $\sum_{i=1}^h \Delta_{(i)}^{-2}$ instead of K/Δ_{\min}^2 . This yields a substantial reduction in sample complexity compared to the unstructured algorithms of Chapters 2 and 3. In the fixed-confidence setting, GEGER achieves identification within $\mathcal{O}(\log(1/\Delta_{(1)}))$ adaptive rounds, while in the fixed-budget setting, it is, to our knowledge, the first algorithm for multi-output linear bandits with near-optimal performance guarantees.

Learning model. We formalize the linear PSI problem. Let $d, h \in \mathbb{N}^*$ with $h \leq K$. Each arm $k \in [K]$ is associated with a known feature vector $x_k \in \mathbb{R}^h$ and an unknown mean reward vector $\mu_k \in \mathbb{R}^d$ satisfying

$$\mu_k = \theta^\top x_k, \quad \theta \in \mathbb{R}^{h \times d},$$

where θ is an unknown parameter matrix. Let $\mathcal{X} := (x_1, \dots, x_K)^\top$ and $[K] := \{1, \dots, K\}$.

Definition 4.1.1 (Pareto dominance). For any two arms $i, j \in [K]$, arm i is *weakly dominated* by j if $\mu_i^c \leq \mu_j^c$ for all $c \in [d]$. It is *dominated* ($i \preceq j$) if it is weakly dominated and strictly smaller in at least one coordinate, and *strictly dominated* ($i \prec j$) if all inequalities are strict.

The (strict) *Pareto set* is then defined as

$$\mathcal{S}^* = \{i \in [K] : \nexists j \in [K] \setminus \{i\}, \mu_i \prec \mu_j\}.$$

At each round t , the learner selects an arm $A_t \in [K]$ and observes

$$y_t = \theta^\top x_{A_t} + \eta_t,$$

where η_t is centered and, conditionally on A_t , has σ -subgaussian marginals¹. The goal is to identify \mathcal{S}^* as efficiently as possible.

In the *fixed-confidence* setting, given $\delta \in (0, 1)$, the learner collects samples up to a stopping time τ and outputs $\widehat{\mathcal{S}}_\tau$ such that $\mathbb{P}(\widehat{\mathcal{S}}_\tau \neq \mathcal{S}^*) \leq \delta$ while minimizing τ . In the *fixed-budget* setting, sampling stops after T rounds, and the objective is to minimize the error probability $e_T := \mathbb{P}(\widehat{\mathcal{S}}_T \neq \mathcal{S}^*)$.

Throughout the chapter, we use the following notation. Δ_n denotes the probability simplex in \mathbb{R}^n . For a positive semidefinite matrix $A \in \mathbb{R}^{n \times n}$ and a vector $x \in \mathbb{R}^n$, we write $\|x\|_A^2 := x^\top A x$, and $x(i)$ denotes the i -th coordinate of x . For scalars $a, b \in \mathbb{R}$, $a \wedge b := \min\{a, b\}$ and $(a)_+ := \max\{a, 0\}$.

4.1.1 Complexity measures for Pareto Set Identification

As discussed in Chapters 2 and 3, the difficulty of Pareto Set Identification (PSI) in the unstructured setting can be quantified through suitable notions of *suboptimality gaps*, which measure how confidently one can classify each arm as dominated or non-dominated. These gaps, first introduced by [Auer et al. 2016](#) and further refined in Chapter 2, rely on the quantities

$$m(i, j) := \min_{c \in [d]} [\mu_j^c - \mu_i^c], \quad M(i, j) := -m(i, j),$$

which respectively capture how much arm j dominates arm i , or how much j must be shifted component-wise to dominate i . In particular, $m(i, j) > 0$ if and only if $i \prec j$, while $M(i, j) > 0$ if and only if $i \not\prec j$.

For each arm i , the *suboptimality gap* Δ_i (see Lemma 1 in [Kone, Kaufmann, et al. 2024](#)) is defined as

$$\Delta_i := \begin{cases} \Delta_i^*, & \text{if } i \notin \mathcal{S}^*, \\ \min_{j \neq i} [M(i, j) \wedge (M(j, i)_+ + (\Delta_j^*)_+)], & \text{otherwise,} \end{cases} \quad (4.1)$$

with

$$\Delta_i^* := \max_{j \in [K]} m(i, j).$$

Intuitively, for a dominated arm i , Δ_i represents the minimal shift required to make it non-dominated, while for an optimal arm i , it measures its closest separation from both other optimal and suboptimal arms.

¹A centered random variable X is σ -subgaussian if $\log \mathbb{E}[\exp(\lambda X)] \leq \lambda^2 \sigma^2 / 2$ for all $\lambda \in \mathbb{R}$.

In the sequel, we assume without loss of generality that $\Delta_1 \leq \Delta_2 \leq \dots \leq \Delta_K$. Following Chapters 2 and 3, we recall the global complexity measures

$$H_1(\nu) := \sum_{i=1}^K \Delta_i^{-2} \quad \text{and} \quad H_2(\nu) := \max_{i \in [K]} \frac{i}{\Delta_i^2},$$

which govern the sample complexity in the fixed-confidence and fixed-budget settings, respectively.

Structured complexity under linearity. In the linear setting, the learner estimates the low-dimensional parameter $\theta \in \mathbb{R}^{h \times d}$ rather than the K individual means, leading to a substantial reduction in complexity. We therefore define the analogous quantities

$$H_{1,\text{lin}}(\nu) := \sum_{i=1}^h \frac{1}{\Delta_i^2}, \quad H_{2,\text{lin}}(\nu) := \max_{i \in [h]} \frac{i}{\Delta_i^2}, \quad (4.2)$$

which characterize the hardness of the linear PSI problem in the fixed-confidence and fixed-budget settings, respectively. Since these depend only on the h smallest gaps, $H_{1,\text{lin}}(\nu) \leq H_1(\nu)$ and $H_{2,\text{lin}}(\nu) \leq H_2(\nu)$, illustrating the gain achieved by leveraging linear structure.

4.1.2 Least-squares estimation and optimal designs

Given n arm selections a_1, \dots, a_n , define $X_n := (x_{a_1} \dots x_{a_n})^\top \in \mathbb{R}^{n \times h}$ and $Y_n := (y_1 \dots y_n)^\top \in \mathbb{R}^{n \times d}$. We define the information matrix as $V_n := X_n^\top X_n = \sum_{i=1}^K N_n(i) x_i x_i^\top \in \mathbb{R}^{h \times h}$ where $N_{n,i}$ denotes the number of observations from arm i among the n samples. More generally, given $w \in \mathbb{R}_+^K$, we define $V(w) := \sum_{i=1}^K w(i) x_i x_i^\top$.

The multi-output regression model can be written in matrix form as $Y_n = X_n \theta + H_n$ where $H_n := (\eta_1 \dots \eta_n)^\top$ is the noise matrix. The least-squares estimate $\hat{\theta}_n$ of the matrix θ is defined as the matrix minimizing the least-squares error $\text{Err}_n(A) := \|X_n A - Y_n\|_F^2$. Computing the gradient of the loss yields $V_n \hat{\theta}_n = X_n^\top Y_n$. If the matrix V_n is nonsingular, the least-squares estimator can be written

$$\hat{\theta}_n = V_n^{-1} X_n^\top Y_n.$$

In the course of our elimination algorithm, we will compute least-squares estimates based on observations from a restricted number of arms, and we will face the case in which V_n is singular. In this case, different choices have been made in prior work on linear bandits: [Alieva et al. 2021](#) define a custom “pseudo-inverse” while [Yang & Tan 2022](#) define new contexts \tilde{x}_i that are projections of the x_i onto a subspace of dimension $\text{rank}(\mathcal{X}_S)$ where $\mathcal{X}_S := (x_i : i \in S)^\top$ and S is the set of arms that are active. We adopt an approach close to the latter, which is described below. Let the singular value decomposition of $(\mathcal{X}_S)^\top$ be USV^\top where U, V are orthogonal matrices and $B := (u_1, \dots, u_m)$ is formed with the first m columns of U where $m = \text{rank}(\mathcal{X}_S)$. We then define

$$V_n^\dagger := B(B^\top V_n B)^{-1} B^\top \quad \text{and} \quad \hat{\theta}_n = V_n^\dagger X_n^\top Y_n. \quad (4.3)$$

The following result addresses the statistical uncertainty of this estimator.

Lemma 4.1.2. *If $\Sigma \in \mathbb{R}^{d \times d}$ is the covariance matrix of the noise η_t and a_1, \dots, a_n are deterministically chosen, then for any $x_i \in \{x_{a_1}, \dots, x_{a_n}\}$, $\text{Cov}(\hat{\theta}_n^\top x_i) = \|x_i\|_{V_n^\dagger}^2 \Sigma$.*

Hence, achieving uniform estimation accuracy across arms amounts to selecting $\{a_1, \dots, a_n\}$ so as to minimize $\max_{i \in S} \|x_i\|_{V_n^\dagger}^2$. Relaxing this combinatorial problem leads to the *G-optimal design*:

$$w_S^* \in \underset{w \in \Delta_{|S|}}{\operatorname{argmin}} \max_{i \in S} \|\tilde{x}_i\|_{\tilde{V}(w)^{-1}}^2, \quad (4.4)$$

where $\tilde{x}_i = B^\top x_i$ are the projected features and $\tilde{V}(w) = \sum_{i \in S} w(i) \tilde{x}_i \tilde{x}_i^\top$, with the convention that $w(i)$ indexes arms in S . Intuitively, w_S^* specifies a sampling distribution over the active arms that yields the most uniform estimation accuracy of their mean responses under the projected model (4.3). A procedure for computing w_S^* is provided in Appendix H of [Kone, Kaufmann, et al. 2025a](#).

4.2 Algorithmic contribution

Our elimination algorithms operate in rounds. They progressively eliminate a portion of arms and classify them as optimal or suboptimal based on empirical estimation of their gaps. In each round, a sampling budget is allocated among the surviving arms based on a G-optimal design.

4.2.1 Optimal designs and gap estimation

At round r , we denote by \mathcal{A}_r the set of arms that are still active. To estimate the means and, henceforth, the gaps, we first compute an estimate of the regression matrix denoted $\hat{\theta}_r$. This estimate is obtained by carefully sampling the arms using an integral rounding of a G-optimal design.

Algorithm 4.1 takes as input a set of arms S , a budget N and chooses some N arms to pull (with repetitions) based on an integer rounding of w_S^* , a continuous G-optimal design over the set $\{\tilde{x}_i, i \in S\}$ of (transformed) features associated with the arms (see Section 4.1.2). Several rounding procedures have been proposed in the literature, and we use that of [Allen-Zhu et al. 2017](#), henceforth referred to as ROUND, which satisfies the following guarantees, as proven in Lemma 12 of [Kone, Kaufmann, et al. 2025a](#).

Lemma 4.2.1 (Lemma 12 of [Kone, Kaufmann, et al. 2025a](#)). *ROUND($N, \tilde{\mathcal{X}}_S, w_S^*, \kappa$) outputs a sequence of arms $(a_1, \dots, a_N) \in S^N$ such that*

$$\max_{i \in S} \|x_i\|_{V_N^\dagger}^2 \leq (1 + 6\kappa) \frac{F_S(w_S^*)}{N},$$

where $F_S(w_S^*)$ is the optimal value of (4.4).

Algorithm 4.1: OPTESTIMATOR: Least-squares estimation from G-optimal design

- Require:** Subset $S \subset [K]$, sample size N , precision κ
- // Feature transformation (Section 4.1.2)
- 1 Compute the transformed features $\tilde{\mathcal{X}}_S = (B^\top x_i, i \in S)$ with B as defined in Section 4.1.2 // G-optimal design over transformed features and sampling
 - 2 Compute a G-optimal design w_S^* over the set $\tilde{\mathcal{X}}_S$
 - 3 Pull $(a_1, \dots, a_N) \leftarrow \text{ROUND}(N, \tilde{\mathcal{X}}_S, w_S^*, \kappa)$ and collect responses y_1, \dots, y_N
 - 4 Compute V_N^\dagger as in Eq. (4.3) and compute the OLS estimator on the samples collected

$$\hat{\theta} \leftarrow V_N^\dagger \sum_{t=1}^N x_{a_t} y_t^\top$$

return: $\hat{\theta}$

By the Kiefer–Wolfowitz theorem (Kiefer & Wolfowitz 1960), $F_S(w_S^*) = h_S$, the dimension of $\text{span}\{x_i : i \in S\}$. This observation is key to the following concentration result.

Lemma 4.2.2. *Let $S \subset [K]$, $\kappa \in (0, 1/3]$ and $N \geq 5h_S/\kappa^2$ where $h_S = \dim(\text{span}\{x_i : i \in S\})$. The output $\hat{\theta}$ of OPTESTIMATOR(S, N, κ) satisfies for all $\varepsilon > 0$ and $i \in S$*

$$\mathbb{P}\left(\|(\theta - \hat{\theta})^\top x_i\|_\infty \geq \varepsilon\right) \leq 2d \exp\left(-\frac{N\varepsilon^2}{2(1+6\kappa)\sigma^2 h_S}\right).$$

Proof of Lemma 4.2.2. By assumption, the noise vector has σ -subgaussian marginals. It is easy (see e.g., proof of Lemma 11 in Kone, Kaufmann, et al. 2025a) to see that for any $i \in S$, the marginals of $(\theta - \hat{\theta})x_i$ are $\sigma\|X_N^\top V_N^\dagger x_i\|_2$ -subgaussian. Then, direct algebra yields

$$\begin{aligned} \|X_N^\top V_N^\dagger x_i\|_2^2 &= x_i^\top V_N^\dagger V_N V_N^\dagger x_i \\ &= x_i^\top (B_S (B_S^\top V_N B_S)^{-1} B_S^\top) V_N (B_S (B_S^\top V_N B_S)^{-1} B_S^\top) x_i \\ &= x_i^\top B_S (B_S^\top V_N B_S)^{-1} B_S^\top x_i \\ &= x_i^\top V_N^\dagger x_i = \|x_i\|_{V_N^\dagger}^2. \end{aligned}$$

Next, by concentration of subgaussian variables (see e.g., Lattimore & Szepesvari 2020)

$$\begin{aligned} \mathbb{P}(\|(\theta - \hat{\theta})^\top x_i\|_\infty \geq \varepsilon) &\leq 2d \exp\left(-\frac{\varepsilon^2}{2\sigma^2 \|x_i\|_{V_N^\dagger}^2}\right) \\ &\leq 2d \exp\left(-\frac{\varepsilon^2}{2\sigma^2 \max_{k \in S} \|x_k\|_{V_N^\dagger}^2}\right), \end{aligned}$$

and the result follows from Lemma 4.2.1. □

Empirical gaps. Once the parameter $\hat{\theta}_r$ has been obtained as an output of Algorithm 4.1 with $S = \mathcal{A}_r$ and an appropriate value of the budget N , we compute estimates of the mean vectors as $\hat{\mu}_{i,r} := \hat{\theta}_r^\top x_i$ and the empirical Pareto set of active arms,

$$S_r := \{i \in \mathcal{A}_r : \nexists j \in \mathcal{A}_r : \hat{\mu}_{i,r} \prec \hat{\mu}_{j,r}\}.$$

In both the fixed-confidence and fixed-budget settings, at round r , after collecting new samples from the surviving arms, GEGE discards a fraction of the arms based on the empirical estimation of their gaps. We first introduce the empirical quantities used to compute the gaps:

$$\begin{aligned} M(i, j; r) &:= \max_{c \in [d]} [\hat{\mu}_{i,r}^c - \hat{\mu}_{j,r}^c], & \text{and} \\ m(i, j; r) &:= \min_{c \in [d]} [\hat{\mu}_{j,r}^c - \hat{\mu}_{i,r}^c]. \end{aligned}$$

We define for any arm $i \in \mathcal{A}_r$,

$$\hat{\Delta}_{i,r}^* := \max_{j \in \mathcal{A}_r} m(i, j; r),$$

and the empirical estimates of the PSI gaps as:

$$\hat{\Delta}_{i,r} := \begin{cases} \hat{\Delta}_{i,r}^* & \text{if } i \in \mathcal{A}_r \setminus S_r \\ \hat{\delta}_{i,r}^* & \text{if } i \in S_r \end{cases} \quad (4.5)$$

with $\hat{\delta}_{i,r}^* := \min_{j \in \mathcal{A}_r \setminus \{i\}} [M(i, j; r) \wedge (M(j, i; r)_+ + (\hat{\Delta}_{j,r}^*)_+)]$; the empirical estimates of the gaps introduced earlier in Section 4.1.1 and already used in Chapter 2.

Unlike BAI, which only requires rejecting suboptimal arms, PSI must explicitly classify each arm as either Pareto-optimal or dominated, since the cardinality of the Pareto set is unknown. These classification mechanisms are detailed in the following sections for the fixed-budget and fixed-confidence settings.

Final output. In both cases, letting \mathcal{A}_r be the set of active arms and \mathcal{B}_r be the set of arms already classified as optimal at the beginning of round r , GEGE outputs $\mathcal{B}_{\tau+1} \cup \mathcal{A}_{\tau+1}$ as the candidate Pareto-optimal set, where τ denotes the final round. Moreover, $\mathcal{A}_{\tau+1}$ contains at most one arm.

GEGE for ε -PSI. Algorithm 4.2 can be easily modified to identify an ε -Pareto set. As introduced in Auer et al. 2016 and studied in Chapter 3, an ε -Pareto set S_ε is such that $S^* \subset S_\varepsilon$ and for any arm $i \in S_\varepsilon$, $\Delta_i^* \leq \varepsilon$: it contains the Pareto set and possibly some suboptimal arms that are (ε) -close to be optimal. To identify an ε -Pareto set, we relax the stopping condition: instead of stopping when it remains only one active arm (*i.e.*, $|\mathcal{A}_r| \leq 1$), we stop when $(|\mathcal{A}_r| \leq 1 \text{ or } \varepsilon_r \leq \varepsilon/4)$ holds. After stopping, the same set is recommended, namely $\mathcal{A}_{\tau+1} \cup \mathcal{B}_{\tau+1}$. The guarantees of GEGE under this modification are discussed in Section E.5 of Kone, Kaufmann, et al. 2025a.

Algorithm 4.2: GEGE: G-optimal Empirical Gap Elimination (fixed-confidence)

Require : risk parameter $\delta \in (0, 1)$; noise proxy σ^2

- 1 Initialize: set of active arms $\mathcal{A}_1 \leftarrow [K]$; set of arms (so far) optimal $\mathcal{B}_1 \leftarrow \emptyset$; set of discarded arms $\mathcal{D}_1 \leftarrow \emptyset$; $r \leftarrow 1$
- 2 **while** $|\mathcal{A}_r| > 1$ **do**
 - // Schedule, confidence split, and local dimension*
 - 3 $\varepsilon_r \leftarrow \frac{1}{2 \cdot 2^r}$; $\delta_r \leftarrow \frac{6\delta}{\pi^2 r^2}$; $h_r \leftarrow \dim(\text{span}\{x_i : i \in \mathcal{A}_r\})$
 - // Round budget via G-optimal design complexity*
 - 4 $t_r \leftarrow \left\lceil \frac{32(1 + 3\varepsilon_r) \sigma^2 h_r \log\left(\frac{|\mathcal{A}_r| d}{2 \delta_r}\right)}{\varepsilon_r^2} \right\rceil$
 - // Compute estimate from near G-optimal allocation over \mathcal{A}_r*
 - 5 $\hat{\theta}_r \leftarrow \text{OPTESTIMATOR}(\mathcal{A}_r, t_r, \varepsilon_r)$
 - 6 Compute S_r and empirical gaps $\hat{\Delta}_{i,r}$ via (4.5)
 - // Accept strong Pareto-optimal arms; discard strong suboptimal arms*
 - 7 $\mathcal{B}_{r+1} \leftarrow \mathcal{B}_r \cup \{i \in S_r : \hat{\Delta}_{i,r} \geq \varepsilon_r\}$
 - 8 $\mathcal{D}_{r+1} \leftarrow \mathcal{D}_r \cup \{i \in \mathcal{A}_r \setminus S_r : \hat{\Delta}_{i,r} \geq \varepsilon_r/2\}$
 - // Shrink active set and advance round*
 - 9 $\mathcal{A}_{r+1} \leftarrow \mathcal{A}_r \setminus (\mathcal{D}_{r+1} \cup \mathcal{B}_{r+1})$
 - 10 $r \leftarrow r + 1$
- 11 **return** $\mathcal{B}_r \cup \mathcal{A}_r$

4.2.2 Fixed-budget algorithm

Algorithm 4.3, operates over $\lceil \log_2(h) \rceil$ rounds, with an equal budget of $T/\lceil \log_2(h) \rceil$ allocated per round. By construction $|\mathcal{A}_{\lceil \log_2(h) \rceil+1}| = 1$. At the end of round r , the $\lceil h/2^r \rceil$ arms with the smallest empirical gaps are kept active while the remaining arms are discarded and classified as Pareto optimal (added to \mathcal{B}_{r+1}) if they are empirically optimal (belonging to set S_r) and deemed suboptimal otherwise. If a tie occurs, we break it in favor of empirically Pareto-optimal arms (*i.e.*, keep arms in S_r when possible). This is crucial to prove the guarantees on the algorithm, as sketched in Section 4.3.

Algorithm 4.3: GEGER: G-optimal Empirical Gap Elimination (fixed-budget)

Require: Total budget T ; feature dimension h

- 1 Initialize: $R \leftarrow \lceil \log_2(h) \rceil$; $t_r \leftarrow \left\lfloor \frac{T}{R} \right\rfloor$: per-round allocation size (equal split across rounds); set of active arms $\mathcal{A}_1 \leftarrow [K]$; $\mathcal{B}_1 \leftarrow \emptyset$; $\mathcal{D}_1 \leftarrow \emptyset$;
- 2 **for** $r = 1$ **to** $\lceil \log_2(h) \rceil$ **do**
 - 3 // G-optimal design over current active set
 $\hat{\theta}_r \leftarrow \text{OPTESTIMATOR}(\mathcal{A}_r, t_r, 1/3)$;
// Execute allocation, update estimates, form Pareto set and gaps
 - 4 Compute S_r (empirical Pareto set) and gaps $\hat{\Delta}_{i,r}$ via (4.5);
// Keep the $\lceil h/2^r \rceil$ smallest-gap arms for next round
 - 5 $m_r \leftarrow \left\lfloor \frac{h}{2^r} \right\rfloor$; $\mathcal{A}_{r+1} \leftarrow$ the m_r arms in \mathcal{A}_r with smallest $\hat{\Delta}_{i,r}$;
// Tie-breaks favor arms in S_r
 - 6 break ties by keeping $i \in S_r$ when possible;
 - 7 $\mathcal{B}_{r+1} \leftarrow \mathcal{B}_r \cup (S_r \cap (\mathcal{A}_r \setminus \mathcal{A}_{r+1}))$;
 - 8 $\mathcal{D}_{r+1} \leftarrow \mathcal{D}_r \cup ((\mathcal{A}_r \setminus \mathcal{A}_{r+1}) \setminus S_r)$;
- 9 $\mathcal{B}_{R+1} \cup \mathcal{A}_{R+1}$

4.3 Main theoretical results

We now turn to the theoretical analysis of the proposed algorithms. In this section, we establish finite-sample guarantees for both the fixed-budget and fixed-confidence versions of GEGER, beginning with the fixed-budget setting.

4.3.1 Fixed-budget

In the fixed-budget setting, Algorithm 4.3 satisfies the following guarantees.

Theorem 4.3.1. *The probability of error of Algorithm 4.3 run with budget $T \geq 45h \log_2 h$ is at most*

$$\exp\left(-\frac{T}{1200\sigma^2 H_{2,\text{lin}} \lceil \log_2 h \rceil} + \log C(h, d, K)\right),$$

where $C(h, d, K) = 2d(K + h + \lceil \log_2 h \rceil)$.

To the best of our knowledge, GEGE is the first algorithm with theoretical guarantees for fixed-budget linear PSI. Our result shows that in this setting, the probability of error scales only with the first h gaps. In Chapter 2, we introduced EGE-SH, an algorithm for fixed-budget PSI in the unstructured setting whose probability of error is essentially upper-bounded by

$$\exp\left(-\frac{T}{288\sigma^2 H_2 \log_2 K} + \log(2d(K-1)|\mathcal{S}^*| \log_2 K)\right).$$

Therefore, GEGE largely improves upon EGE-SH when $K \gg h$.

When $K = h$ and x_1, \dots, x_K are the canonical \mathbb{R}^h -basis, both algorithms coincide², thus, GEGE can be seen as a generalization of EGE-SH.

The following lower bound for linear PSI in the fixed-budget setting shows that GEGE is optimal in the worst case, up to constants and a $\log_2(h)$ factor.

Theorem 4.3.2. *Let \mathbb{W}_H be the set of instances with complexity $H_{2,\text{lin}}$ smaller than H . For any budget T , letting $\widehat{S}_T^{\text{alg}}$ be the output of an algorithm alg, it holds that*

$$\inf_{\text{alg}} \sup_{\nu \in \mathbb{W}_H} \mathbb{P}_{\nu}(\widehat{S}_T^{\text{alg}} \neq \mathcal{S}^*(\nu)) \geq \frac{1}{4} \exp\left(-\frac{2T}{H\sigma^2}\right).$$

4.3.2 Fixed-confidence

We now turn to the fixed-confidence regime, where the learner adaptively collects samples until it can identify the Pareto set with confidence $1 - \delta$.

Theorem 4.3.3. *With probability at least $1 - \delta$, Algorithm 4.2 identifies the Pareto set using at most*

$$\log_2(2/\Delta_1) + \mathcal{O}\left(\sum_{i=2}^h \frac{\sigma^2}{\Delta_i^2} \log\left(\frac{Kd}{\delta} \log_2\left(\frac{2}{\Delta_i}\right)\right)\right)$$

samples and $\lceil \log_2(1/\Delta_1) \rceil$ rounds.

This result shows that the complexity of Algorithm 4.2 scales only with the first h gaps. In particular, when $K \gg h$ using our algorithm substantially reduces the sample complexity of PSI. In Table 4.1, we compare the sample complexity of GEGE to that of existing fixed-confidence PSI algorithms, showing that GEGE enjoys stronger guarantees than its competitors. We emphasize that both Kim et al. 2023 and Zuluaga, Krause, et al.

Table 4.1: Sample complexity up to constant multiplicative terms of different algorithms for PSI in the fixed-confidence setting.

Algorithm	Upper-Bound on τ_δ	Linear PSI
Zuluaga, Krause, et al. 2016	$\left(\frac{h^2}{\Delta_{\min}^2}\right) \log^3\left(\frac{Kd}{\delta}\right)$	✓
Kone, Kaufmann, et al. 2023	$\sum_{i=1}^K \frac{1}{\Delta_i^2} \log\left(\frac{Kd}{\delta} \log\left(\frac{1}{\Delta_i}\right)\right)$	✗
Kim et al. 2023	$\frac{h}{\Delta_{\min}^2} \log\left(\frac{d(h\sqrt{K})}{\delta\Delta_{\min}^2}\right)$	✓
GEGE (Ours)	$\sum_{i=1}^h \frac{1}{\Delta_i^2} \log\left(\frac{Kd}{\delta} \log\left(\frac{1}{\Delta_i}\right)\right)$	✓

2016 use uniform sampling and do not exploit an optimal design, which prevents them from reaching the guarantees given in Theorem 4.3.3.

The result below upper bounds the sample complexity by leveraging properties of the optimal design.

Theorem 4.3.4. *With probability at least $1 - \delta$, the sample complexity of Algorithm 4.2 is upper bounded by*

$$\mathcal{O}\left(\log \frac{1}{\Delta_1} \inf_{w \in \Delta_K} \max_{k \in [K]} \frac{2\sigma^2 \|x_k\|_{V(w)}^2}{\Delta_k^2} \log \frac{Kd}{\Delta_1 \delta}\right).$$

This provides an alternative characterization of the sample complexity, though it is not directly comparable to Theorem 4.3.3 in general.

We state a lower bound, showing that our algorithm is essentially minimax optimal for linear PSI.

Theorem 4.3.5. *For any $K, d, h \in \mathbb{N}$, there exists a set $\mathcal{M}(K, d, h)$ of linear PSI instances such that for $\nu \in \mathcal{M}(K, d, h)$ and for any δ -correct algorithm for linear PSI, with probability at least $1 - \delta$,*

$$\tau_\delta^A \geq \Omega\left(H_{1, \text{lin}}(\nu) \log \frac{1}{\delta}\right).$$

Remark 4.1. When $K = h$ and x_1, \dots, x_K form the canonical \mathbb{R}^h basis, we recover the classical PSI problem. We note that, unlike its fixed-budget version, GEGE does not coincide with an existing PSI identification algorithm. Instead, it matches the optimal guarantees of Kone, Kaufmann, et al. 2023 while needing only $\lceil \log(1/\Delta_1) \rceil$ rounds of adaptivity, which is the first fixed-confidence PSI algorithm having this property. Such a batched algorithm may be desirable in some applications, e.g., in clinical trials where measuring different biological indicators of efficacy can take time.

²Up to the fact that EGE-SH does not discard samples between rounds

4.3.3 Sketch of proofs

Before sketching our proof strategy, we highlight a key property of PSI that makes the analysis different from classical BAI settings. Let a be a (Pareto) suboptimal arm. From (4.1), there exists $a^* \in \mathcal{S}^*$ such that $\Delta_a = m(a, a^*)$ and importantly, a^* could be the unique arm dominating a . Therefore, discarding a^* before a may result in the latter appearing as optimal in the remaining rounds, thus leading to the misidentification of the Pareto set.

To avoid this, an elimination algorithm for PSI should guarantee that if a suboptimal arm a is active, then a^* is also active. We introduce the following event

$$\mathcal{P}_r := \{\forall s \leq r : \forall i \in (\mathcal{S}^*)^c, i \in \mathcal{A}_s \Rightarrow i^* \in \mathcal{A}_s\}.$$

An important aspect of our proofs is to control the occurrence of \mathcal{P}_∞ (by convention, if \mathcal{P}_t holds and $\mathcal{A}_s = \emptyset$ for any $s \geq t$, then \mathcal{P}_∞ holds). The first step of the proof is to show that when \mathcal{P}_r holds, we can control the deviations of the empirical gaps, which is essential to guarantee the correctness of GEGE and to control its sample complexity in the fixed-confidence setting. We now define for $\eta > 0$, the good event

$$\mathcal{E}^r(\eta) := \left\{ \forall i, j \in \mathcal{A}_r : \left\| (\hat{\theta}_r - \theta)^\top (x_i - x_j) \right\|_\infty \leq \eta \right\}. \quad (4.6)$$

Letting $n_r = |\mathcal{A}_r|$ and λ a constant to be specified, we introduce

$$\mathcal{E}_{\text{fb}}^\lambda := \bigcap_{r=1}^{\lceil \log_2(h) \rceil} \mathcal{E}^r(\lambda \Delta_{n_{r+1}+1}), \quad \text{and} \quad \mathcal{E}_{\text{fc}} := \bigcap_{r=1}^{\infty} \mathcal{E}^r(\varepsilon_r/2).$$

We then prove by concentration and induction the following key result.

Proposition 4.3.6. *Let $\lambda \in (0, 1/5)$ and assume \mathcal{E}_{fc} (resp. $\mathcal{E}_{\text{fb}}^\lambda$ in fixed-budget) holds. Then at any round r , \mathcal{P}_r holds and for all arm $i \in \mathcal{A}_r$,*

$$\widehat{\Delta}_{i,r} - \Delta_i \geq \begin{cases} -\eta_r & \text{if } i \in \mathcal{S}^* \\ -\eta_r/2 & \text{else,} \end{cases} \quad \text{where} \quad \eta_r := \begin{cases} 2\lambda \Delta_{n_{r+1}+1} & \text{(fixed-budget)} \\ \varepsilon_r & \text{(fixed-confidence).} \end{cases}$$

Building on this result, we show that the recommendation of Algorithm 4.3 is correct on $\mathcal{E}_{\text{fb}}^\lambda$, so its probability of error is upper-bounded by $\inf_{\lambda \in (0, 1/5)} \mathbb{P}(\mathcal{E}_{\text{fb}}^\lambda)$. We conclude the proof of Theorem 4.3.1 by upper-bounding this probability (see Section 4.5).

Similarly, using Proposition 4.3.6 we prove the correctness of Algorithm 4.2 on \mathcal{E}_{fc} : at any round r , $\mathcal{B}_r \subset \mathcal{S}^*$ and $\mathcal{D}_r \subset (\mathcal{S}^*)^c$. To further upper bound its sample complexity, we need an additional result to control the size of \mathcal{A}_r .

Lemma 4.3.7. *The following statement holds for Algorithm 4.2 on the event \mathcal{E}_{fc} : for all $p \in [K]$, after $\lceil \log(1/\Delta_p) \rceil$ rounds it remains less than p active arms. In particular, GEGE stops after at most $\lceil \log(1/\Delta_1) \rceil$ rounds.*

The proof of this lemma is given in section 4.5.3. To get the sample complexity bound of Theorem 4.3.3, some extra arguments are needed. We sketch some elements below (the

full proof is given in section 4.5). Assume \mathcal{E}_{fc} holds and let τ_δ be the sample complexity of Algorithm 4.2. Lemma 4.3.7 yields $\tau_\delta \leq \sum_{r=1}^{\lceil \log(1/\Delta_1) \rceil} \mathcal{O}(h_r/\varepsilon_r^2)$ with $h_r \leq |\mathcal{A}_r|$.

Using Lemma 4.3.7, we introduce "checkpoints rounds" between which we control $|\mathcal{A}_r|$ and thus h_r . Let the sequence $(\alpha_s)_{s \geq 0}$ defined as $\alpha_0 = 0$ and $\alpha_s = \lceil \log_2(1/\Delta_{\lfloor h/2^s \rfloor}) \rceil$, for $s \geq 1$. Simple calculation yields $\alpha_{\lfloor \log_2(h) \rfloor} = \lceil \log_2(1/\Delta_1) \rceil$ and $\{1, \dots, \lceil \log_2(1/\Delta_1) \rceil\} = \cup_{s=1}^{\lfloor \log_2(h) \rfloor} \llbracket 1 + \alpha_{s-1}, \alpha_s \rrbracket$. Therefore

$$\tau_\delta \leq \sum_{s=1}^{\lfloor \log_2(h) \rfloor} \sum_{r=\alpha_{s-1}+1}^{\alpha_s} \mathcal{O}(|\mathcal{A}_r|/\varepsilon_r^2).$$

Now by Lemma 4.3.7, for $r > \alpha_s$, $|\mathcal{A}_r| \leq \lfloor h/2^s \rfloor$, so essentially $\tau_\delta \leq \sum_{s=1}^{\lfloor \log_2(h) \rfloor} \mathcal{O}(4^{\alpha_s} \lfloor h/2^s \rfloor)$. Carefully re-indexing this sum and addressing a few more technicalities, we obtain the result in Theorem 4.3.3. Showing that $\mathbb{P}(\mathcal{E}_{\text{fc}}) \geq 1 - \delta$, from Lemma 4.2.2 completes the proof of Theorem 4.3.3 (see Section 4.5)

On the other side, combining the Kiefer–Wolfowitz theorem

$$\begin{aligned} \tau_\delta &\leq \sum_{r=1}^{\lceil \log(1/\Delta_1) \rceil} \mathcal{O}\left(\frac{\sigma^2 h_r}{\varepsilon_r^2} \log \frac{K d 2^r}{\delta}\right) \\ &\leq \sum_{r=1}^{\lceil \log(1/\Delta_1) \rceil} \mathcal{O}\left(\inf_{w \in \Delta_K} \max_{k \in \mathcal{A}_r} \|x_k\|_{V(w)^{-1}}^2 \frac{\sigma^2}{\varepsilon_r^2} \log \frac{K d 2^r}{\delta}\right) \quad (\text{Kiefer–Wolfowitz}). \end{aligned}$$

Next, note that if $k \in \mathcal{A}_r$ then $\Delta_k \leq 2\varepsilon_{r-1}$ (otherwise it would have been discarded at the end of round $r-1$). Therefore, $1/\varepsilon_r^2 \leq 4/\Delta_k^2$ for any arm $k \in \mathcal{A}_r$. Putting these displays together yields

$$\begin{aligned} \tau_\delta &\leq \sum_{r=1}^{\lceil \log(1/\Delta_1) \rceil} \mathcal{O}\left(\sigma^2 \inf_{w \in \Delta_K} \max_{k \in \mathcal{A}_r} \frac{\|x_k\|_{V(w)^{-1}}^2}{\Delta_k^2} \log \frac{K d 2^r}{\delta}\right) \\ &\leq \mathcal{O}\left(\log \frac{1}{\Delta_1} \inf_{w \in \Delta_K} \max_{k \in [K]} \frac{2\sigma^2 \|x_k\|_{V(w)^{-1}}^2}{\Delta_k^2} \log \frac{K d}{\Delta_1 \delta}\right). \end{aligned}$$

4.4 Numerical study and discussion

We now empirically evaluate the proposed algorithms on both synthetic and real-world instances. Our experiments aim to validate the theoretical guarantees derived in Section 4.3, and to illustrate the statistical and computational advantages of exploiting the linear structure in Pareto set identification.

Experimental setup. In the fixed-budget setting, we compare GEGER-FB with two unstructured PSI algorithms from Chapter 2—EGE-SH and EGE-SR—as well as a uniform sampling

baseline. In the fixed-confidence setting, we benchmark GEGER-FC against APE, the fully adaptive algorithm introduced in Chapter 3, and PAL (Zuluaga, Sergent, et al. 2013), a Gaussian-process-based method instantiated here with a linear kernel. For the synthetic instances, we fix feature vectors x_1, \dots, x_h and add additional arms x_{h+1}, \dots, x_K . A common regression matrix θ is shared across all instances. Given $K \geq h$, we construct a linear PSI instance ν_K by augmenting the base feature set with the additional arms x_{h+1}, \dots, x_K , chosen such that the first h arms have identical minimal gaps. Consequently, the complexity terms $H_{1,\text{lin}}$ and $H_{2,\text{lin}}$ remain constant as K grows. Unless stated otherwise, we set $h = 8$ and $d = 2$.

For the real-world dataset, we evaluate the algorithms on the NoC dataset (Almer et al. 2011b), a bi-objective hardware design benchmark involving $d = 2$ performance metrics—energy consumption and runtime—across $K = 259$ implementations, each described by $h = 4$ features for hardware implementation (chip area, bandwidth, etc.)

For each algorithm, we report (i) the empirical probability of misidentification in the fixed-budget case, and (ii) the empirical distribution of the sample complexity in the fixed-confidence case, both averaged over 500 independent runs. We fix $\delta = 0.01$ for fixed-confidence experiments and $T = \lceil H_{2,\text{lin}} \rceil$ for fixed-budget ones.

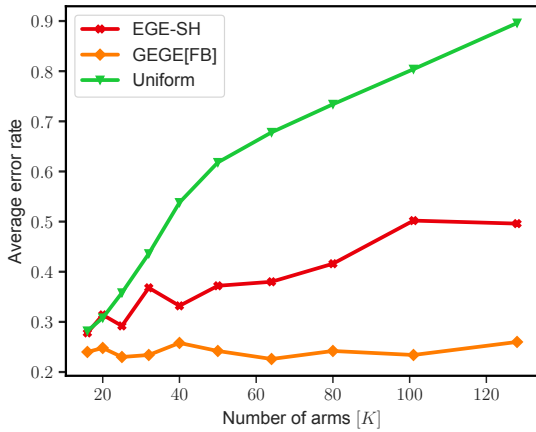


Figure 4.1: Average misidentification rate vs. K (synthetic data).

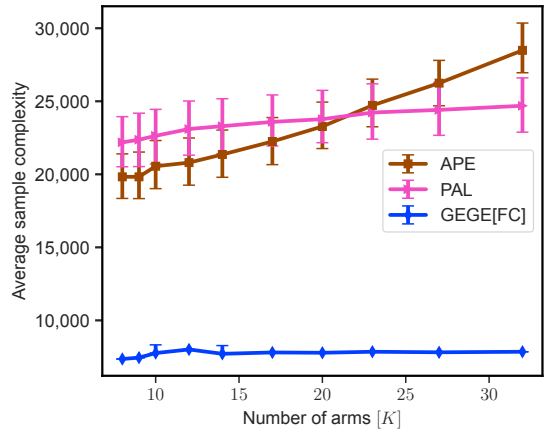


Figure 4.2: Average sample complexity vs. K (synthetic data).

Summary. Theoretical results (Theorems 4.3.1 and 4.3.3) predict that the sample complexity of GEGER depends on the number of arms K only through logarithmic factors. This trend is confirmed empirically: in Figure 4.1, the misidentification rate of GEGER-FB remains almost constant as K increases, whereas that of EGE-SH/SR rises sharply. Similarly, in Figure 4.2, the sample complexity of GEGER-FC shows only mild dependence on K , unlike the unstructured baseline.

On the real-world NoC experiment, GEGER consistently outperforms all competitors. In the fixed-confidence regime (Figure 4.4), it requires substantially fewer samples to achieve δ -correct identification than both APE and PAL. In the fixed-budget regime (Figure 4.3),

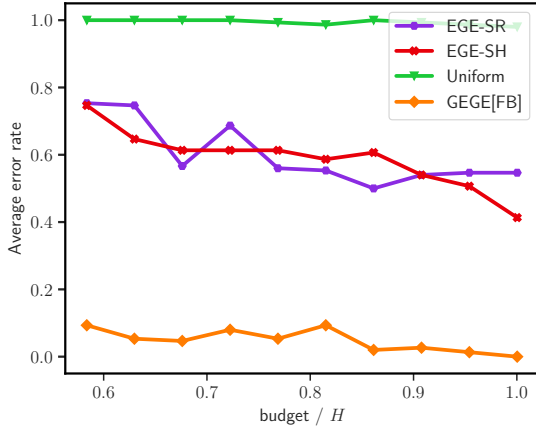


Figure 4.3: Average misidentification rate vs. T (NoC experiment).

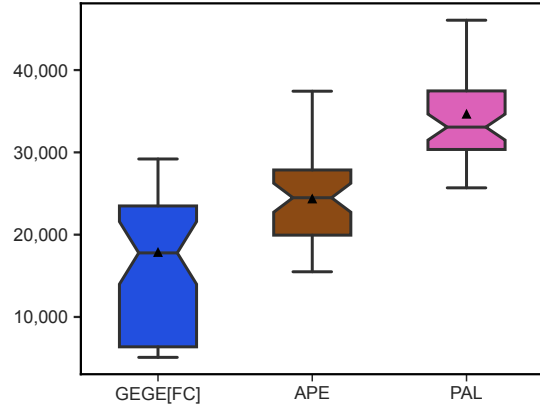


Figure 4.4: Empirical sample complexity (NoC experiment).

the probability of misidentification is reduced by up to 50% compared to EGE-SH. Notably, EGE-SH requires $T \geq K \log_2 K \approx 2000$ to operate on this instance, while GEGE-FB succeeds with $T \geq \lceil \log_2 h \rceil$. Each experiment runs in under 10 seconds for instances with up to $K = 500$ and $d = 8$.

Discussion. Overall, GEGE provides the first near-optimal algorithms for Pareto set identification in multi-output linear bandits. By leveraging optimal design principles, it efficiently estimates both mean responses and dominance gaps, yielding strong theoretical guarantees and consistent empirical performance across settings. In the fixed-budget regime, GEGE-FB achieves near-minimax guarantees; in the fixed-confidence regime, GEGE-FC remains the only algorithm with provably limited adaptivity.

While the complexity terms of GEGE depend only on the h smallest gaps, logarithmic dependencies in K remain due to standard union bounds. Following the empirical-process techniques of [Katz-Samuels, Jain, et al. 2020](#), future work could explore whether these logarithmic factors can be eliminated. Finally, extending the optimal design approach to high-dimensional or kernelized feature spaces—along the lines of [Camilleri, Jamieson, et al. 2021](#)—would be a promising direction for PSI in rich function classes.

4.5 Additional proofs

In the subsequent sections, r will always denote a round of GEGE, which should be clear from the context. We then denote by \mathcal{A}_r active arms at round r and by $\hat{\theta}_r$ the empirical estimate of θ at round r , computed by a call to Algorithm 4.1. By convention we let $\max_{\emptyset} = -\infty$. For any suboptimal arm i , there exists a Pareto-optimal arm i^* (not necessarily unique) such that $\Delta_i = m(i, i^*)$. More generally given a suboptimal i we denote by i^* any arm of $\operatorname{argmax}_{j \in \mathcal{S}^*} m(i, j)$. At a round r , we let

$$\mathcal{P}_r := \{\forall s \in \{1, \dots, r\}, \forall i \in \mathcal{A}_s, i \in (\mathcal{S}^*)^c \Rightarrow i^* \in \mathcal{A}_s\}, \quad (4.7)$$

with $\mathcal{P} = \mathcal{P}_\infty$.

In particular, for a suboptimal arm i , $i^* \in \mathcal{A}_s$ should be understood as

$$\mathcal{A}_s \cap (\operatorname{argmax}_{j \in \mathcal{S}^*} m(i, j)) \neq \emptyset$$

. If for some τ , \mathcal{P}_τ is true and $\mathcal{A}_{\tau+1} = \emptyset$ then we say that \mathcal{P} holds.

4.5.1 Empirical gaps: Proof of Proposition 4.3.6

In both the fixed-budget and fixed-confidence setting, the proof proceeds by induction on the round r .

To establish this first result, we need the following intermediate lemmas, proved at the end of the section.

Lemma 4.5.1. *At any round r and for any arms $i, j \in \mathcal{A}_r$ it holds that*

$$\begin{aligned} |M(i, j; r) - M(i, j)| &\leq \|(\hat{\theta}_r - \theta)^\top(x_i - x_j)\|_\infty \text{ and} \\ |m(i, j; r) - m(i, j)| &\leq \|(\hat{\theta}_r - \theta)^\top(x_i - x_j)\|_\infty. \end{aligned}$$

Lemma 4.5.2. *At any round r , for any suboptimal arm $i \in \mathcal{A}_r$, if $i^* \in \mathcal{A}_r$ and i^* does not empirically dominate i then $\Delta_i^* < \|(\hat{\theta}_r - \theta)^\top(x_i - x_{i^*})\|_\infty$.*

Deviations of the gaps when \mathcal{P}_r holds. In this part, we control the deviations of the empirical gaps when proposition \mathcal{P}_r holds.

Lemma 4.5.3. *Assume that the proposition \mathcal{P}_r holds at some round r . Then for any arm $i \in \mathcal{A}_r$ it holds that*

$$\left| (\widehat{\Delta}_{i,r}^*)_+ - (\Delta_i^*)_+ \right| \leq \left| \widehat{\Delta}_{i,r}^* - \Delta_i^* \right| \leq \gamma_{i,r}$$

where $\gamma_{i,r} := \max_{j \in \mathcal{A}_r} \|(\hat{\theta}_r - \theta)^\top(x_i - x_j)\|_\infty$.

Proof. This inequality is a direct consequence of Lemma 4.5.1 and the relation $|x_+ - y_+| \leq |x - y|$ which holds for any $x, y \in \mathbb{R}$. Note that for a Pareto-optimal arm i we trivially have $(\Delta_i^*)_+ = 0 = (\max_{j \in \mathcal{A}_r} m(i, j))_+$. And for a suboptimal arm $i \in \mathcal{A}_r$, as $i^* \in \mathcal{A}_r$ (from proposition \mathcal{P}_r) we have $\Delta_i^* = m(i, i^*) = \max_{j \in \mathcal{A}_r} m(i, j)$. Thus for any arm $i \in \mathcal{A}_r$ we have

$$\begin{aligned} \left| (\widehat{\Delta}_{i,r}^*)_+ - (\Delta_i^*)_+ \right| &= \left| \left(\max_{j \in \mathcal{A}_r} m(i, j; r) \right)_+ - \left(\max_{j \in \mathcal{A}_r} m(i, j) \right)_+ \right|, \\ &\leq \left| \max_{j \in \mathcal{A}_r} m(i, j; r) - \max_{j \in \mathcal{A}_r} m(i, j) \right|, \\ &\leq \max_{j \in \mathcal{A}_r} |m(i, j; r) - m(i, j)|, \\ &\leq \max_{j \in \mathcal{A}_r} \left\| (\hat{\theta}_r - \theta)^\top(x_i - x_j) \right\|_\infty = \gamma_{i,r}, \end{aligned}$$

where the last inequality follows from Lemma 4.5.1. \square

Lemma 4.5.4. *If the proposition \mathcal{P}_r holds at some round r then for any arm $i \in \mathcal{A}_r$,*

$$\widehat{\Delta}_{i,r} - \Delta_i \geq \begin{cases} -2\gamma_r & \text{if } i \in \mathcal{S}^*, \\ -\gamma_{i,r} & \text{else,} \end{cases}$$

where $\gamma_{i,r} := \max_{j \in \mathcal{A}_r} \|(\widehat{\theta}_r - \theta)^\top (x_i - x_j)\|_\infty$ and $\gamma_r := \max_{i \in \mathcal{A}_r} \gamma_{i,r}$.

Proof. We first prove the result a suboptimal arm i . From the proposition \mathcal{P}_r , as $i \in \mathcal{A}_r$ we have $i^* \in \mathcal{A}_r$ so $\Delta_i = \max_{j \in \mathcal{A}_r} m(i, j)$ and we recall that

$$\widehat{\Delta}_{i,r} := \max(\widehat{\Delta}_{i,r}^*, \widehat{\delta}_{i,r}^*). \quad (4.8)$$

Note that by reverse triangle we have for any arm $i \in \mathcal{A}_r$ (suboptimal or not)

$$\begin{aligned} \left| \left(\max_{j \in \mathcal{A}_r} m(i, j; r) \right) - \left(\max_{j \in \mathcal{A}_r} m(i, j) \right) \right| &\leq \max_{j \in \mathcal{A}_r} |m(i, j; r) - m(i, j)|, & (4.9) \\ &\leq \max_{j \in \mathcal{A}_r} \left\| (\widehat{\theta}_r - \theta)^\top (x_i - x_j) \right\|_\infty = \gamma_{i,r}. & (4.10) \end{aligned}$$

where the last inequality follows from Lemma 4.5.1. If i a suboptimal arm ($i \notin \mathcal{S}^*$) then as $\Delta_i = \Delta_i^*$, it follows

$$\widehat{\Delta}_{i,r} - \Delta_i \geq \widehat{\Delta}_{i,r}^* - \Delta_i^*$$

therefore

$$\begin{aligned} \widehat{\Delta}_{i,r} - \Delta_i &\geq -|\widehat{\Delta}_{i,r}^* - \Delta_i^*| \\ &= -\left| \left(\max_{j \in \mathcal{A}_r} m(i, j; r) \right) - \left(\max_{j \in \mathcal{A}_r} m(i, j) \right) \right| \\ &\geq -\gamma_{i,r} \quad (\text{see (4.10)}). \end{aligned}$$

Now we assume i is a Pareto-optimal arm ($i \in \mathcal{S}^*$) so that

$$\Delta_i = \delta_i^*.$$

Combining with Eq. (4.8) yields

$$\widehat{\Delta}_{i,r} - \Delta_{i,r} \geq \widehat{\delta}_{i,r}^* - \delta_i^*,$$

where we recall that

$$\widehat{\delta}_{i,r}^* = \min_{j \in \mathcal{A}_r \setminus \{i\}} [M(i, j; r) \wedge (M(j, i; r)_+ + (\widehat{\Delta}_{j,r}^*)_+)]$$

and

$$\delta_i^* := \min_{j \in [K] \setminus \{i\}} [M(i, j) \wedge (M(j, i)_+ + (\Delta_j^*)_+)].$$

As for any $x, y \in \mathbb{R}$ we have $|x^+ - y^+| \leq |x - y|$, the following holds for any $i, j \in \mathcal{A}_r$

$$|M(j, i; r)^+ - M(j, i)^+| \leq |M(j, i; r) - M(j, i)| \quad (4.11)$$

$$\leq \gamma_{j,r}. \quad (4.12)$$

From Lemma 4.5.3 we have for any $j \in \mathcal{A}_r$

$$(\widehat{\Delta}_{j,r}^*)_+ - (\Delta_j^*)_+ \geq -\gamma_{j,r}. \quad (4.13)$$

Combining (4.12) and (4.13) yields for any $j \in \mathcal{A}_r$

$$M(j, i; r)_+ + (\widehat{\Delta}_{j,r}^*)_+ \geq M(j, i)_+ + (\Delta_j^*)_+ - 2\gamma_{j,r}, \quad (4.14)$$

which in addition to $M(j, i; r) \geq M(j, i) - \gamma_{j,r}$ yields

$$[M(i, j; r) \wedge (M(j, i; r)_+ + (\widehat{\Delta}_{j,r}^*)_+)] \geq [M(i, j) \wedge (M(j, i)_+ + (\Delta_j^*)_+)] - 2\gamma_{j,r},$$

for any arm $j \in \mathcal{A}_r$. Thus taking the min over $i \in \mathcal{A}_r$ yields

$$\begin{aligned} \widehat{\delta}_{i,r}^* &= \min_{j \in \mathcal{A}_r \setminus \{i\}} [M(i, j; r) \wedge (M(j, i; r)_+ + (\widehat{\Delta}_{j,r}^*)_+)] \\ &\geq \min_{j \in \mathcal{A}_r \setminus \{i\}} [M(i, j) \wedge (M(j, i)_+ + (\Delta_j^*)_+)] - 2\gamma_r, \\ &\geq \min_{j \in [K] \setminus \{i\}} [M(i, j) \wedge (M(j, i)_+ + (\Delta_j^*)_+)] - 2\gamma_r, \\ &= \delta_i^* - 2\gamma_r \end{aligned}$$

which concludes the proof of this lemma. \square

Building on this result, we show that \mathcal{P}_∞ holds in the fixed-confidence and fixed-budget settings.

Fixed-budget setting. We recall the definition of the good event for any $\lambda > 0$.

$$\mathcal{E}_{\text{fb}}^{r,\lambda} = \left\{ \forall i, j \in \mathcal{A}_r : \|(\widehat{\theta}_r - \theta)^\top (x_i - x_j)\|_\infty \leq \lambda \Delta_{n_{r+1}+1} \right\}$$

and $\mathcal{E}_{\text{fb}}^\lambda := \bigcap_{r=1}^{\lceil \log_2(h) \rceil} \mathcal{E}_{\text{fb}}^{r,\lambda}$. We prove that proposition \mathcal{P}_∞ holds on the event $\mathcal{E}_{\text{fb}}^\lambda$ for some any $\lambda \in (0, 1/5)$.

Lemma 4.5.5. *The proposition holds \mathcal{P}_∞ on the event $\mathcal{E}_{\text{fb}}^\lambda$ for any $\lambda \in (0, 1/5)$: at any round $r \in \{1, \dots, \lceil \log_2 h \rceil + 1\}$ and for any arm $i \in \mathcal{A}_r \cap (\mathcal{S}^*)^c, i^* \in \mathcal{A}_r$.*

Proof. We prove \mathcal{P}_∞ by induction on the round r . In the sequel we assume $\mathcal{E}_{\text{fb}}^\lambda$ holds. We also assume \mathcal{P}_r is true until some round r . As $\mathcal{E}_{\text{fb}}^\lambda$ holds, we have by application of Lemma 4.5.4: for any arm $i \in \mathcal{A}_r$,

$$\widehat{\Delta}_{i,r} - \Delta_i \geq \begin{cases} -2\lambda \Delta_{n_{r+1}+1} & \text{if } i \in \mathcal{S}^* \\ -\lambda \Delta_{n_{r+1}+1} & \text{else.} \end{cases} \quad (4.15)$$

We shall prove that if a Pareto-optimal arm i is discarded at the end of round r then there exists no arm suboptimal $j \in \mathcal{A}_{r+1}$ such that $j^* = i$. Since i is removed and $|\mathcal{A}_{r+1}| = n_{r+1}$ there exists $k_r \in \mathcal{A}_{r+1} \cup \{i\}$ such that

$$\Delta_{k_r} \geq \Delta_{n_{r+1}+1}. \quad (4.16)$$

If i is empirically suboptimal then as it is discarded we have

$$\widehat{\Delta}_{i,r} = \widehat{\Delta}_{i,r}^* \geq \widehat{\Delta}_{k,r}$$

for any arm $k \in \mathcal{A}_{r+1}$. So $\widehat{\Delta}_{i,r}^* \geq \widehat{\Delta}_{k,r}$ thus using (4.15) and (4.16) it comes that

$$\begin{aligned} \max_{q \in \mathcal{A}_r \setminus \{i\}} m(i, q) &\geq \Delta_{n_{r+1}+1} - 3\lambda \Delta_{n_{r+1}+1} \\ &= (1 - 3\lambda) \Delta_{n_{r+1}+1} \end{aligned}$$

and the latter inequality is not possible for $\lambda < 1/3$ as the LHS of the inequality is negative as i is a Pareto-optimal arm.

Next we assume that i is empirically optimal. We claim that j is not dominated by i . To see this, first note that as $j \in \mathcal{A}_{r+1}$ we have

$$\widehat{\Delta}_{i,r} \geq \widehat{\Delta}_{j,r} \quad (4.17)$$

so that as i is empirically optimal, if j was empirically dominated by i we would have

$$\widehat{\Delta}_{i,r} \leq M(j, i; r)_+ + (\widehat{\Delta}_{j,r}^*)_+ = \widehat{\Delta}_{j,r}. \quad (4.18)$$

Combining (4.17) and (4.18) yields $\widehat{\Delta}_{i,r} = \widehat{\Delta}_{j,r}$, i is empirically optimal and j is empirically suboptimal. However, our breaking rule ensures that this case cannot occur. Therefore, j is not dominated by i . But, by assumption, j is such that $j^* = i$ and we have proved that i does not empirically dominate j so by Lemma 4.5.2

$$\Delta_j \leq \|(\hat{\theta}_r - \theta)^\top (x_j - x_i)\|_\infty$$

which on the event \mathcal{E}_{fb} yields

$$\Delta_j \leq \lambda \Delta_{n_{r+1}+1}. \quad (4.19)$$

On the other side, as i is discarded as an empirically optimal arm we have

$$\widehat{\Delta}_{i,r} = \widehat{\delta}_{i,r}^* \geq \widehat{\Delta}_{k,r}$$

for any arm $k \in \mathcal{A}_{r+1}$. Since $k_r \in \mathcal{A}_{r+1} \cup \{i\}$ it comes $\widehat{\delta}_{i,r}^* \geq \widehat{\Delta}_{k_r,r}$ thus using (4.15) and (4.16) yields

$$M(j, i)_+ + \Delta_j \geq \Delta_{n_{r+1}+1} - 4\lambda \Delta_{n_{r+1}+1}$$

which further combined with (4.19) yields

$$M(j, i)_+ \geq (1 - 5\lambda)\Delta_{n_{r+1}+1}.$$

However, as $j^* = i$ we have $M(j, i)_+ = 0$ so the latter inequality is not possible as long as $\lambda < 1/5$. Combining these, we have proved that if \mathcal{P}_r holds, then for any Pareto-optimal arm i which is removed at the end of round r , there does not exist an arm $j \in \mathcal{A}_{r+1}$ such that $j^* = i$. So \mathcal{P}_{r+1} holds. Finally noting that \mathcal{P}_r trivially holds for $r = 1$ we conclude that \mathcal{P}_∞ holds on the event $\mathcal{E}_{\text{fb}}^\lambda$ for any $\lambda < 1/5$. \square

Combining this result with Lemma 4.5.4 and assuming $\mathcal{E}_{\text{fb}}^\lambda$ holds then yields at any round $r \in \{1, \dots, \lceil \log_2 h \rceil\}$ and for any arm $i \in \mathcal{A}_r$:

$$\widehat{\Delta}_{i,r} - \Delta_i \geq \begin{cases} -2\lambda\Delta_{n_{r+1}+1} & \text{if } i \in \mathcal{S}^* \\ -\lambda\Delta_{n_{r+1}+1} & \text{else,} \end{cases} \quad (4.20)$$

which proves Proposition 4.3.6 in the fixed-budget setting.

Fixed-confidence setting. We recall below the good events we study in the fixed-confidence setting:

$$\mathcal{E}_{\text{fc}}^r = \left\{ \forall i, j \in \mathcal{A}_r : \|(\hat{\theta}_r - \theta)^\top(x_i - x_j)\|_\infty \leq \varepsilon_r/2 \right\}$$

and $\mathcal{E}_{\text{fc}} := \bigcap_{r=1}^\infty \mathcal{E}_{\text{fc}}^r$.

Lemma 4.5.6. *The proposition \mathcal{P}_∞ holds on the event \mathcal{E}_{fc} : at any round r for any arm $i \in \mathcal{A}_r \cap (\mathcal{S}^*)^c$, $i^* \in \mathcal{A}_r$.*

Proof. We prove the proposition by induction on the round r . Note that the proposition \mathcal{P}_r trivially holds for $r = 1$. Assume the property holds until the beginning of some round r . Let $i \in \mathcal{S}^*$ be an optimal arm and assume i is discarded at the end of round r . We will prove that there exists no suboptimal arm $j \in \mathcal{A}_{r+1}$ such that $j^* = i$. Recall that when i is discarded, we have either $i \in S_r$ (empirically optimal) or $i \notin S_r$ (empirically suboptimal). We analyze both cases below. If $i \notin S_r$ then it holds that

$$\widehat{\Delta}_{i,r} \geq \varepsilon_r/2,$$

then, as $i \notin S_r$ it follows that $\widehat{\Delta}_{i,r} = \widehat{\Delta}_i^* := \max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j; r)$, so

$$\max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j; r) \geq \varepsilon_r/2$$

which using Lemma 4.5.1 and assuming event $\mathcal{E}_{\text{fc}}^r$ holds would yield

$$\max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j) > 0.$$

The latter inequality is not possible as $i \in \mathcal{S}^*$ is a Pareto-optimal arm. Therefore, on $\mathcal{E}_{\text{fc}}^r$, when $i \in \mathcal{S}^*$ is discarded we have $i \in S_r$.

Next, we analyze the case $i \in S_r$: that is i is discarded and classified as optimal. In this case it follows from the definition of $\widehat{\Delta}_{i,r}$ that

$$\min_{j \in \mathcal{A}_r \setminus \{i\}} [M(j, i; r)_+ + (\widehat{\Delta}_{j,r}^*)_+] \geq \varepsilon_r. \quad (4.21)$$

Let $j \in \mathcal{A}_{r+1} \cap (\mathcal{S}^*)^c$ be such that $j^* = i$. If j is empirically optimal then $(\widehat{\Delta}_{j,r}^*)_+ = 0$ thus $M(j, i; r)_+ \geq \varepsilon_r$. On the contrary, if j is empirically suboptimal, then because it has not been removed at the end of round r it holds that

$$\widehat{\Delta}_{j,r}^* < \varepsilon_r/2,$$

which combined with (4.21) yields $M(j, i; r)_+ > \varepsilon_r/2$. Thus, in both cases we have $M(j, i; r)_+ > \varepsilon_r/2$ which using Lemma 4.5.1 and assuming event $\mathcal{E}_{\text{fc}}^r$ would imply that

$$M(j, i)_+ > 0,$$

which is impossible as, by assumption $j^* = i$, so j is dominated by i .

Together, these imply that on \mathcal{E}_{fc} , if \mathcal{P}_r holds then \mathcal{P}_{r+1} holds. Since the property trivially holds for $r = 1$ we have proved that the property \mathcal{P}_r holds at any round when \mathcal{E}_{fc} holds. \square

Combining this result with Lemma 4.5.4 proves that, on the event \mathcal{E}_{fc} , for any round r and for any arm $i \in \mathcal{A}_r$

$$\widehat{\Delta}_{i,r} - \Delta_i \geq \begin{cases} -\varepsilon_r & \text{if } i \in \mathcal{S}^* \\ -\varepsilon_r/2 & \text{else,} \end{cases} \quad (4.22)$$

which proves Proposition 4.3.6 in the fixed-confidence setting.

4.5.2 Probability of error: Proof of Theorem 4.3.1

Proof of Theorem 4.3.1. We first prove the correctness of GEGE on the event $\mathcal{E}_{\text{fb}}^\lambda$ for some λ small enough. Let us assume $\mathcal{E}_{\text{fb}}^\lambda$ holds which by Proposition 4.3.6 implies that \mathcal{P}_∞ holds and at round r , we have for any arm $i \in \mathcal{A}_r$

$$\widehat{\Delta}_{i,r} - \Delta_i \geq \begin{cases} -2\lambda\Delta_{n_{r+1}+1} & \text{if } i \in \mathcal{S}^* \\ -\lambda\Delta_{n_{r+1}+1} & \text{else.} \end{cases} \quad (4.23)$$

We recall the definition of the good event for any $\lambda > 0$,

$$\mathcal{E}_{\text{fb}}^{r,\lambda} = \left\{ \forall i, j \in \mathcal{A}_r : \|(\widehat{\theta}_r - \theta)^\top(x_i - x_j)\|_\infty \leq \lambda\Delta_{n_{r+1}+1} \right\}$$

and $\mathcal{E}_{\text{fb}}^\lambda := \bigcap_{r=1}^{\lceil \log_2(h) \rceil} \mathcal{E}_{\text{fb}}^{r,\lambda}$. Applying Lemma 4.5.1 on this event then yields for all arms $i, j \in \mathcal{A}_r$,

$$|M(i, j; r) - M(i, j)| \leq \lambda \Delta_{n_{r+1}+1} \text{ and} \quad (4.24)$$

$$|m(i, j; r) - m(i, j)| \leq \lambda \Delta_{n_{r+1}+1}. \quad (4.25)$$

Let i be an arm discarded at the end of round r . Since i is discarded and $|\mathcal{A}_{r+1}| = n_{r+1}$ there exists $k_r \in \mathcal{A}_{r+1} \cup \{i\}$ such that

$$\Delta_{k_r} \geq \Delta_{n_{r+1}+1}. \quad (4.26)$$

If $i \notin S_r$ that is i is empirically suboptimal then

$$\widehat{\Delta}_{i,r} = \widehat{\Delta}_{i,r}^* \geq \widehat{\Delta}_{k_r,r},$$

then, recalling that

$$\widehat{\Delta}_{i,r}^* := \max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j; r)$$

and further applying (4.23) to k_r and using (4.25) yields

$$\max_{j \in \mathcal{A}_r \setminus \{i\}} m(i, j) \geq (1 - 3\lambda) \Delta_{n_{r+1}+1}$$

which for $\lambda < 1/3$ implies that $\max_{j \in \mathcal{A}_r} m(i, j) > 0$, that is there exists $j \in \mathcal{A}_r$ such that $\mu_i < \mu_j$ so i is a suboptimal arm. Next, assume $i \in S_r$ (i.e., i is empirically Pareto-optimal). In this case we have $\widehat{\Delta}_{i,r} = \widehat{\delta}_{i,r}^* \geq \widehat{\Delta}_{k_r,r}$.

We recall that

$$\widehat{\delta}_{i,r}^* = \min_{j \in \mathcal{A}_r \setminus \{i\}} [M(i, j; r) \wedge (M(j, i; r)_+ + (\widehat{\Delta}_{i,r}^*)_+)].$$

Applying (4.23) to k_r and using (4.24), it follows that

$$\min_{j \in \mathcal{A}_r \setminus \{i\}} M(i, j) \geq (1 - 3\lambda) \Delta_{n_{r+1}+1}.$$

Thus, for $\lambda < 1/3$, we have $\min_{j \in \mathcal{A}_r \setminus \{i\}} M(i, j) > 0$. Therefore, no active arm at round r dominates i (based on their true means), which, together with proposition \mathcal{P}_∞ , yields that i is a Pareto-optimal arm (otherwise, we would have $i^* \in \mathcal{A}_r$ that dominates i).

Combining these, we have proved that for any $\lambda < 1/5$ (we need $\lambda < 1/5$ for \mathcal{P}_∞ to hold), Algorithm 4.3 does not make any error on the event $\mathcal{E}_{\text{fb}}^\lambda$. It then follows that the probability of error of GEGE is at most

$$\inf_{\lambda \in (0, 1/5)} \mathbb{P}((\mathcal{E}_{\text{fb}}^\lambda)^c) \quad (4.27)$$

Now we upper-bound Eq. (4.27), which will conclude the proof. Let $\lambda \in (0, 1/5)$ be fixed. We have by union bound

$$\begin{aligned} \mathbb{P}((\mathcal{E}_{\text{fb}}^\lambda)^c) &\leq \sum_{r=1}^{\lceil \log_2 h \rceil} \mathbb{E} \left[\mathbb{P}((\mathcal{E}_{\text{fb}}^{r,\lambda})^c | \mathcal{A}_r) \right] \\ &\leq \sum_{r=1}^{\lceil \log_2 h \rceil} \mathbb{E} \left[\sum_{i \in \mathcal{A}_r} \mathbb{P}(\|(\hat{\theta}_r - \theta)^\top x_i\|_\infty > \frac{1}{2} \lambda \Delta_{n_{r+1}+1} | \mathcal{A}_r) \right] \end{aligned}$$

Note that for i fixed, we can use Lemma 4.2.2 with $\kappa = 1/3$ and the conditions of this theorem are satisfied as the budget per phase is $T/\log_2(h) \geq 45h$ (recall from the theorem that GEGE is run with $T \geq 45h \log_2(h)$). Thus, applying this theorem yields

$$\begin{aligned} \mathbb{P}((\mathcal{E}_{\text{fb}}^\lambda)^c) &\leq 2d \sum_{r=1}^{\lceil \log_2 h \rceil} n_r \mathbb{E} \left[\exp \left(-\frac{\lambda^2 \Delta_{n_{r+1}+1}^2 T}{24\sigma^2 h_r \log_2 h} \right) \right] \\ &\leq 2d \sum_{r=1}^{\lceil \log_2 h \rceil} n_r \exp \left(-\frac{\lambda^2 T \Delta_{n_{r+1}+1}^2}{24\sigma^2 \min(h, n_r) \lceil \log_2 h \rceil} \right), \quad \text{as } h_r \leq \min(n_r, h). \end{aligned}$$

Then, note that

$$\begin{aligned} \frac{\Delta_{n_{r+1}+1}^2}{\min(h, n_r)} &= \frac{\Delta_{\lceil h/2^r \rceil + 1}^2}{\lceil h/2^{r-1} \rceil} \\ &= \frac{\Delta_{\lceil h/2^r \rceil + 1}^2}{\lceil h/2^r \rceil + 1} \frac{\lceil h/2^r \rceil + 1}{\lceil h/2^{r-1} \rceil} \\ &\geq \frac{\Delta_{\lceil h/2^r \rceil + 1}^2}{\lceil h/2^r \rceil + 1} \frac{h/2^r + 1}{h/2^{r-1} + 1} \\ &\geq \frac{\Delta_{\lceil h/2^r \rceil + 1}^2}{\lceil h/2^r \rceil + 1} \frac{1}{2}, \end{aligned}$$

which follows as $(x+1)/(2x+1) = (1/2) + (1/2)/(2x+1) \geq 1/2$ for $x \geq 0$. Therefore,

$$\begin{aligned} \frac{\Delta_{n_{r+1}+1}^2}{\min(h, n_r)} &\geq \frac{1}{2} \frac{\Delta_{\lceil h/2^r \rceil + 1}^2}{\lceil h/2^r \rceil + 1} \\ &\geq \frac{1}{2H_{2,\text{lin}}}. \end{aligned}$$

Thus,

$$\begin{aligned} \mathbb{P}((\mathcal{E}_{\text{fb}}^\lambda)^c) &\leq 2 \exp \left(-\frac{\lambda^2 T}{48\sigma^2 H_{2,\text{lin}} \lceil \log_2 h \rceil} + \log(d) \right) \sum_{r=1}^{\lceil \log_2 h \rceil} n_r \\ &\leq 2(K + h + \lceil \log_2 h \rceil) \exp \left(-\frac{\lambda^2 T}{48\sigma^2 H_{2,\text{lin}} \lceil \log_2 h \rceil} + \log(d) \right) \end{aligned}$$

Finally, it follows that

$$\inf_{\lambda \in (0, 1/5)} \mathbb{P}((\mathcal{E}_{\text{fb}}^\lambda)^c) \leq 2(K + h + \lceil \log_2 h \rceil) \exp \left(-\frac{T}{1200\sigma^2 H_{2,\text{lin}} \lceil \log_2 h \rceil} + \log(d) \right),$$

which concludes the proof. \square

4.5.3 Sample complexity: Proof of Theorem 4.3.3

We prove the theoretical guarantees in the fixed-confidence setting. We prove the correctness of Algorithm 4.2 and we prove the sample complexity bound of Theorem 4.3.3 and some key lemmas. We first prove the correctness of the fixed-confidence variant of GEGE.

Proof of the correctness. We need to prove that the final recommendation of Algorithm 4.2 is correct: that is we should show that: at any round r , $\mathcal{B}_r \subset \mathcal{S}^*$ and $\mathcal{D}_r \subset (\mathcal{S}^*)^c$.

Lemma 4.5.7. *On the event \mathcal{E}_{fc} , Algorithm 4.2 identifies the correct Pareto set.*

Proof of Lemma 4.5.7. In this part let τ denotes the stopping time of Algorithm 4.2. We assume \mathcal{E}_{fc} holds.

Using Proposition 4.3.6: for any round $r \leq \tau$ for any (Pareto) suboptimal $i \in \mathcal{A}_r$ we have $i^* \in \mathcal{A}_r$. We then prove the correctness of the algorithm as follows. Let i be an arm that is removed at the end of some round r . Assume $i \in S_r$ then, as i is discarded and empirically optimal we have $\widehat{\Delta}_{i,r} = \widehat{\delta}_i^* \geq \varepsilon_r$. In particular, it holds that

$$\min_{j \in \mathcal{A}_r \setminus \{i\}} M(i, j; r) \geq \varepsilon_r$$

which using Lemma 4.5.1 on the event \mathcal{E}_{fc} yields

$$\min_{j \in \mathcal{A}_r \setminus \{i\}} M(i, j) > \varepsilon_r/2 > 0,$$

that is no active arm dominates i . Combined with proposition \mathcal{P}_∞ (cf Lemma 4.5.6) the latter inequality yields $i \in \mathcal{S}^*$. Now assume we have $i \notin S_r$: i is discarded and it is empirically suboptimal. Then

$$\widehat{\Delta}_{i,r} = \max_{j \in \mathcal{A}_r} m(i, j; r) \geq \varepsilon_r/2,$$

so using Lemma 4.5.1 again on event \mathcal{E}_{fc} it follows that there exists $j \in \mathcal{A}_r$ such that $m(i, j) > 0$: that is $i \notin \mathcal{S}^*$. In all cases, we have proved that if \mathcal{E}_{fc} holds then for any arm i discarded at some round r ,

$$i \in \mathcal{B}_{r+1} \iff i \in \mathcal{S}^*.$$

Note that if \mathcal{A}_r is non-empty, then it contains a single arm and because \mathcal{P}_∞ holds, this arm is also Pareto optimal. \square

Thus, Algorithm 4.2 is correct on \mathcal{E}_{fc} . Before proving Theorem 4.3.3 we need Lemma 4.3.7 to control the size of the active set \mathcal{A}_r in the fixed-confidence setting.

Size of the active set. We prove the following result that controls the size of the active set.

Proof of Lemma 4.3.7. By Lemma 4.5.6 we have on the event \mathcal{E}_{fc} : for any round r and for any arm $i \in \mathcal{A}_r$,

$$\widehat{\Delta}_{i,r} - \Delta_i \geq \begin{cases} -\varepsilon_r & \text{if } i \in \mathcal{S}^* \\ -\varepsilon_r/2 & \text{else.} \end{cases}$$

Then let $p \in [K]$ and let us assume an arm $i \in \{p, \dots, K\}$ is still active at round $r = \lceil \log_2(1/\Delta_p) \rceil$. We have $\widehat{\Delta}_{i,r} \geq \Delta_i - \varepsilon_r$ with $\varepsilon_r = 1/2^{r+1}$ and $\Delta_i \geq \Delta_p$ which combined with $\widehat{\Delta}_{i,r} \geq \Delta_i - \varepsilon_r$ yields

$$\widehat{\Delta}_{i,r} \geq \Delta_p - \varepsilon_r. \quad (4.28)$$

As $r = \lceil \log_2(1/\Delta_p) \rceil$, it holds that $2\varepsilon_r \leq \Delta_p$ so Eq. (4.28) yields $\widehat{\Delta}_{i,r} \geq \varepsilon_r$ thus i will be discarded at the end of round r that is any arm $i \in \{p, \dots, K\}$ will be discarded at the end of round $\lceil \log_2(1/\Delta_p) \rceil$. \square

We now prove the main lemma on the sample complexity of GEGE in the fixed-confidence setting.

Bound on the sample complexity of Theorem 4.3.3. We provide an upper bound on the sample complexity of the algorithm.

Proof. We assume \mathcal{E}_{fc} holds. The correctness of Algorithm 4.2 is then proven in Lemma 4.5.7 and Lemma 4.3.7 upper-bounds the number of rounds before termination. It remains to bound the sample complexity of the algorithm on \mathcal{E}_{fc} and compute $\mathbb{P}(\mathcal{E}_{\text{fc}})$ to conclude.

By Lemma 4.3.7 an upper-bound on $|\mathcal{A}_r|$ for some specific rounds. Interestingly we can bound the sample complexity between consecutive "checkpoints rounds". In what follows, we rewrite the complexity as a sum of number of pulls between these intermediate "checkpoints rounds". Let us introduce the sequence $\{\alpha_s : s \geq 0\}$ defined as $\alpha_0 = 0$ and for any $s \geq 1$, $\alpha_s = \lceil \log_2(1/\Delta_{\lfloor h/2^s \rfloor}) \rceil$. We assume *w.l.o.g* that the sequence is non-decreasing and that the gaps are bounded in $(0, 1)$ (otherwise, we could start the sequence $(\alpha)_s$ from arms with gap smaller than 1). Simple calculation shows that $\alpha_{\lfloor \log_2(h) \rfloor} = \lceil \log_2(1/\Delta_1) \rceil$ and

$$\{1, \dots, \lceil \log_2(1/\Delta_1) \rceil\} = \bigcup_{s=1}^{\lfloor \log_2(h) \rfloor} \llbracket 1 + \alpha_{s-1}, \alpha_s \rrbracket. \quad (4.29)$$

Introducing

$$T_r := \frac{32(1 + 3\varepsilon_r)\sigma^2 h_r}{\varepsilon_r^2} \log\left(\frac{dn_r}{\delta_r}\right),$$

where $n_r = |\mathcal{A}_r|$, we have $t_r = \lceil T_r \rceil$, so $t_r \leq T_r + 1$. Using (4.29) then leads to

$$\begin{aligned} \sum_{r=1}^{\lceil \log_2(1/\Delta_1) \rceil} T_r &= \sum_{s=0}^{\lfloor \log_2(h) \rfloor - 1} \sum_{r=\alpha_s+1}^{\alpha_{s+1}} T_r \\ &=: \sum_{s=0}^{\lfloor \log_2(h) \rfloor - 1} N_s \end{aligned}$$

where $N_s = \sum_{r=\alpha_s+1}^{\alpha_{s+1}} T_r$ is "the number of arms pulls" between round $(\alpha_s + 1)$ and α_{s+1} .

Next we bound the term N_s for $s \in \{0, \dots, \lfloor \log_2(h) \rfloor - 1\}$. We recall that $h_r \leq \min(h, n_r)$ as, $n_r = |\mathcal{A}_r|$ is the number of active arms at round r and h_r is the dimension of the space spanned by the features of the active arms. Using Lemma 4.3.7 on \mathcal{E}_{fc} , it holds that for $r \geq \alpha_s + 1$

$$n_r \leq \begin{cases} K & \text{if } s = 0 \\ \lfloor h/2^s \rfloor - 1 & \text{if } s \geq 1 \end{cases} \quad (4.30)$$

Therefore for $s \in \{0, \dots, \lfloor \log_2(h) \rfloor - 1\}$ and for any $r \geq \alpha_s + 1$, we simply have $\min(h, n_r) \leq \lfloor h/2^s \rfloor$, so $h_r \leq \lfloor h/2^s \rfloor$ and even $h_r \leq \lfloor h/2^s \rfloor - 1$ if $s > 0$. In particular, it holds that

$$h_r \leq 2\lfloor h/2^{s+1} \rfloor \quad \text{for } r \geq \alpha_s + 1.$$

It then follows that

$$\tilde{N}_s = N_s / (32(1 + 3\varepsilon_1)\sigma^2) = \sum_{r=\alpha_s+1}^{\alpha_{s+1}} T_r / (32(1 + 3\varepsilon_1)\sigma^2) \quad (4.31)$$

$$\leq 2\lfloor h/2^{s+1} \rfloor \log\left(\frac{Kd}{\delta_{\alpha_{s+1}}}\right) \sum_{r=\alpha_s+1}^{\alpha_{s+1}} \frac{1}{\varepsilon_r^2} \quad (4.32)$$

$$= 8\lfloor h/2^{s+1} \rfloor \log\left(\frac{Kd}{\delta_{\alpha_{s+1}}}\right) \sum_{r=\alpha_s+1}^{\alpha_{s+1}} 4^r \quad (4.33)$$

$$\leq 8\lfloor h/2^{s+1} \rfloor \log\left(\frac{Kd}{\delta_{\alpha_{s+1}}}\right) \sum_{r=1}^{\alpha_{s+1}} 4^r \quad (4.34)$$

$$= \frac{32\lfloor h/2^{s+1} \rfloor}{3} \log\left(\frac{Kd}{\delta_{\alpha_{s+1}}}\right) (4^{\alpha_{s+1}} - 1) \quad (4.35)$$

then further using that

$$\alpha_s \geq \begin{cases} \log_2(1/\Delta_{\lfloor h/2^s \rfloor}) & \text{if } s \geq 1 \\ 0 & \text{if } s = 0 \end{cases}$$

yields

$$4^{\alpha_{s+1}} \leq \frac{1}{\Delta_{\lfloor h/2^{s+1} \rfloor}^2}$$

which combined with (4.35) yields

$$\tilde{N}_s \leq \frac{32\sigma^2\lfloor h/2^{s+1} \rfloor}{3\Delta_{\lfloor h/2^{s+1} \rfloor}^2} \log\left(\frac{Kd}{\delta_{\alpha_{s+1}}}\right). \quad (4.36)$$

We can now bound $N = \sum_s N_s$ in terms of the suboptimality gaps:

$$\tilde{N} = \sum_{s=0}^{\lfloor \log_2 h \rfloor - 1} \tilde{N}_s \quad (4.37)$$

$$\leq \frac{32\sigma^2}{3} \sum_{s=0}^{\lfloor \log_2 h \rfloor - 1} \frac{\lfloor h/2^{s+1} \rfloor}{\Delta_{\lfloor h/2^{s+1} \rfloor}^2} \log \left(\frac{\pi^2 K d \lceil \log_2(1/\Delta_{\lfloor h/2^{s+1} \rfloor}) \rceil^2}{6\delta} \right), \quad (4.38)$$

$$= \frac{32\sigma^2}{3} \sum_{s=1}^{\lfloor \log_2 h \rfloor} \frac{\lfloor h/2^s \rfloor}{\Delta_{\lfloor h/2^s \rfloor}^2} \log \left(\frac{\pi^2 K d \lceil \log_2(1/\Delta_{\lfloor h/2^s \rfloor}) \rceil^2}{6\delta} \right) \quad (4.39)$$

Then, recalling that by assumption $\Delta_1 \leq \dots \leq \Delta_K$, one can observe that the mapping from $[K]$ to $(0, \infty)$,

$$u \mapsto \frac{1}{\Delta_u^2} \log \left(\frac{\pi^2 K d \lceil \log_2(1/\Delta_u) \rceil^2}{6\delta} \right)$$

is non-increasing and it is easy to check that

$$\lfloor h/2^s \rfloor - \lceil \lfloor h/2^s \rfloor / 2 \rceil + 1 \geq \frac{1}{2} \lfloor h/2^s \rfloor$$

therefore

$$\frac{\lfloor h/2^s \rfloor}{\Delta_{\lfloor h/2^s \rfloor}^2} \log \left(\frac{\pi^2 K d \lceil \log_2(1/\Delta_{\lfloor h/2^s \rfloor}) \rceil^2}{6\delta} \right) \leq 2 \sum_{u=\lceil \lfloor h/2^s \rfloor / 2 \rceil}^{\lfloor h/2^s \rfloor} \frac{1}{\Delta_u^2} \log \left(\frac{\pi^2 K d \lceil \log_2(1/\Delta_u) \rceil^2}{6\delta} \right) \quad (4.40)$$

Combining (4.39) and (4.40) yields

$$N \leq \frac{64\sigma^2}{3} \sum_{s=1}^{\lfloor \log_2 h \rfloor} \sum_{u=\lceil \lfloor h/2^s \rfloor / 2 \rceil}^{\lfloor h/2^s \rfloor} \frac{1}{\Delta_u^2} \log \left(\frac{\pi^2 K d \lceil \log_2(1/\Delta_u) \rceil^2}{6\delta} \right) \quad (4.41)$$

Now let us introduce for any s , the set of integers $\mathcal{I}_s = \lceil \lceil \lfloor h/2^s \rfloor / 2 \rceil, \lfloor h/2^s \rfloor \rceil$. We have

$$\bigcup_{s=1}^{\lfloor \log_2 h \rfloor} \mathcal{I}_s \subset \{2, \dots, h\}.$$

We show that for any $p, q \in \{1, \dots, \lfloor \log_2(h) \rfloor\}$ if $|p - q| \geq 2$ then $\mathcal{I}_p \cap \mathcal{I}_q = \emptyset$. Assuming $p \leq q$ we claim that

$$\lfloor h/2^{p+2} \rfloor < \lceil \lfloor h/2^p \rfloor / 2 \rceil \quad (4.42)$$

Assume otherwise, then $\lfloor h/2^{p+2} \rfloor \geq \lceil \lfloor h/2^p \rfloor / 2 \rceil \geq \lfloor h/2^p \rfloor / 2$ so

$$h/2^{p+1} \geq \lfloor h/2^p \rfloor$$

which is impossible since for any $p \in \{0, \dots, \lfloor \log_2(h) \rfloor - 1\}$, $h/2^p \geq 1$. Therefore we have proved (4.42) and for any $q \geq p + 2$ it holds that

$$\lfloor h/2^q \rfloor \leq \lfloor h/2^{p+2} \rfloor < \lceil \lfloor h/2^p \rfloor / 2 \rceil$$

thus $\mathcal{I}_q \cap \mathcal{I}_p = \emptyset$ and for any $i \in \{2, \dots, h\}$, i belongs to no more than 2 of the subsets $\mathcal{I}_1, \dots, \mathcal{I}_{\lfloor \log_2 h \rfloor}$, thus we have

$$\tilde{N} \leq \frac{64}{3} \sigma^2 \sum_{s=1}^{\lfloor \log_2 h \rfloor} \sum_{u=\lceil h/2^s \rceil/2}^{\lfloor h/2^s \rfloor} \frac{1}{\Delta_u^2} \log \left(\frac{\pi^2 K d \lceil \log_2(1/\Delta_u) \rceil^2}{6\delta} \right) \quad (4.43)$$

$$\leq \frac{128}{3} \sigma^2 \sum_{i=2}^h \frac{1}{\Delta_i^2} \log \left(\frac{\pi^2 K d \lceil \log_2(1/\Delta_i) \rceil^2}{6\delta} \right) \quad (4.44)$$

$$\leq \frac{128}{3} \sigma^2 \sum_{i=2}^h \frac{1}{\Delta_i^2} \log \left(\frac{\pi^2 K d \log_2(2/\Delta_i)^2}{6\delta} \right) \quad (4.45)$$

$$\leq \frac{256}{3} \sigma^2 \sum_{i=2}^h \frac{1}{\Delta_i^2} \log \left(\frac{K d}{\delta} \log_2 \left(\frac{2}{\Delta_i} \right) \right) \quad (4.46)$$

Then, from Lemma 4.5.6 it holds that with probability at least $1 - \delta$ the sample complexity N^δ of GEGE is upper-bounded as

$$\log_2(2/\Delta_1) + \mathcal{O} \left(\sum_{i=2}^h \frac{\sigma^2}{\Delta_i^2} \log \left(\frac{K d}{\delta} \log_2 \left(\frac{1}{\Delta_i} \right) \right) \right),$$

□

where $\mathcal{O}(\cdot)$ hides universal multiplicative constant. Therefore, we have shown the sample complexity bound and the correctness on \mathcal{E}_{fc} . Thus, proving that $\mathbb{P}(\mathcal{E}_{fc}) \geq 1 - \delta$ will conclude the proof.

Probability of the good event \mathcal{E}_{fc} .

Proof. At round r ,

$$\mathbb{P}((\mathcal{E}_{fc}^r)^c \mid \mathcal{A}_r) \leq \sum_{i \in \mathcal{A}_r} \mathbb{P} \left(\|(\hat{\theta}_r - \theta)^\top x_i\|_\infty > \varepsilon_r/4 \mid \mathcal{A}_r \right)$$

Then, note that at round r , Algorithm 4.2 calls OPTESTIMATOR with precision $\varepsilon_r/2$ and budget t_r and by design we have $t_r \geq 20h_r/\varepsilon_r^2$, so using Lemma 4.2.2, it follows

$$\begin{aligned} \mathbb{P}((\mathcal{E}_{fc}^r)^c \mid \mathcal{A}_r) &\leq 2d \exp \left(-\frac{t_r \varepsilon_r^2}{32(1 + 3\varepsilon_r)\sigma^2 h_r} \right) \\ &\leq \delta_r / |\mathcal{A}_r| \end{aligned}$$

which follows by plugging in the value of t_r . Therefore, by union bound over \mathcal{A}_r and r it holds that $\mathbb{P}(\mathcal{E}_{fc}) \geq 1 - \sum_{r \geq 1} \delta_r \geq 1 - \delta$. □

This concludes the proof of Theorem 4.3.3.

4.5.4 Lower bounds

We now prove minimax lower bounds in both fixed-confidence and fixed-budget settings. We recall the lower bound below for unstructured PSI in the fixed confidence setting.

Theorem 4.5.8 (Theorem 17 of [Auer et al. 2016](#)). *For any set of operating points $\mu_i \in [1/4, 3/4]^d$, $i = 1, \dots, K$, there exist distributions $(\mathcal{D}_i)_{1 \leq i \leq K}$ such that with probability at least $1 - \delta$, any δ -correct algorithm for linear PSI requires at least*

$$\Omega \left(\sum_{i=1}^K \frac{1}{\tilde{\Delta}_i^2} \log(\delta^{-1}) \right)$$

samples to identify the Pareto set. Where for any suboptimal arm $\tilde{\Delta}_i = \Delta_i$ and for an optimal arm $\tilde{\Delta}_i = \delta_i^+$.

In particular, there exist instances where $\Delta_i = \delta_i^+$ for any Pareto-optimal arm i . Thus, this result shows that H_1 is a good proxy to measure the complexity of PSI in the fixed-confidence setting.

The proof of such results often relies on the celebrated change of distribution technique (see e.g. [Kaufmann, Cappé, et al. 2016](#)) which, given the instance $\nu := (\nu_1, \dots, \nu_K)$, shifts the mean of ν_i for an arm i while keeping the others fixed constant. However, in linear PSI, the arms' means are correlated through Θ . So, in general, [Theorem 4.5.8](#) does not directly apply to linear PSI. We recall below our lower bound for linear PSI in the fixed-confidence setting.

Theorem 4.3.5. *For any $K, d, h \in \mathbb{N}$, there exists a set $\mathcal{M}(K, d, h)$ of linear PSI instances such that for $\nu \in \mathcal{M}(K, d, h)$ and for any δ -correct algorithm for linear PSI, with probability at least $1 - \delta$,*

$$\tau_\delta^A \geq \Omega \left(H_{1, \text{lin}}(\nu) \log \frac{1}{\delta} \right).$$

Proof of [Theorem 4.3.5](#). The idea of the proof is to transform an unstructured bandit instance into a linear PSI instance. Let ν be a bandit instance with $K \geq 2$ arms and dimension $d \geq 1$ and with means $\mu_1, \dots, \mu_K \in [0, 1]^d$. Let e_1, \dots, e_h denote the canonical basis of \mathbb{R}^h . We define a linear PSI instance ν_{lin} with features

$$x_i = \begin{cases} e_i & \text{if } i \leq h \\ \mathbf{0} & \text{else.} \end{cases}$$

We assume that the learner knows that $\mu_i \in [0, 1]^d$ for any arm i . We claim that with this information, an "efficient" algorithm for PSI should not pull arms from $\{h + 1, \dots, K\}$. To see this, first note that these arms will be suboptimal, so $\mathcal{S}^* \subset [h]$. Moreover, even if an arm $i \in \{h + 1, \dots, K\}$ dominates another arm $j \in \{1, \dots, h\}$, as j is not Pareto-optimal there exists another arm $j^* \in \mathcal{S}^* \subset \{1, \dots, h\}$ which dominates j with a larger margin,

so is "cheaper" to pull. Therefore, the complexity of ν_{lin} reduces to the complexity of a linear bandit $\tilde{\nu}_{\text{lin}}$ with only h arms. As the features in x_1, \dots, x_h form the canonical \mathbb{R}^h basis, $\tilde{\nu}_{\text{lin}}$ reduces to an unstructured bandit instance with (uncorrelated) means $\tilde{\mu}_i = \theta^\top x_i$, $i = 1, \dots, h$. Therefore, by choosing $\mu_1, \dots, \mu_h \in [1/4, 3/4]^d$, we can apply Theorem 4.5.8 to $\tilde{\nu}_{\text{lin}}$. \square

The result proven above holds for a class of instances $\mathcal{M}(K, d, h)$ with the covariates defined as above and with matrix coefficients in $[1/4, 3/4]^d$. For the fixed-budget setting, [Kone, Kaufmann, et al. 2024](#) proved a lower bound for a class of instances. We recall their result below after introducing some notation.

Their lower bound applies to the class of instances \mathcal{M} defined as follows. \mathcal{M} contains the instances such that each suboptimal arm i is only dominated by a Pareto-optimal arm denoted by i^* and that for each optimal arm j there exists a unique suboptimal arm which is dominated by j , denoted by \underline{j} . Moreover, for any instance in \mathcal{M} the authors require its Pareto-optimal arms not to be close to the suboptimal arms they don't dominate: for any suboptimal arm i and Pareto-optimal arm j such that $\mu_i \not\leq \mu_j$,

$$M(i, j) \geq 3 \max(\Delta_i, \Delta_{\underline{j}}).$$

Let $\nu := (\nu_1, \dots, \nu_K)$ be an unstructured instance whose means belongs to \mathcal{M} and with isotropic multi-variate normal arms $\nu_i \sim \mathcal{N}(\mu_i, \sigma^2 I)$. For every $i \in [K]$, define the alternative instance $\nu^{(i)} := (\nu_1, \dots, \nu_i^{(i)}, \dots, \nu_K)$ in which *only* the mean of arm i is shifted:

$$\mu_i^{(i)} := \begin{cases} \mu_i - 2\Delta_i \tilde{e}_{d_i} & \text{if } i \in \mathcal{S}^*(\nu), \\ \mu_i + 2\Delta_i \tilde{e}_{d_i} & \text{else,} \end{cases} \quad (4.47)$$

where $\tilde{e}_1, \dots, \tilde{e}_d$ denotes the canonical basis of \mathbb{R}^d and for any arm i , $d_i := \operatorname{argmin}_{c \in [d]} [\mu_{i^*}^c - \mu_i^c]$. Defining $\nu^{(0)} := \nu$, the theorem below holds.

Theorem 4.5.9 (Theorem 5 of [Kone, Kaufmann, et al. 2024](#)). *Let $\nu = (\nu_1, \dots, \nu_K)$ be an instance in \mathcal{M} with means $\mu := (\mu_1 \dots \mu_K)^\top$ and $\nu_i \sim \mathcal{N}(\mu_i, \sigma^2 I)$. For any algorithm \mathcal{A} , there exists $i \in \{0, \dots, K\}$ such that $H(\nu^{(i)}) \leq H(\nu)$ and the probability of error \mathcal{A} on $\nu^{(i)}$ is at least*

$$\frac{1}{4} \exp\left(-\frac{2T}{\sigma^2 H(\nu^{(i)})}\right).$$

As explained above for the fixed-confidence setting. The proof of this lower bound also uses the change of distribution lemma. In the instances $\nu^{(i)}$ introduced above, only the mean of arm i must change w.r.t. $\nu^{(0)}$. Therefore, Theorem 4.5.9 does not apply to general instances in linear PSI. We recall our lower bound for linear PSI in the fixed-budget setting.

Theorem 4.3.2. *Let \mathbb{W}_H be the set of instances with complexity $H_{2, \text{lin}}$ smaller than H . For any budget T , letting \hat{S}_T^{alg} be the output of an algorithm alg , it holds that*

$$\inf_{\text{alg}} \sup_{\nu \in \mathbb{W}_H} \mathbb{P}_\nu(\hat{S}_T^{\text{alg}} \neq \mathcal{S}^*(\nu)) \geq \frac{1}{4} \exp\left(-\frac{2T}{H\sigma^2}\right).$$

Proof of Theorem 4.3.2. Let H be fixed and recall that $\mathbb{W}_H : \{\nu_{\text{lin}} : H_{2,\text{lin}}(\nu) \leq H\}$ is the set of linear PSI instances with complexity less than H . The proof of Theorem 4.3.2 follows similar lines to Theorem 4.3.5. Let ν be an un-structured bandit instance with $K \geq 2$ arms, dimension $d \geq 1$, with means $\mu_1, \dots, \mu_K \in [0, 1]^d$ and such that $H_2(\nu) \leq H$. We construct a linear PSI instance ν_{lin} from an unstructured multi-dimensional instance ν by setting $x_i := e_i$ for any $i \leq h$ and for $i > h$, $x_i = \mathbf{0}$ where e_1, \dots, e_h is the canonical \mathbb{R}^h -basis. We also assume that the agent knows that $\mu_i \in [0, 1]^d$ for any arm i .

For ν_{lin} , the arms $\{h+1, \dots, K\}$ are necessarily suboptimal, so $\mathcal{S}^* \subset [h]$. Thus, to identify the Pareto set, an efficient algorithm should reduce to pulling arms in $\{1, \dots, h\}$. Indeed, as explained in the proof of Theorem 4.3.5 even if an arm $i \in \{h+1, \dots, K\}$ dominates another arm $j \in \{1, \dots, h\}$, as j is not Pareto-optimal there exists another arm $j^* \in \mathcal{S}^* \subset \{1, \dots, h\}$ which is "cheaper" to pull as it dominates j with a larger margin. ν_{lin} reduces to a linear bandit $\tilde{\nu}_{\text{lin}}$ with only h arms and since the features x_1, \dots, x_h forms the canonical basis of \mathbb{R}^h , $\tilde{\nu}_{\text{lin}}$ is an un-structured bandit instance with (un-correlated) means $\tilde{\mu}_i = \theta^\top x_i$, $i = 1, \dots, h$. Therefore, by choosing $\tilde{\nu} := (\nu_1, \dots, \nu_h)$ that belongs to \mathcal{M} , we can apply Theorem 4.5.9 which yields

$$\max_{i \in \{0, \dots, K\}} \mathbb{P}_{\tilde{\nu}^{(i)}}(S_T^{\text{alg}} \neq \mathcal{S}^*(\tilde{\nu}^{(i)})) \geq \frac{1}{4} \exp\left(-\frac{2T}{H\sigma^2}\right)$$

where by construction $\tilde{\nu}^{(i)}$ (see construction above) is also a linear PSI instance. Then, further noting that $H \geq H_2(\nu) \geq H_2(\tilde{\nu})$ and by Theorem 4.5.9 for any $i \leq h$ $H_{2,\text{lin}}(\tilde{\nu}) \geq H_2(\tilde{\nu}^{(i)})$. Then recalling that ν_{lin} is equivalent to $\tilde{\nu}$ it comes

$$\inf_{\text{alg}} \sup_{\nu \in \mathbb{W}_H} \mathbb{P}_\nu(S_T^{\text{alg}} \neq \mathcal{S}^*(\nu)) \geq \frac{1}{4} \exp\left(-\frac{2T}{H\sigma^2}\right),$$

which is the claimed result. \square

4.5.5 Concentration lemmas

Proof of Lemma 4.5.1. We have

$$\begin{aligned} |M(i, j; r) - M(i, j)| &= \left| \max_c [\hat{\mu}_{i,r}^c - \hat{\mu}_{j,r}^c] - \max_c [\mu_i^c - \mu_j^c] \right|, \\ &\stackrel{(i)}{\leq} \max_c |(\hat{\mu}_{i,r}^c - \hat{\mu}_{j,r}^c) - (\mu_i^c - \mu_j^c)|, \\ &= \|(\hat{\mu}_{i,r} - \hat{\mu}_{j,r}) - (\mu_i - \mu_j)\|_\infty, \\ &= \|(\hat{\theta}_r - \theta)^\top (x_i - x_j)\|_\infty. \end{aligned}$$

where (i) follows from the reverse triangle inequality. The second part of the lemma is a direct consequence of the relation $M(i, j) = -m(i, j)$ as well as $M(i, j; r) = -m(i, j; r)$ that holds for any pair of arms i, j . \square

Proof of Lemma 4.5.2. Since i^* does not empirically dominate i it holds that $M(i, i^*; r) > 0$ so $M(i, i^*; r) - M(i, i^*) > -M(i, i^*)$. Then noting that

$$-M(i, i^*) = m(i, i^*) = \Delta_i$$

yields $M(i, i^*; r) - M(i, i^*) > \Delta_i$. Therefore

$$\begin{aligned} \Delta_i = \Delta_i^* &< M(i, i^*; r) - M(i, i^*) \\ &\leq \|(\hat{\theta}_r - \theta)^\top(x_i - x_{i^*})\|_\infty, \end{aligned}$$

where the last inequality is a consequence of Lemma 4.5.1. □

Chapter 5

Posterior Sampling for Pareto Set Identification

In previous chapters, we studied Pareto Set Identification (PSI) through deterministic, confidence-sequence-based algorithms that offer strong finite-sample guarantees but suffer from three main limitations: the confidence level δ must be fixed *a priori*, asymptotic optimality is not achieved as $\delta \rightarrow 0$, and correlations between objectives are typically ignored.

This chapter introduces a posterior-sampling approach that overcomes these issues. Focusing on the unstructured multivariate Gaussian setting with possibly correlated objectives, we propose PSIPS (*Pareto Set Identification with Posterior Sampling*), an algorithm that uses posterior draws in both the sampling rule and the stopping rule. PSIPS removes the need to pre-specify δ , achieves asymptotic optimality from both Bayesian and frequentist perspectives, and naturally exploits correlations between objectives.

Beyond its theoretical guarantees, PSIPS remains computationally tractable by avoiding costly optimization oracles. This chapter is based on joint work with Marc Jourdan and Émilie Kaufmann, published in the proceedings of *AISTATS 2025*, where PSIPS is further analyzed in the more general linear PSI setting.

5.1	Introduction	122
5.2	PSI with Posterior Sampling	125
5.2.1	The Posterior Sampling (PS) Stopping Rule	125
5.2.2	Game-based sampling rule	128
5.3	Main theoretical results	129
5.3.1	Sketch of proofs	130
5.4	Numerical study and discussion	134
5.5	Additional proofs	137
5.5.1	Stopping rule	137
5.5.2	Sample complexity	146
5.5.3	Posterior convergence	149

5.1 Introduction

In previous chapters, we investigated Pareto Set Identification (PSI) using deterministic confidence-sequence-based methods that provide rigorous finite-time guarantees. While effective in practice, these approaches exhibit three key limitations. First, they require fixing the confidence parameter δ before data collection, forcing practitioners to predefine an acceptable risk level. Second, although these methods achieve sample complexity of $\mathcal{O}(H(\nu) \log(H(\nu)d/\delta))$ where $H(\nu) := \sum_{k \in [K]} \Delta_k^{-2}$ depends on instance-specific gaps, they fail to achieve *asymptotic optimality*—their sample complexity does not match the information-theoretic lower bound as $\delta \rightarrow 0$. Evidence from best-arm identification suggests LUCB-style algorithms incur suboptimal constants in this regime (Kejriwal et al. 2025). Third, while not assuming independence between objectives, they do not explicitly exploit correlations common in applications such as efficacy-toxicity trade-offs in clinical trials.

This chapter addresses these limitations through a *posterior sampling* approach. We propose PSIPS (Posterior Sampling for Pareto Set Identification), an algorithm that uses random draws from a posterior distribution for both exploration and stopping decisions. This algorithmic design eliminates the need to pre-specify δ , naturally accounts for correlations across objectives through the covariance structure Σ , and achieves asymptotic optimality. Specifically, we establish that $\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}[\tau_\delta]}{\log(1/\delta)}$ matches the information-theoretic lower bound, answering a fundamental question left open by previous chapters.

We focus on the unstructured Gaussian setting where each arm $i \in [K]$ yields independent d -dimensional observations from $\mathcal{N}(\mu_i, \Sigma)$ with known covariance Σ . This formulation explicitly captures correlations between objectives, enabling the learner to transfer information between outcomes and reduce sample complexity. These properties make PSIPS particularly relevant for adaptive clinical trials, where multiple correlated outcomes must be analyzed simultaneously, for instance, in dose-finding trials where higher efficacy correlates with higher toxicity.

Related work. The algorithms of Chapters 3 and 4, along with Auer et al. 2016, rely on confidence intervals requiring δ to be specified in advance. While practical, these methods target finite-time guarantees rather than asymptotic optimality and do not exploit the covariance structure. Crepon et al. 2024 achieved asymptotic optimality through gradient-based optimization inspired by the information-theoretic lower bound, but their approach requires solving $\mathcal{O}(K^d)$ convex programs per iteration and assumes independent objectives ($\Sigma = \sigma^2 I_d$). More broadly, asymptotically optimal algorithms based on optimization or game theory (Garivier & Kaufmann 2016; Degenne, W. M. Koolen, et al. 2019; Degenne, Ménard, et al. 2020) provide strong guarantees but remain computationally expensive in multi-objective settings.

Posterior sampling approaches for pure exploration extend Thompson Sampling’s success from regret minimization (Thompson 1933; Kaufmann, Korda, et al. 2012) to identification problems. Russo 2016 introduced Top-Two Thompson Sampling for best-arm identification with optimal posterior contraction but no stopping rule. S. Wang & Zhu 2022 proposed TS-EXPLORE using posterior samples for exploration and stopping, though its gap-based

sampling remains asymptotically suboptimal. The PEPS algorithm (Z. Li et al. 2024) achieves optimal Bayesian convergence in linear bandits but lacks fixed-confidence guarantees. Most multi-objective algorithms assume independent objectives (Auer et al. 2016; Zuluaga, Krause, et al. 2016; Katz-Samuels & Scott 2018), limiting efficiency when outcomes are correlated.

In summary, existing PSI algorithms either rely on confidence-sequence-based methods that require a pre-specified δ , or they achieve asymptotic optimality at the cost of computational intractability and restrictive independence assumptions.

PSIPS integrates posterior sampling into both sampling and stopping, achieving asymptotic optimality while remaining computationally efficient and naturally exploiting objective correlations.

Learning model. We consider the standard multi-armed bandit framework in the multi-objective setting with dimension $d \geq 2$. The collection of means is denoted by $\mu = (\mu_i)_{i \in [K]} \in \mathbb{R}^{K \times d}$. Throughout this chapter, we focus on the *Gaussian case*:

$$\nu_i = \mathcal{N}(\mu_i, \Sigma), \quad i \in [K],$$

where the covariance matrix Σ is known and identical across arms. The off-diagonal entries of Σ encode dependence—for instance, between toxicity and efficacy in dose-finding trials—allowing the learner to transfer information across outcomes. Such correlations are common in clinical and engineering applications and play a key role in improving sample efficiency.

At each time t , the learner selects an arm $A_t \in [K]$ based on past observations and obtains a sample $Z_t \in \mathbb{R}^d$ with conditional distribution ν_{A_t} given A_t . An *algorithm* specifies both a *sampling rule* that determines A_t and a *recommendation rule* that outputs a candidate Pareto set $\widehat{S}_t \subseteq [K]$ before each arm pull. Formally, both A_t and \widehat{S}_t are \mathcal{H}_t -measurable, where the σ -algebra $\mathcal{H}_t := \sigma(U_1, A_1, Z_1, \dots, A_{t-1}, Z_{t-1}, U_t)$ represents the *history* before time t , and $U_t \sim \mathcal{U}([0, 1])$ captures any exogenous randomness used by the algorithm.

For completeness, we recall that an arm i is said to be (Pareto) dominated by an arm j , written $\mu_i \prec \mu_j$, if $\mu_i^c \leq \mu_j^c$ for all $c \in [d]$ and the inequality is strict for at least one coordinate. The *Pareto set* is then

$$\mathcal{S}^*(\mu) := \{i \in [K] : \nexists j \in [K] \setminus \{i\}, \mu_i \prec \mu_j\},$$

that is, the collection of arms not strictly dominated by any other.

The empirical allocation over arms is $N_t := (N_{t,a})_{a \in [K]}$, where $N_{t,a} := \sum_{s \in [t-1]} \mathbb{1}(A_s = a)$, so that $\frac{N_t}{t-1} \in \Delta_K := \{w \in \mathbb{R}_+^K : \sum_{a \in [K]} w_a = 1\}$ belongs to the probability simplex. For vector $x \in \mathbb{R}^d$, $x(c)$ and x^c denote c -th coordinate. In the sequel, we let $\mathcal{I} := \mathbb{R}^d$.

Lower Bound. To formalize the notion of asymptotic optimality, we now characterize the minimal number of samples required by any δ -correct algorithm to reliably identify the Pareto set. A δ -correct algorithm must be able to statistically distinguish the true mean matrix μ from all *alternative instances* that would induce a different Pareto set. We define

this set of alternatives as

$$\text{Alt}(\mathcal{S}^*(\mu)) := \{\lambda \in \mathbb{R}^{K \times d} : \mathcal{S}^*(\lambda) \neq \mathcal{S}^*(\mu)\}.$$

Intuitively, these are the configurations of mean vectors that would lead the learner to recommend an incorrect Pareto set. The requirement to separate μ from $\text{Alt}(\mathcal{S}^*)$ leads to a fundamental information-theoretic lower bound on the expected stopping time (Garivier & Kaufmann 2016; Garivier, Ménard, et al. 2019; Crepon et al. 2024).

Lemma 5.1.1. *Let \mathcal{D} be the class of d -dimensional Gaussian distributions with known covariance Σ . For any δ -correct algorithm on \mathcal{D}^K , and any bandit instance $\nu \in \mathcal{D}^K$ we have*

$$\mathbb{E}_\nu[\tau_\delta] \geq T^*(\nu) \log\left(\frac{1}{2.4\delta}\right),$$

where $T^*(\nu)$, called the characteristic time, is defined by

$$T^*(\nu)^{-1} := \sup_{w \in \Delta_K} \inf_{\lambda \in \text{Alt}(\mathcal{S}^*(\mu))} \left[\sum_{k=1}^K \frac{w_k}{2} \|\mu_k - \lambda_k\|_{\Sigma^{-1}}^2 \right].$$

The set of maximizers of the outer optimization, denoted by $w^*(\nu)$, corresponds to the asymptotically optimal allocation of samples across arms.

This result establishes that no δ -correct algorithm can, in expectation, identify the Pareto set faster than $T^*(\nu) \log(1/(2.4\delta))$ samples. An algorithm is said to be *asymptotically optimal* if it matches *exactly* this lower bound asymptotically, *i.e.*, if

$$\forall \nu \in \mathcal{D}^K, \lim_{\delta \rightarrow 0} \frac{\mathbb{E}_\nu[\tau_\delta]}{\log \frac{1}{\delta}} = T^*(\nu). \quad (5.1)$$

Chernoff stopping rules. In fixed-confidence pure exploration, stopping rules determine when to terminate sampling. The classical *Chernoff* or *generalized likelihood ratio* (GLR) stopping rule (Garivier & Kaufmann 2016) stops when

$$\tau_\delta^{\text{GLR}} = \inf\{t \geq 1 \mid \text{GLR}(t) > \beta(t, \delta)\}, \quad \text{with}$$

$$\text{GLR}(t) = \inf_{\lambda \in \text{Alt}(\widehat{\mathcal{S}}_t)} \sum_{i=1}^K \frac{N_{t,i}}{2} \|\hat{\mu}_{t,i} - \lambda_i\|_{\Sigma^{-1}}^2. \quad (5.2)$$

While the GLR rule ensures δ -correctness and can achieve asymptotic optimality with a proper sampling rule, its evaluation is computationally demanding.

To the best of our knowledge, computing (5.2) is only tractable for the setting with independent objectives, *i.e.*, diagonal Σ . Crepon et al. 2024 proposed the best-known procedure in that restrictive setting, which requires solving $\mathcal{O}(K d^3 |\widehat{\mathcal{S}}_t|^d)$ convex problems at each round. Their computationally costly procedure is difficult to extend to correlated objectives.

Sampling rules. Given a stopping rule, the sampling rule should make it stop as soon as possible. While gap-based sampling rules are suboptimal in their asymptotic allocation, the asymptotic optimality of Track-and-Stop (Garivier & Kaufmann 2016) comes at the cost of computational intractability since it computes $w_t \in w^*(\hat{\mu}_t)$ at time t . The game-based approach arises from viewing the lower bound as the value of a two-player zero-sum game (Degenne, W. M. Koolen, et al. 2019). Using two no-regret learning algorithms against each other yields a saddle-point algorithm that sequentially approximates $T^*(\nu)^{-1}$. The popularized instance uses the best response as a min learner, *i.e.*,

$$\lambda^{\text{BR}} : w \mapsto \operatorname{argmin}_{\lambda \in \text{Alt}(\hat{S}_t)} \sum_{k=1}^K \frac{w_k}{2} \|\hat{\mu}_{t,k} - \lambda_k\|_{\Sigma^{-1}}^2, \quad (5.3)$$

whose computational cost is the same as the GLR statistic (5.2). For the PS stopping rule, we rely on an inflated posterior sampling (PS) learner for the min learner. Therefore, our algorithm can tackle structured settings with correlated objectives and has a low computational cost.

Contributions. As the main contribution, we propose the **PSIPS** algorithm, which is the first computationally efficient, asymptotically optimal PSI algorithm with correlated objectives. **PSIPS** builds on posterior sampling in both the PS stopping rule and the min learner of the game-based sampling rule. By removing the oracle calls to (5.2) or (5.3), **PSIPS** deals with the PSI structure in a computationally efficient way.

From a frequentist perspective, it is shown to be asymptotically optimal for PSI with correlated objectives (Theorem 5.3.1). From a Bayesian perspective, the posterior probability that **PSIPS** misidentifies the Pareto set decays exponentially fast (as a function of t) with the optimal rate $1/T^*(\nu)$ (Theorem 5.3.2). Our experiments on both synthetic and real-world instances showcase the superior performance of our algorithm in terms of sample complexity and computational cost.

5.2 PSI with Posterior Sampling

We propose the **PSIPS** (PSI with Posterior Sampling) algorithm for PSI with correlated objectives. **PSIPS** combines the PS stopping rule and the game-based sampling rule using a PS min learner.

5.2.1 The Posterior Sampling (PS) Stopping Rule

We introduce the posterior sampling (PS) stopping rule. To specify the PS stopping rule, a budget function $M : \mathbb{N} \times (0, 1) \rightarrow \mathbb{R}_+$ and an inflation function $c : \mathbb{N} \times (0, 1) \rightarrow \mathbb{R}_+$ are given. The PS stopping rule stops when $M(t-1, \delta)$ independent realizations from a posterior distribution inflated by $c(t-1, \delta)$ all yield the same Pareto set as the current empirical estimate \hat{S}_t , suggesting that $S^*(\mu) = \hat{S}_t$. Conditioned on \mathcal{H}_t , let $(v_t^m)_{m \in [M(t, \delta)]}$ be *i.i.d.* draws

Algorithm 5.1: PSIPS: Pareto Set Identification with Posterior Sampling

Require: budget function M ; inflation function c ; exploration allocation $w_{\text{exp}} > 0$; fored exploration rate $\alpha > 0$; posterior inflation rate function η ; AdaHedge learner $\mathcal{L}^{[K]}$

- 1 Initialize: pull each arm once and compute $\hat{\mu}_1$
- 2 **for** $t = 1, 2, \dots$ **do**
 - // Current empirical Pareto set*
 - 3 $\hat{S}_t \leftarrow \mathcal{S}^*(\hat{\mu}_t)$;
 - 4 Initialize: $m \leftarrow 0$, $m_t \leftarrow +\infty$, and $m_{t,\delta} \leftarrow +\infty$;
 - // Posterior sampling and stopping rule*
 - 5 **while** $\max\{m_t, m_{t,\delta}\} = +\infty$ **do**
 - 6 $m \leftarrow m + 1$; sample $v_t^m \sim \bigotimes_{k=1}^K \mathcal{N}(0_d, \Sigma/N_{t,k})$;
 - 7 **if** $m_{t,\delta} = +\infty$ **then**
 - 8 $\theta_{t,k}^m \leftarrow \hat{\mu}_{t,k} + \sqrt{c(t-1, \delta)} v_{t,k}^m \quad \forall k$;
 - 9 **if** $\theta_t^m \in \text{Alt}(\hat{S}_t)$ **then**
 - 10 $m_{t,\delta} \leftarrow m$;
 - 11 **else if** $m \geq M(t-1, \delta)$ **then**
 - 12 **return** \hat{S}_t ;
 - 13 **if** $m_t = +\infty$ **then**
 - 14 $\lambda_{t,k}^m \leftarrow \hat{\mu}_{t,k} + \eta_t^{-1/2} v_{t,k}^m \quad \forall k$;
 - 15 **if** $\lambda_t^m \in \text{Alt}(\hat{S}_t) \cap \mathcal{I}_t^K$ **then**
 - 16 $m_t \leftarrow m$;
 - // Learner update and mixture exploration*
 - 17 Get w_t from learner $\mathcal{L}^{[K]}$ and set $\lambda_t \leftarrow \lambda_t^{m_t}$;
 - 18
$$\tilde{w}_t \leftarrow (1 - \gamma_t)w_t + \gamma_t w_{\text{exp}}, \quad \gamma_t = t^{-\alpha}.$$
 - // Sampling and observation*
 - 19 Sample $A_t \sim \tilde{w}_t$ and collect $Z_t \sim \nu_{A_t}$;
 - // Feed back to learner*
 - 20 Feed gain
 - $$g_t(w) = \sum_{k=1}^K \frac{w_k}{2} \|\lambda_{t,k} - \hat{\mu}_{t,k}\|_{\Sigma^{-1}}^2 \quad \text{to } \mathcal{L}^{[K]}$$

from the centered posterior distribution $\Pi_t := \bigotimes_{k=1}^K \mathcal{N}(0_d, \Sigma/N_{t,k})$. Each v_t^m is a collection of \mathbb{R}^d vectors $(v_{t,1}^m, \dots, v_{t,K}^m)$. For all $m \in [M(t-1, \delta)]$, let $\theta_{t,k}^m := \hat{\mu}_{t,k} + \sqrt{c(t-1, \delta)}v_{t,k}^m$. Then,

$$\tau_\delta^{\text{PS}} := \inf\{t \geq 1 \mid \forall m \in [M(t-1, \delta)], \mathcal{S}^*(\theta_t^m) = \widehat{S}_t\}.$$

When $\tau_\delta^{\text{PS}} < +\infty$, we recommend $\widehat{S}_{\tau_\delta^{\text{PS}}} := \mathcal{S}^*(\hat{\mu}_{\tau_\delta^{\text{PS}}})$. When $t < \tau_\delta^{\text{PS}}$, let $m_{t,\delta} := \inf\{m \mid \theta_t^m \in \text{Alt}(\widehat{S}_t)\}$.

Computational cost. To reduce the computational cost of **PSIPS**, we draw $(v_t^m)_m$ sequentially as the PS stopping condition is often infringed before $M(t, \delta)$ realizations are drawn (see Figure 5.6).

Compared to the oracle call to (5.2), testing that $\theta_t^m \notin \text{Alt}(\widehat{S}_t)$ has a lower computational cost. It scales as the sum of the costs of membership to $\text{Alt}(\widehat{S}_t)^c$, which is at most $\mathcal{O}(d|\widehat{S}_t| \max\{|\widehat{S}_t|, K - |\widehat{S}_t|\})$ since, thanks to Lemma 1 in Crepon et al. 2024

$$\text{Alt}(\widehat{S}_t)^c = \left(\bigcap_{(i,j) \in \widehat{S}_t^2, i \neq j} \bigcup_{c \in [d]} \{\lambda \mid \lambda_i^c \geq \lambda_j^c\} \right) \cap \left(\bigcap_{i \notin \widehat{S}_t} \bigcup_{j \in \widehat{S}_t} \bigcap_{c \in [d]} \{\lambda \mid \lambda_j^c \geq \lambda_i^c\} \right). \quad (5.4)$$

To update the candidate answer, we check whether $\hat{\mu}_t \in \text{Alt}(\widehat{S}_{t-1})^c$, in which case $\widehat{S}_t = \widehat{S}_{t-1}$. Otherwise, we compute $\widehat{S}_t = \mathcal{S}^*(\hat{\mu}_t)$. While a naive implementation has cost at most $\mathcal{O}(dK^2)$, Kung et al. 1975 proposed an algorithm having a cost scaling as $\mathcal{O}(K(\log K)^{\max\{1, d-2\}})$.

Correctness. Lemma 5.2.2 exhibits choices of budget $M(t, \delta)$ and inflation $c(t, \delta)$ ensuring δ -correctness in the setting with independent or correlated objectives. Letting $X \sim \mathcal{N}(0, 1)$, we denote by $R(x) := \frac{\mathbb{P}(X > x)}{f_X(x)}$, the Mills ratio (Mills 1926) of X , with f_X , the density of X .

We introduce the concentration event $\mathcal{E}_\delta := \bigcap_{t \in \mathbb{N}} \mathcal{E}_{t,\delta}$ with

$$\mathcal{E}_{t,\delta} := \left\{ \sum_{k=1}^K N_{t,k} \|\mu_k - \hat{\mu}_{t,k}\|_{\Sigma^{-1}}^2 \leq 2\beta(t-1, \delta) \right\}. \quad (5.5)$$

The lemma below gives choices of $\beta(t, \delta)$ such that $\mathbb{P}_\nu(\mathcal{E}_\delta^c) \leq \delta$. These thresholds satisfy

$$\lim_{\delta \rightarrow 0} \beta(\cdot, \delta) / \log(1/\delta) = 1,$$

and $\beta(t, \cdot) =_{+\infty} \mathcal{O}(\log \log t)$. This result is proven in Lemma 3 of Kone, Jourdan, et al. 2025.

Lemma 5.2.1 (Kone, Jourdan, et al. 2025). *Let $s > 1$, ζ be the Riemann ζ function and $\overline{W}_{-1}(x) := -W_{-1}(-e^{-x})$ for all $x \geq 1$, where W_{-1} is the negative branch of the Lambert W function. Let \mathcal{E}_δ be defined as in (5.5). Then, we have $\mathbb{P}_\nu(\mathcal{E}_\delta^c) \leq \delta$ by taking*

$$\beta(t, \delta) = \frac{dK}{2} \overline{W}_{-1} \left(\frac{2}{dK} \log \frac{e^{Ks} \zeta(s)^K}{\delta} + \frac{2s}{d} \log \left(1 + \frac{d}{2s} \log \frac{t}{K} \right) + 1 \right).$$

We can now state the calibration of M and c to ensure δ -correctness of the PS stopping rule.

Lemma 5.2.2. *Let $\delta \in (0, 1)$, $s > 1$, and let ζ denote the Riemann zeta function. Define*

$$r(\delta, n) := \left(\frac{1}{\sqrt{2\pi}} R \left(\sqrt{\frac{2}{n} \log(1/\delta)} \right) \right)^n,$$

where R is the Mills ratio of $\mathcal{N}(0, 1)$. Regardless of the sampling rule, the PS stopping rule with parameters

$$c(t, \delta) = \frac{\beta(t, \delta/2)}{\log(1/\delta)}, \quad M(t, \delta) = \left\lceil \frac{\log(2t^s \zeta(s)/\delta)}{\delta \cdot q(t, \delta)} \right\rceil,$$

ensures δ -correctness for any $\delta \in (0, 1)$, where $q(t, \delta)$ is defined as follows:

- If Σ is diagonal:

$$q(t, \delta) = \min \left\{ r(\delta, d), r(\delta, d + |\widehat{S}_{t+1}|) \right\}.$$

- If Σ is non-diagonal:

$$q(t, \delta) = \det(\Sigma \bar{\sigma})^{-1/2} \min \left\{ r \left(\delta^{\frac{d_\Sigma}{d}}, d \right), r \left(\delta^{\frac{d_\Sigma + |\widehat{S}_{t+1}|}{d + |\widehat{S}_{t+1}|}}, d + |\widehat{S}_{t+1}| \right) \right\},$$

where $\bar{\sigma} = \|\Sigma^{-1}\|$ and $d_\Sigma = \|1_d\|_{(\bar{\sigma}\Sigma)^{-1}}^2$.

Moreover, these choices satisfy

$$\limsup_{\delta \rightarrow 0} \frac{c(t, \delta) \log M(t, \delta)}{\log(1/\delta)} \leq 1.$$

5.2.2 Game-based sampling rule

We introduce a game-based sampling rule (Degenne & W. Koolen 2019) using a PS min learner. By combining a max learner playing w_t and a min learner playing λ_t , it yields a saddle-point algorithm which approximates $T^*(\nu)^{-1}$.

Max learner. We opt for AdaHedge (De Rooij et al. 2014) as the max learner. We add forced exploration by mixing the played allocation w_t with an exploration allocation $w_{\text{exp}} \in \Delta_K$ with positive weights. Formally, we pull $A_t \sim \tilde{w}_t$ with $\tilde{w}_t := (1 - \gamma_t)w_t + \gamma_t w_{\text{exp}}$ where $\gamma_t = 1/t^\alpha$ with $\alpha \in (0, 1)$. Forced exploration ensures that \widehat{S}_t converges to $S^*(\mu)$ despite the initial fluctuation of the min learning space $\text{Alt}(\widehat{S}_t)$. The cost of forced-exploration is only logarithmic for α close to 1.

Min learner. As in [Z. Li et al. 2024](#), we propose a min learner based on posterior samples. The PS min learner draws independent realizations from a posterior distribution inflated by some parameter $\eta_t > 0$ until one disagrees with \widehat{S}_t . The parameter η_t will act as the learning rate as the posterior distribution is interpreted as the sampling distribution of continuous exponential weights for the quadratic loss.

Formally, let $(v_t^m)_m$ be *i.i.d.* draws from Π_t . For all $m \geq 1$, let $\lambda_{t,k}^m := \widehat{\mu}_{t,k} + \eta_t^{-1/2} v_{t,k}^m \forall k$. Then,

$$\lambda_t := \lambda_t^{m_t} \text{ with } m_t = \inf\{m \mid \lambda_t^m \in \mathcal{I}_t^K \cap \text{Alt}(\widehat{S}_t)\}, \quad (5.6)$$

where $\mathcal{I}_t^K \subset \mathcal{I}^K$ is a bounded set and we recall that $\Pi_t := \bigotimes_{k=1}^K \mathcal{N}(0_d, \Sigma/N_{t,k})$. Intuitively, $\lambda_t^{m_t}$ is a randomized approximation of the best-response oracle $\lambda^{\text{BR}}(N_t)$ as in (5.3) (projected onto \mathcal{I}_t^K).

We re-use the realizations $(v_t^m)_m$ in the PS min learner of our sampling rule with a different inflation parameter. Equivalently, we could draw λ_t from the distribution $\bigotimes_{k=1}^K \mathcal{N}(\widehat{\mu}_{t,k}, \eta_t^{-1} \Sigma/N_{t,k})$ truncated to $\mathcal{I}_t^K \cap \text{Alt}(\widehat{S}_t)$. The inflation η_t is chosen as a function of an upper bound on the magnitude of the stochastic loss, and \mathcal{I}_t^K is defined to control the magnitude of this loss. We refer the reader to [Kone, Jourdan, et al. 2025](#) for more details on how η_t and \mathcal{I}_t^K are precisely defined sequentially.

Computational cost. At time $t < \tau_\delta^{\text{PS}}$, we sequentially re-use the realizations $(v_t^m)_{m \in [m_{t,\delta}]}$ drawn by the PS stopping rule with a different inflation parameter. When $m_{t,\delta} < m_t$, we draw fresh realizations sequentially. Despite the different inflations, m_t is expected to be close to $m_{t,\delta}$, *i.e.*, $\mathcal{O}(M(t, \delta))$ (see Figure 5.6). The computational cost of (5.6) is lower than a best-response oracle to (5.2).

5.3 Main theoretical results

We show that under a generic condition on (c, M) , which is satisfied in Lemma 5.2.2, the expected sample complexity of **PSIPS** is asymptotically optimal (Theorem 5.3.1). From a Bayesian perspective, the posterior probability of misidentifying the Pareto set decays at an asymptotically tight exponential rate (Theorem 5.3.2).

Sample complexity. Theorem 5.3.1 shows that **PSIPS** is asymptotically optimal when using suitable choices of budget and inflation. In comparison, known upper bounds on $\mathbb{E}_\nu[\tau_\delta]$ were either only worst-case optimal (see Chapter 3), or restricted to the setting with independent objectives ([Crepon et al. 2024](#)).

Theorem 5.3.1. *For any budget M and inflation c such that $\limsup_{\delta \rightarrow 0} \frac{c(t,\delta) \log M(t,\delta)}{\log(1/\delta)} \leq 1$, **PSIPS** satisfies that*

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_\nu[\tau_\delta^{\text{PS}}]}{\log \frac{1}{\delta}} \leq T^*(\nu) \quad \text{and} \quad \mathbb{P} \left(\limsup_{\delta \rightarrow 0} \frac{\tau_\delta^{\text{PS}}}{\log \frac{1}{\delta}} \leq T^*(\nu) \right) = 1, \quad \forall \nu \in \mathcal{D}^K.$$

Lemma 5.2.2 gives choices of (c, M) satisfying the above condition and ensuring δ -correctness. Theorem 5.3.1 also holds on the set $\tilde{\mathcal{D}}$ of multivariate Σ -subgaussian, where $\kappa \in \tilde{\mathcal{D}}$ with mean μ implies that $\mathbb{E}_{X \sim \kappa}[e^{u^\top(X-\mu)}] \leq e^{\frac{1}{2}u^\top \Sigma u}$ for all $u \in \mathbb{R}^d$ (Degenne & Perchet 2016). This includes distributions with Bernoulli marginals and distributions with bounded support. Optimality is only achieved for multivariate Gaussian distributions.

Posterior probability of misidentification. Theorem 5.3.2 shows that the posterior probability that PSIPS misidentifies the Pareto set decays exponentially fast (as a function of t) with a rate $1/T^*(\nu)$, which is shown to be asymptotically optimal by a lower bound.

In comparison, similar Bayesian guarantees in the literature were restricted to BAI (e.g., TTTS in Russo 2016 and PEPS in Z. Li et al. 2024).

Theorem 5.3.2. *Let $\tilde{\Pi}_t := \bigotimes_{k=1}^K \mathcal{N}(\hat{\mu}_{t,k}, \Sigma/N_{t,k})$ be the posterior distribution under a flat Gaussian prior (without inflation). For any bandit instance $\nu \in \mathcal{D}^K$, it holds with probability 1 that*

$$\limsup_{t \rightarrow +\infty} -\frac{1}{t} \log \mathbb{P}_{\tilde{\Pi}_t | \mathcal{H}_t}(\text{Alt}(\mathcal{S}^*)) \leq T^*(\nu)^{-1},$$

for any algorithm, and PSIPS almost surely satisfies that

$$\liminf_{t \rightarrow +\infty} -\frac{1}{t} \log \mathbb{P}_{\tilde{\Pi}_t | \mathcal{H}_t}(\text{Alt}(\mathcal{S}^*)) \geq T^*(\nu)^{-1}.$$

5.3.1 Sketch of proofs

Proof sketch of Lemma 5.2.2

To prove Lemma 5.2.2 we first establish the following result.

Lemma 5.3.3. *For all c, M , the PS stopping rule satisfies*

$$\mathbb{P}_\nu(\tau_\delta^{\text{PS}} < +\infty, \hat{S}_{\tau_\delta^{\text{PS}}} \neq \mathcal{S}^*) \leq \delta/2 + \mathbb{E} \left[\sum_{t \geq 1} \mathbb{1}(\mathcal{E}_{t, \delta/2}) \mathbb{1}(\hat{S}_t \neq \mathcal{S}^*) \cdot \exp\left(-M(t-1, \delta) \mathbb{P}_\nu\left(\theta_t \in \text{Alt}(\hat{S}_t) | \mathcal{H}_t\right)\right) \right]. \quad (5.7)$$

Proof. Using the definition of the PS stopping rule, we obtain

$$\begin{aligned}
\mathbb{P}_\nu \left(\tau_\delta^{\text{PS}} < \infty, \widehat{S}_{\tau_\delta^{\text{PS}}} \neq \mathcal{S}^* \right) &\leq \delta/2 + \mathbb{P}_\nu \left(\mathcal{E}_{\delta/2} \cap \{ \tau_\delta^{\text{PS}} < \infty, \widehat{S}_{\tau_\delta^{\text{PS}}} \neq \mathcal{S}^* \} \right) \\
&\leq \delta/2 + \sum_{t \geq 1} \mathbb{P}_\nu \left(\mathcal{E}_{t, \delta/2} \cap \{ \tau_\delta^{\text{PS}} = t, \widehat{S}_{\tau_\delta^{\text{PS}}} \neq \mathcal{S}^* \} \right) \\
&= \delta/2 + \mathbb{E} \left[\sum_{t \geq 1} \mathbb{1}(\mathcal{E}_{t, \delta/2}) \mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \mathbb{P}_\nu(\tau_\delta^{\text{PS}} = t | \mathcal{H}_t) \right] \\
&\leq \delta/2 + \mathbb{E} \left[\sum_{t \geq 1} \mathbb{1}(\mathcal{E}_{t, \delta/2}) \mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \cdot \right. \\
&\quad \left. \mathbb{P}_\nu \left(\forall m \leq M(t-1, \delta), \theta_t^m \notin \text{Alt}(\widehat{S}_t) | \mathcal{H}_t \right) \right] \\
&= \delta/2 + \mathbb{E} \left[\sum_{t \geq 1} \mathbb{1}(\mathcal{E}_{t, \delta/2}) \mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \cdot \right. \\
&\quad \left. \mathbb{P}_\nu \left(\theta_t \notin \text{Alt}(\widehat{S}_t) | \mathcal{H}_t \right)^{M(t-1, \delta)} \right] \\
&= \delta/2 + \mathbb{E} \left[\sum_{t \geq 1} \mathbb{1}(\mathcal{E}_{t, \delta/2}) \mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \cdot \right. \\
&\quad \left. \left(1 - \mathbb{P}_\nu \left(\theta_t \in \text{Alt}(\widehat{S}_t) | \mathcal{H}_t \right) \right)^{M(t-1, \delta)} \right]
\end{aligned}$$

which follows since the samples $\theta_t, \theta_t^1, \dots, \theta_t^m$ are *i.i.d.* given \mathcal{H}_t . Further recalling that $1 - x \leq \exp(-x)$, it follows

$$\begin{aligned}
\mathbb{P}_\nu \left(\tau_\delta^{\text{PS}} < \infty, \widehat{S}_{\tau_\delta^{\text{PS}}} \neq \mathcal{S}^* \right) &\leq \delta/2 + \mathbb{E} \left[\sum_{t \geq 1} \mathbb{1}(\mathcal{E}_{t, \delta/2}) \mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \cdot \right. \\
&\quad \left. \exp \left(-M(t-1, \delta) \mathbb{P}_\nu \left(\theta_t \in \text{Alt}(\widehat{S}_t) | \mathcal{H}_t \right) \right) \right]. \tag{5.8}
\end{aligned}$$

□

Lemma 5.3.3 shows that choosing c, M such that

$$\mathbb{E} \left[\sum_{t \geq 1} \mathbb{1}(\mathcal{E}_{t, \delta/2}) \mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \exp \left(-M(t-1, \delta) \mathbb{P}_\nu \left(\theta_t \in \text{Alt}(\widehat{S}_t) | \mathcal{H}_t \right) \right) \right] \leq \delta/2$$

will ensure the correctness of the PS stopping rule. However, it is not possible to pick $M(t-1, \delta) \propto \mathbb{P}_\nu(\theta_t \in \text{Alt}(\widehat{S}_t) | \mathcal{H}_t)^{-1}$ as the latter quantity is intractable due to multidimensional integration over a non-trivial domain. Instead, we will derive tractable lower bounds on $\mathbb{1}(\mathcal{E}_{t, \delta/2}) \mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \mathbb{P}_\nu(\theta_t \in \text{Alt}(\widehat{S}_t) | \mathcal{H}_t)$. This will ensure δ -correctness of the PS

stopping rule, regardless of the sampling rule. Let $\widehat{\Pi}_t := \bigotimes_{k=1}^K \mathcal{N}(\widehat{\mu}_{t,k}, c(t-1, \delta)\Sigma/N_{t,k})$. From the above display, it is sufficient to lower bound

$$\mathbb{1}(\mathcal{E}_{t,\delta/2} \cap \{\widehat{\mathcal{S}}_t \neq \mathcal{S}^*\}) \exp\left(-M(t-1, \delta)\mathbb{P}_{\widehat{\Pi}_t|\mathcal{H}_t}(\text{Alt}(\widehat{\mathcal{S}}_t))\right)$$

by $\frac{\delta}{2\zeta(s)t^s}$ since their sum is smaller than $\delta/2$. Therefore, we should lower bound $\mathbb{P}_{\widehat{\Pi}_t|\mathcal{H}_t}(\text{Alt}(\widehat{\mathcal{S}}_t))$ under the event $\mathcal{E}_{t,\delta/2} \cap \{\widehat{\mathcal{S}}_t \neq \mathcal{S}^*\}$ to conclude the proof.

Calibration of $M(t, \delta)$. Having $\widehat{\mathcal{S}}_t \neq \mathcal{S}^*$ implies that $\mu \in \text{Alt}(\widehat{\mathcal{S}}_t)$, see (5.4). Either (1) there exists $(i, j) \in \widehat{\mathcal{S}}_t$ with $i \neq j$ such that $\mu_i(c) < \mu_j(c)$ for all $c \in [d]$, and we define $\mathcal{J}_t = \{j\}$. Or (2) there exists $i \notin \widehat{\mathcal{S}}_t$ and $c \in [d]^{|\widehat{\mathcal{S}}_t|}$ such that $\mu_j(c_j) < \mu_i(c_j)$ for all $j \in \widehat{\mathcal{S}}_t$, and $\mathcal{J}_t = \widehat{\mathcal{S}}_t$.

We can show that

$$\mathbb{1}(\widehat{\mathcal{S}}_t \neq \mathcal{S}^*) \mathbb{P}_{\widehat{\Pi}_t|\mathcal{H}_t}(\text{Alt}(\widehat{\mathcal{S}}_t)) \geq \prod_{k \in \{i\} \cup \mathcal{J}_t} \mathbb{P}_{X \sim \mathcal{N}(0_d, \Sigma)}\left(X > \sqrt{\frac{N_{t,k}}{c(t-1, \delta)}}(\mu_k - \widehat{\mu}_{t,k})\right),$$

where $X > x$ denotes $X(c) > x(c)$ for all $c \in [d]$.

When Σ is diagonal, we have $\mathbb{P}_{X \sim \mathcal{N}(0_d, \Sigma)}(X > x) = \prod_{c \in [d]} \mathbb{P}_{X \sim \mathcal{N}(0, \Sigma_{c,c})}(X > x(c))$. When Σ is not diagonal, we derive a lower bound on $\mathbb{P}_{X \sim \mathcal{N}(0_d, \Sigma)}(X > x)$ for any vector $x \in \mathbb{R}^d$ (Lemma 5.5.3).

To further lower bound those quantities, we introduce the ratio of the tail distribution to the Gaussian density function, known as the Mills ratio (Mills 1926). Under the event $\mathcal{E}_{t,\delta/2}$, we have $\sum_{k=1}^K N_{t,k} \|\mu_k - \widehat{\mu}_{t,k}\|_{\Sigma^{-1}}^2 \leq 2\beta(t-1, \delta/2)$. Using the fact that the Mills ratio is non-increasing and log-convex (Baricz 2008), we obtain

$$\mathbb{1}(\mathcal{E}_{t,\delta/2} \cap \{\widehat{\mathcal{S}}_t \neq \mathcal{S}^*\}) \mathbb{P}_{\widehat{\Pi}_t|\mathcal{H}_t}(\text{Alt}(\widehat{\mathcal{S}}_t)) \geq \delta q(t-1, \delta)$$

by taking $c(t, \delta) = \beta(t, \delta/2)/\log(1/\delta)$. This concludes the proof by considering $M(t, \delta) = \left\lceil \frac{\log(2t^s \zeta(s)/\delta)}{\delta q(t, \delta)} \right\rceil$.

In contrast, the δ -correctness of the GLR stopping rule is obtained quite simply by concentration. Indeed, we have $\text{GLR}(t) \leq \sum_{k=1}^K \frac{N_{t,k}}{2} \|\widehat{\mu}_{t,k} - \mu_k\|_{\Sigma^{-1}}^2 \leq \beta(t-1, \delta)$ under $\mathcal{E}_{t,\delta} \cap \{\widehat{\mathcal{S}}_t \neq \mathcal{S}^*\}$, hence $\mathbb{P}_\nu(\mathcal{E}_\delta \cap \{\tau_\delta^{\text{GLR}} < +\infty, \widehat{\mathcal{S}}_{\tau_\delta^{\text{GLR}}} \neq \mathcal{S}^*\}) = 0$. Under a similar event, proving the δ -correctness of the PS stopping rule requires controlling the randomness of $\widehat{\Pi}_t|\mathcal{H}_t$ by deriving anti-concentration results on $\mathbb{P}_{\widehat{\Pi}_t|\mathcal{H}_t}(\text{Alt}(\widehat{\mathcal{S}}_t))$.

These results are derived in Section 5.5.1 where Lemma 5.2.2 is proven.

Proof sketch of Theorem 5.3.1

Studying the expected sample complexity of an algorithm using the PS stopping rule requires controlling the randomness of $\widehat{\Pi}_t|\mathcal{H}_t$ by deriving concentration results on $\mathbb{P}_{\widehat{\Pi}_t|\mathcal{H}_t}(\text{Alt}(\widehat{\mathcal{S}}_t))$

since a direct union bound yields that

$$\mathbb{P}_\nu(\tau_\delta^{\text{PS}} > t) \leq M(t-1, \delta) \mathbb{E}_\nu[\mathbb{P}_{\hat{\Pi}_t|\mathcal{H}_t}(\text{Alt}(\hat{S}_t))].$$

Now, we invoke the following lemma to write $\text{Alt}(\hat{S}_t)$ as a finite union of convex sets.

Lemma 5.3.4 (Convex decomposition of alternative set; [Crepon et al. 2024](#); [Kone, Jourdan, et al. 2025](#)). *For any parameter $\mu \in \mathcal{I}^K$ with $p := |S^*(\mu)|$, the alternative set $\text{Alt}(\mu)$ can be decomposed as the union of $n := p(p-1) + (K-p)d^p$ convex subsets.*

We then combine this lemma with the result below on Gaussian concentration on convex sets.

Lemma 5.3.5 (Gaussian tail bound for convex sets; [D. Lu & W. V. Li 2009](#)). *Let $X \sim \mathcal{N}(\mu, A)$ with $\mu \in \mathbb{R}^n$ and covariance matrix A . For any convex set $C \subset \mathbb{R}^n$,*

$$\mathbb{P}(X \in C) \leq \frac{1}{2} \exp\left(-\frac{1}{2} \inf_{\lambda \in C} \|\lambda - \mu\|_{A^{-1}}^2\right).$$

As $\hat{\Pi}_t|\mathcal{H}_t$ is a Kd -sized Gaussian vector with block diagonal covariance, this result relates $\mathbb{P}_{\hat{\Pi}_t|\mathcal{H}_t}(\text{Alt}(\hat{S}_t))$ to the GLR statistic defined in (5.2).

Lemma 5.3.4 proves that there exists n_t and convex sets C_1, \dots, C_{n_t} such that $\text{Alt}(\hat{S}_t) = \cup_{i \in [n_t]} C_i$, and $\hat{\Pi}_t|\mathcal{H}_t = \otimes_{k=1}^K \mathcal{N}(\hat{\mu}_{t,k}, c(t-1, \delta)\Sigma/N_{t,k})$, so, thanks to Lemma 5.3.5, and convexity of C_i , it follows that

$$\mathbb{P}(\theta_t \in C_i|\mathcal{H}_t) \leq \frac{1}{2} \exp\left(-\inf_{\theta \in C_i} \left[\sum_{k=1}^K N_{t,k} \frac{\|\hat{\mu}_{t,k} - \theta_k\|_{\Sigma^{-1}}^2}{2c(t-1, \delta)} \right]\right).$$

Therefore, by union bound we have

$$\begin{aligned} \mathbb{P}_{\hat{\Pi}_t|\mathcal{H}_t}(\text{Alt}(\hat{S}_t)) &:= \mathbb{P}(\theta_t \in \cup_{i \in [n_t]}(C_i)|\mathcal{H}_t) \\ &\leq \frac{1}{2} \sum_{i \in [n_t]} \exp\left(-\inf_{\lambda \in C_i} \left[\sum_{k=1}^K N_{t,k} \frac{\|\hat{\mu}_{t,k} - \lambda_k\|_{\Sigma^{-1}}^2}{2c(t-1, \delta)} \right]\right) \\ &\leq \frac{n_t}{2} \exp\left(-\inf_{\lambda \in \text{Alt}(\hat{S}_t)} \left[\sum_{k=1}^K N_{t,k} \frac{\|\hat{\mu}_{t,k} - \lambda_k\|_{\Sigma^{-1}}^2}{2c(t-1, \delta)} \right]\right) \\ &= \frac{n_t}{2} \exp\left(-\frac{\text{GLR}(t)}{c(t-1, \delta)}\right) \end{aligned}$$

and from Lemma 5.3.4, $n_t = p_t(p_t-1) + (K-p_t)d^{p_t}$, with $p_t := |\hat{S}_t|$. From the above displays, we have

$$\mathbb{P}_{\hat{\Pi}_t|\mathcal{H}_t}(\text{Alt}(\hat{S}_t)) \leq \alpha_0 \exp(-\text{GLR}(t)/c(t-1, \delta)),$$

where $\alpha_0 \leq K(K + d^K)/2$. Now the event $\{\tau_\delta^{\text{PS}} > t\}$ is linked to the GLR statistic; it suffices to lower bound $\text{GLR}(t)$ to conclude the proof.

Studying the saddle-point convergence of our algorithm shows that there exist $(\Xi_t)_{t \geq 1}$ and $T_0 \in \mathbb{N}$ such that for all $t > T_0$, when Ξ_t holds, $\text{GLR}(t) \geq (t-1)/T^*(\nu) - f(t)$ where $f(t) = {}_{+\infty} o(t)$ and $\sum_{t \in \mathbb{N}} \mathbb{P}(\mathcal{E}_t^c) = C_0 < +\infty$. Then, we have

$$\mathbb{E}_\nu[\tau_\delta^{\text{PS}}] \leq T_0 + C_0 + \sum_{t > T_0} \alpha_0 \exp\left(-\frac{t-1}{T^*(\nu)c(t-1, \delta)} + \frac{f(t)}{c(t-1, \delta)} + \log M(t-1, \delta)\right).$$

Using that $\limsup_{\delta \rightarrow 0} \frac{c(t, \delta) \log M(t, \delta)}{\log(1/\delta)} \leq 1$, a direct inversion of the above sum concludes the proof.

The complete proof is given at the end of the chapter.

5.4 Numerical study and discussion

Experimental setup. We evaluate the performance of PSIPS on real-world-inspired and synthetic instances. As benchmarks, we consider the following PSI algorithms: APE (Chapter 3), the gradient-based algorithm of Crepon et al. 2024 denoted as GAPSI, the *oracle* algorithm which samples arms according to the optimal weights, *i.e.*, $A_t \sim w^*(\nu)$, and round-robin *uniform* sampling (RR).

APE relies on confidence bounds to pull an arm at each round and to define its stopping rule. While APE can be extended to the structured setting, it does not exploit the correlation between objectives as its confidence intervals only consider the marginals' variance. GAPSI computes a supergradient of the GLR (5.2) and uses the GLR stopping rule. Likewise, GAPSI has no efficient implementation or guarantees for correlated objectives.

Since GAPSI is computationally expensive, we exclude it from the experiment in the linear setting (due to the large number of arms), and we run it as heuristic in experiments with correlated objectives.

In our experiments, we use the heuristics $c(t, \delta) = 1 + \frac{\log \log t}{\log(1/\delta)}$ and $M(t, \delta) = \frac{1}{\delta} \log \frac{t}{\delta}$. For a fair comparison, APE uses $\beta(t, \delta) = \log(1/\delta) + \log \log t$ in its bounds, hence improving its performance. We report the averaged results over 500 independent runs with boxplots or shaded areas for standard deviation. The observed empirical error is an order of magnitude lower than δ .

CovBoost trial. As in prior PSI work, we use the dataset from Munro et al. 2021 to simulate a bandit instance for PSI. The dataset contains average responses for patient cohorts in a COVID-19 trial, including 20 vaccine strategies. Among the indicators recorded to measure the efficacy of each strategy, following Kone, Kaufmann, et al. 2023; Crepon et al. 2024, we keep three of them: titers of neutralizing antibodies and immunoglobulin G and the wild-type cellular response (see Table 7.1). The ideal vaccine strategy would maximize all three indicators, but the average response for the considered metrics reveals a Pareto set of two vaccine strategies in the trial.

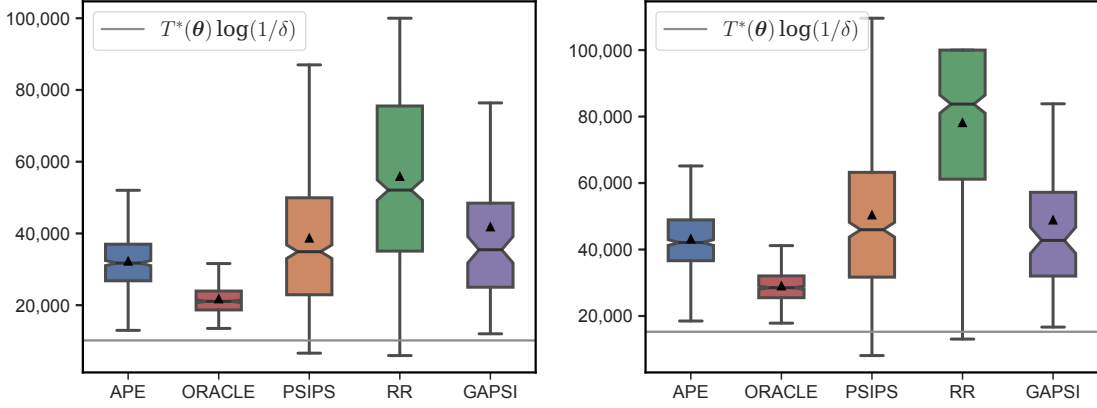


Figure 5.1: Empirical stopping times on the covid19 experiment with $\delta = 0.01$ (left) and $\delta = 0.001$ (right).

Figure 5.1 shows that **PSIPS** has a high variability of stopping time, which was expected due to the use of posterior sampling. While **PSIPS** outperforms uniform sampling, it performs on average slightly worse than **APE** and **oracle** on the CovBoost instance. For $\delta = 0.1$, we reported an average sample complexity of 20456 for **PSIPS** compared to 17909 reported for **GAPSI** by [Crepon et al. 2024](#). However, due to the high computational cost of their algorithm, we averaged its performance over 100 runs.

Correlated Objectives. We evaluate the impact of correlated objectives on the performance of **PSIPS**. Consider a 5-arm Gaussian instance with means $\mu_1 = (0.73, 1.20)^\top$, $\mu_2 = (0.45, -0.63)^\top$, $\mu_3 = (0.63, 1.28)^\top$, $\mu_4 = (0.94, 2.31)^\top$, and $\mu_5 = (2.08, 1.48)^\top$. The noise covariance matrix Σ_ρ has unit diagonal and constant off-diagonal entries equal to $\rho \in (-1, 1)$, controlling the correlation between objectives ($\rho > 0$ for positive and $\rho < 0$ for negative correlation). Using an approximate convex solver and the exponentiated gradient algorithm, we estimate the theoretical complexity $T_{\Sigma_\rho}^*(\nu)$ defined in Lemma 5.1.1. As shown in Figure 5.3, negative correlations reduce the theoretical complexity (characteristic time) compared to the uncorrelated case ($\rho = 0$), reflecting the additional information gained from anti-correlated objectives.

Figure 5.2 shows the sample complexity of **PSIPS** versus ρ , including **GAPSI** as heuristic. It reveals that **PSIPS** has a decreasing sample complexity for stronger (negative) correlation between objectives, reducing by up to a factor of 3 relative to the uncorrelated case ($\rho = 0$). Due to its asymptotic optimality, **PSIPS** inherits the properties of $T^*(\nu)$, which appears to decrease with negative ρ on this specific instance. Both **GAPSI** and **APE**'s performance are independent of ρ .

Robustness on random instances. To evaluate the robustness of **PSIPS**, we measure its performance on 250 randomly picked Bernoulli (marginals) and Gaussian instances with $K = 5$, $d = 2$ and $\Sigma = I_2/2$. Figure 5.4 demonstrates the robustness of **PSIPS**, which remains competitive on subgaussian instances.

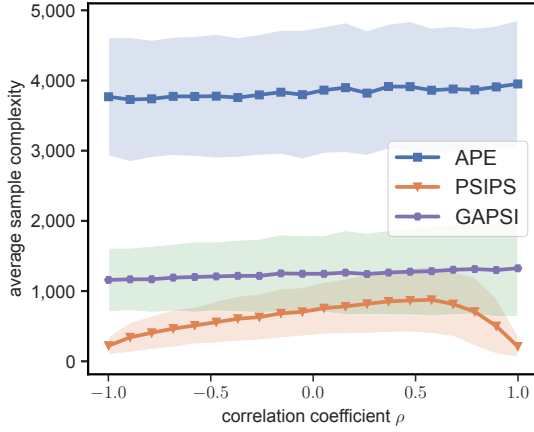


Figure 5.2: Empirical impact of the correlation ρ on the sample complexity.

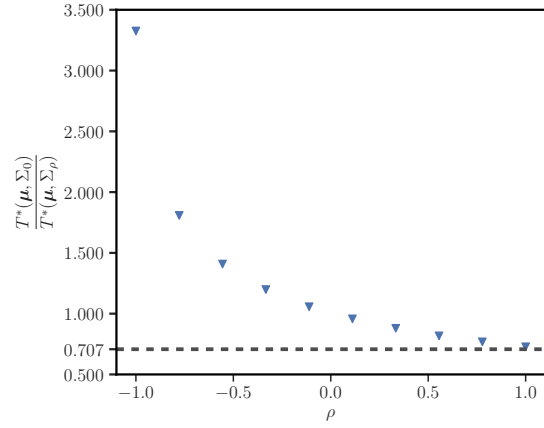


Figure 5.3: Theoretical impact of the correlation coefficient on the complexity T^* .

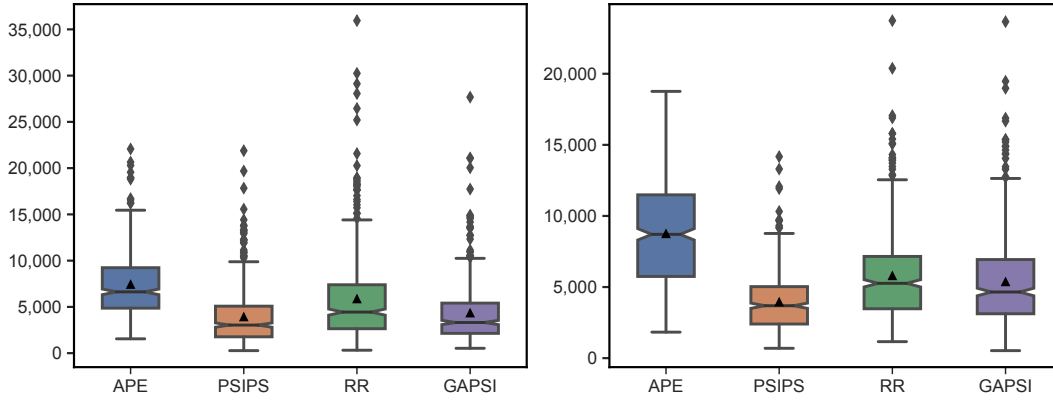


Figure 5.4: Empirical stopping time on random Gaussian (left) and Bernoulli (right) instances.

Computational cost. We evaluate the number of rejection samples used by **PSIPS** on a Gaussian instance with $\Sigma = I_2/2$ defined as $\mu_1 = (1, 1)^\top$ and $\mu_i = R_{\pi/5}\mu_{i-1}$ for $i \in \{2, 3, 4, 5\}$ where $R_{\pi/5}$ is the $\pi/5$ rotation matrix. Without actually stopping, we record both m_t and $m_{t,\delta}$ at each round, averaged over 1000 runs. We also report the cumulative runtime of **PSIPS** and **APE** on the CovBoost experiment.

Figure 5.5 shows that the computational cost of **PSIPS** is comparable to **APE**, which is orders of magnitudes smaller than **GAPSI** (see Figure 4 in [Crepon et al. 2024](#)). Figure 5.6 reveals that, while finding an alternative is initially faster, more rejections are needed before finding an alternative when the posterior concentrates. It is precisely at this time that the PS stopping rule triggers.

Discussion. We proposed the **PSIPS** algorithm for Pareto set identification with correlated objectives. By leveraging posterior sampling in both the stopping and the sampling rules, **PSIPS** is a computationally efficient algorithm that deals with structure and correlation, without the costly oracle calls required by existing algorithms. We show that **PSIPS** is asymptotically optimal both from a frequentist and Bayesian perspective. Moreover, **PSIPS**

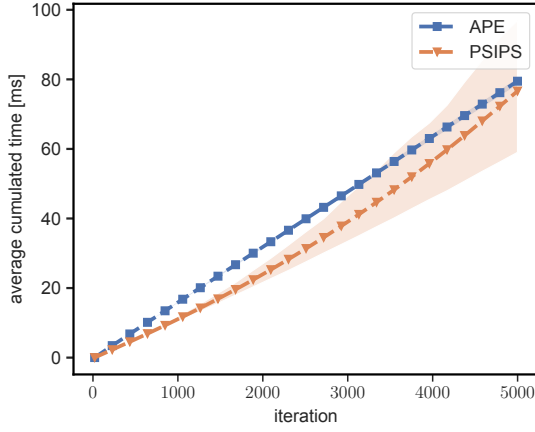


Figure 5.5: Average runtime for the first T iterations in the COVID-19 experiment.

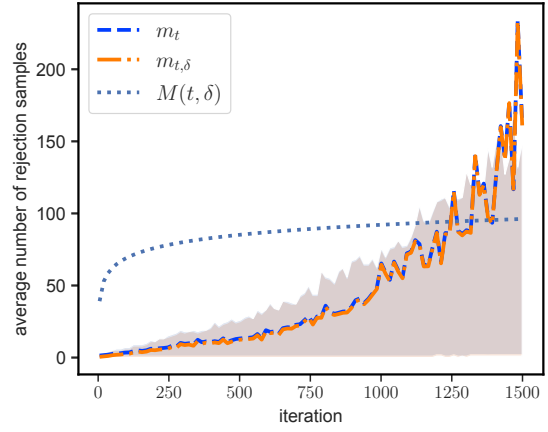


Figure 5.6: Average per-round number of rejection samples $m_t, m_{t,\delta}$.

has competitive empirical performance.

Despite being a theoretically convenient distributional assumption, using Gaussian distributions with known covariance matrix is too restrictive for many applications. Therefore, generalizing our results to other classes of distributions is another interesting direction for future work. For example, when the covariance matrices are unknown and different for each arm, the estimation of the Pareto set is intertwined with estimating the correlation structure.

5.5 Additional proofs

We prove detail additional results in this section. Complementary proofs can be found in [Kone, Jourdan, et al. 2025](#).

5.5.1 Stopping rule

In this section, we prove the correctness of the PS (Posterior Sampling) stopping rule. We recall the definition of τ_δ^{PS} :

$$\tau_\delta^{\text{PS}} := \inf \left\{ t \mid \forall m \in [M(t-1, \delta)], \lambda_t^m \notin \text{Alt}(\widehat{S}_t) \right\},$$

where, conditioned on $\mathcal{H}_t, \theta_t, \theta_t^1, \dots, \theta_t^m$ are *i.i.d.* samples from $\mathcal{N}(\hat{\mu}_{t,k}, c(t-1, \delta)\Sigma/N_{t,k})$. To prove the correctness of this stopping rule, we will have to control the randomness in the posterior, which we do by introducing the concentration event $\mathcal{E}_\delta := \bigcap_{t \in \mathbb{N}} \mathcal{E}_{t,\delta}$ with

$$\mathcal{E}_{t,\delta} := \left\{ \sum_{k=1}^K N_{t,k} \|\mu_k - \hat{\mu}_{t,k}\|_{\Sigma^{-1}}^2 \leq 2\beta(t-1, \delta) \right\},$$

for which, correct values of the thresholds are given in the lemma below.

The following result shows that to prove the δ -correctness, it is sufficient to show some anti-concentration of the posterior under the event \mathcal{E} . $\mathbb{P}_\nu(\mathcal{E}_{\delta/2} \cap \{\tau_\delta^{\text{PS}} < +\infty, \widehat{S}_{\tau_\delta^{\text{PS}}} \neq \mathcal{S}^*\}) \leq \delta/2$.

Next, we calibrate c, M to ensure the correctness of the PS stopping rule. The quantity to control is

$$\mathbb{1}(\mathcal{E}_{t,\delta/2}) \mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \mathbb{P}_\nu(\theta_t \in \text{Alt}(\widehat{S}_t) | \mathcal{H}_\infty) \text{ with } (\theta_t - \widehat{\theta}_t) | \mathcal{H}_t \sim \bigotimes_{k=1}^K \mathcal{N}(0_d, c(t-1, \delta) \Sigma / N_{t,k}),$$

and

$$\begin{aligned} \text{Alt}(\widehat{S}_t) = & \left(\bigcup_{i \neq j \in \widehat{S}_t^2} \{ \lambda := (\lambda_1, \dots, \lambda_K) \in \mathcal{I}^K | \lambda_i \prec \lambda_j \} \right) \cup \\ & \left(\bigcup_{i \notin \widehat{S}_t} \{ \lambda := (\lambda_1, \dots, \lambda_K) \in \mathcal{I}^K | \forall j \in \widehat{S}_t, \lambda_i \not\prec \lambda_j \} \right). \end{aligned}$$

Indeed, recalling that $\text{Alt}(S)$ is the set of parameters for which the Pareto set differs from S . To change the Pareto set, we either add an arm or remove one from it. For a parameter instance $\mu := (\mu_1 \dots \mu_K)^\top$, if $i \in S$ and there exists $j \in S$ such that $\mu_i \prec \mu_j$, then as i is a dominated arm in the instance μ , we will have $S^*(\mu) \neq S$. Similarly, if for some $i \notin S$ it holds that for all $j \in S$, $\mu_i \not\prec \mu_j$, then we cannot have $S^*(\mu) = S$, otherwise, as $i \notin S^*(\mu)$, an arm from $S^*(\mu) = S$ would have dominated i . This is formally shown in Lemma 5.3.4, where it is further proven that

$$\begin{aligned} \text{Alt}(S) = & \left(\bigcup_{i \in S^c} \bigcup_{\bar{d}^i \in [d]^S} \{ \lambda | \forall j \in S, \lambda_i(d^i(\sigma(j))) \geq \lambda_j(d^i(\sigma(j))) \} \right) \cup \\ & \left(\bigcup_{(i,j) \in S^2, i \neq j} \{ \lambda \in \mathcal{I}^K | \lambda_i \prec \lambda_j \} \right) \end{aligned}$$

where σ is any permutation that maps S onto $\{1, \dots, |S|\}$. In the sequel, we let

$$\begin{aligned} \text{Alt}^-(S) & := \bigcup_{(i,j) \in S^2, i \neq j} \{ \lambda | \lambda_i \prec \lambda_j \} \quad \text{and} \\ \text{Alt}^+(S) & := \bigcup_{i \in S^c} \bigcup_{\bar{d}^i \in [d]^S} \{ \lambda | \forall j \in S : \lambda_i(d^i(\sigma(j))) \geq \lambda_j(d^i(\sigma(j))) \}, \end{aligned}$$

so that $\text{Alt}(S) = \text{Alt}^+(S) \cup \text{Alt}^-(S)$. We prove the following lemma.

Lemma 5.5.1. *Letting $Y \sim \mathcal{N}(0_d, \Sigma)$, it holds that*

$$\begin{aligned} \mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \mathbb{P}_{\widehat{\Pi}_t | \mathcal{H}_t}(\text{Alt}(\widehat{S}_t)) &\geq \max \left\{ \max_{i \notin \widehat{S}_t, \bar{d} \in [d]^{|\widehat{S}_t|}} \mathbb{1} \left(\forall j \in \widehat{S}_t, \mu_i(\bar{d}_j^\sigma) \geq \mu_j(\bar{d}_j^\sigma) \right) \right. \\ &\quad \cdot \mathbb{P}_Y \left(Y > \sqrt{\frac{N_{t,i}}{c(t-1, \delta)}} (\mu_i - \hat{\mu}_{t,i}) \right) \cdot \prod_{j \in \widehat{S}_t} \mathbb{P}_{X \sim \mathcal{N}(0,1)} \left(X > \sqrt{\frac{N_{t,j}}{c(t-1, \delta)}} \|\mu_j - \hat{\mu}_{t,j}\|_{\Sigma^{-1}} \right), \\ &\quad \left. \max_{i \neq j \in \widehat{S}_t^2} \mathbb{1}(\mu_i \prec \mu_j) \mathbb{P}_Y \left(Y \geq \sqrt{\left(\frac{1}{N_{t,i}} + \frac{1}{N_{t,j}} \right)^{-1} \frac{1}{c(t-1, \delta)}} ((\hat{\mu}_{t,i} - \hat{\mu}_{t,j}) - (\mu_i - \mu_j)) \right) \right\}. \end{aligned}$$

Proof.

$$\begin{aligned} \mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \mathbb{P}_{\widehat{\Pi}_t | \mathcal{H}_t}(\text{Alt}^-(\widehat{S}_t)) &\geq \mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \max_{(i,j) \in \widehat{S}_t, i \neq j} \mathbb{P}_{\widehat{\Pi}_t | \mathcal{H}_t}(\{\mu \in \mathcal{I}^K : \mu_i \prec \mu_j\}) \\ &\geq \max_{(i,j) \in \widehat{S}_t, i \neq j} \mathbb{1}(\mu_i \prec \mu_j) \mathbb{P}_{\widehat{\Pi}_t | \mathcal{H}_t}(\{\mu \in \mathcal{I}^K : \mu_i \prec \mu_j\}) \\ &= \max_{(i,j) \in \widehat{S}_t, i \neq j} \mathbb{1}(\mu_i \prec \mu_j) \mathbb{P}(\theta_{t,i} \prec \theta_{t,j} | \mathcal{H}_t), \end{aligned}$$

where we have $\theta_{t,i} | \mathcal{H}_t \sim \mathcal{N}(\hat{\mu}_{t,i}, c(t-1, \delta)\Sigma/N_{t,i})$. Therefore, we have

$$\begin{aligned} \mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \mathbb{P}_{\widehat{\Pi}_t | \mathcal{H}_t}(\text{Alt}^-(\widehat{S}_t)) &\geq \max_{(i,j) \in \widehat{S}_t, i \neq j} \mathbb{1}(\mu_i \prec \mu_j) \mathbb{P}((\theta_{t,i} - \theta_{t,j}) - (\hat{\mu}_{t,i} - \hat{\mu}_{t,j}) \\ &\quad \prec -(\hat{\mu}_{t,i} - \hat{\mu}_{t,j}) | \mathcal{H}_t) \\ &\geq \max_{(i,j) \in \widehat{S}_t, i \neq j} \mathbb{1}(\mu_i \prec \mu_j) \mathbb{P}((\theta_{t,i} - \theta_{t,j}) - (\hat{\mu}_{t,i} - \hat{\mu}_{t,j}) \prec \\ &\quad (\mu_i - \mu_j) - (\hat{\mu}_{t,i} - \hat{\mu}_{t,j}) | \mathcal{H}_t) \end{aligned}$$

then, observe that $(\theta_{t,i} - \theta_{t,j}) | \mathcal{H}_t \sim \mathcal{N}(\hat{\mu}_{t,i} - \hat{\mu}_{t,j}, (N_{t,i}^{-1} + N_{t,j}^{-1})c(t-1, \delta)\Sigma)$, so that letting $Y \sim \mathcal{N}(0_d, \Sigma)$, we have

$$\begin{aligned} &\mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \mathbb{P}_{\widehat{\Pi}_t | \mathcal{H}_t}(\text{Alt}^-(\widehat{S}_t)) \\ &\geq \max_{i \neq j \in \widehat{S}_t^2} \mathbb{1}(\mu_i \prec \mu_j) \mathbb{P}_{Y \sim \mathcal{N}(0_d, \Sigma)} \left(Y \prec \sqrt{\left(\frac{1}{N_{t,i}} + \frac{1}{N_{t,j}} \right)^{-1} \frac{1}{c(t-1, \delta)}} ((\mu_i - \mu_j) - (\hat{\mu}_{t,i} - \hat{\mu}_{t,j})) \right) \\ &\geq \max_{i \neq j \in \widehat{S}_t^2} \mathbb{1}(\mu_i \prec \mu_j) \mathbb{P}_{Y \sim \mathcal{N}(0_d, \Sigma)} \left(Y \geq \sqrt{\left(\frac{1}{N_{t,i}} + \frac{1}{N_{t,j}} \right)^{-1} \frac{1}{c(t-1, \delta)}} ((\hat{\mu}_{t,i} - \hat{\mu}_{t,j}) - (\mu_i - \mu_j)) \right). \end{aligned}$$

We now prove a similar result on $\text{Alt}^+(\widehat{S}_t)$. We have

$$\begin{aligned}
 & \mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \mathbb{P}_{\widehat{\Pi}_t | \mathcal{H}_t}(\text{Alt}^+(\widehat{S}_t)) \\
 & \geq \max_{i \notin \widehat{S}_t, \bar{d} \in [d]^{|\widehat{S}_t|}} \mathbb{1}(\forall j \in \widehat{S}_t, \mu_i(\bar{d}_j^\sigma) \geq \mu_j(\bar{d}_j^\sigma)) \mathbb{P}(\forall j \in \widehat{S}_t, \theta_{t,i}(\bar{d}_j^\sigma) \geq \theta_{t,j}(\bar{d}_j^\sigma) | \mathcal{H}_t) \\
 & \geq \max_{i \notin \widehat{S}_t, \bar{d} \in [d]^{|\widehat{S}_t|}} \mathbb{1}(\forall j \in \widehat{S}_t, \mu_i(\bar{d}_j^\sigma) \geq \mu_j(\bar{d}_j^\sigma)) \mathbb{P}(\forall j \in \widehat{S}_t, (\theta_{t,i} - \mu_i)(\bar{d}_j^\sigma) \geq (\theta_{t,j} - \mu_j)(\bar{d}_j^\sigma) | \mathcal{H}_t) \\
 & \geq \max_{i \notin \widehat{S}_t, \bar{d} \in [d]^{|\widehat{S}_t|}} \mathbb{1}(\forall j \in \widehat{S}_t, \mu_i(\bar{d}_j^\sigma) \geq \mu_j(\bar{d}_j^\sigma)) \mathbb{P}((\theta_{t,i} - \mu_i) > 0_d | \mathcal{H}_t) \prod_{j \in \widehat{S}_t} \mathbb{P}_\nu((\theta_{t,j} - \mu_j)(\bar{d}_j^\sigma) > 0 | \mathcal{H}_t)
 \end{aligned}$$

then, noting again that $(\theta_{t,i} - \hat{\mu}_{t,i}) | \mathcal{H}_t \sim \mathcal{N}(0, c(t-1, \delta)\Sigma/N_{t,i})$, it follows that

$$\begin{aligned}
 & \mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \mathbb{P}_{\widehat{\Pi}_t | \mathcal{H}_t}(\text{Alt}^+(\widehat{S}_t)) \geq \max_{i \notin \widehat{S}_t, \bar{d} \in [d]^{|\widehat{S}_t|}} \mathbb{1}(\forall j \in \widehat{S}_t, \mu_i(\bar{d}_j^\sigma) \geq \mu_j(\bar{d}_j^\sigma)) \cdot \\
 & \mathbb{P}_Y \left(Y > \sqrt{\frac{N_{t,i}}{c(t-1, \delta)}} (\mu_i - \hat{\mu}_{t,i}) \right) \prod_{j \in \widehat{S}_t} \mathbb{P}_\nu((\theta_{t,j} - \mu_j)(\bar{d}_j^\sigma) > 0 | \mathcal{H}_t). \tag{5.9}
 \end{aligned}$$

Further noting that by Cauchy-Schwarz inequality,

$$(\hat{\mu}_{t,j} - \mu_j)(\bar{d}_j^\sigma) \leq \|\tilde{e}_{\bar{d}_j^\sigma}\|_\Sigma \|\mu_j - \hat{\mu}_{t,j}\|_{\Sigma^{-1}},$$

it follows for any $j \in \widehat{S}_t$ that

$$\begin{aligned}
 \mathbb{P}((\theta_{t,j} - \mu_j)(\bar{d}_j^\sigma) > 0 | \mathcal{H}_t) &= \mathbb{P}((\theta_{t,j} - \hat{\mu}_{t,j})(\bar{d}_j^\sigma) > (\mu_j - \hat{\mu}_{t,j})(\bar{d}_j^\sigma) | \mathcal{H}_t) \\
 &\geq \mathbb{P}((\theta_{t,j} - \hat{\mu}_{t,j})(\bar{d}_j^\sigma) > \|\tilde{e}_{\bar{d}_j^\sigma}\|_\Sigma \|\mu_j - \hat{\mu}_{t,j}\|_{\Sigma^{-1}} | \mathcal{H}_t) \\
 &= \mathbb{P}_{X \sim \mathcal{N}(0,1)} \left(X > \sqrt{\frac{N_{t,j}}{c(t-1, \delta)}} \|\mu_j - \hat{\mu}_{t,j}\|_{\Sigma^{-1}} \right),
 \end{aligned}$$

which follows from $(\theta_{t,j} - \hat{\mu}_{t,j})(\bar{d}_j^\sigma) | \mathcal{H}_t \sim \mathcal{N}(0, c(t-1, \delta)\|\tilde{e}_{\bar{d}_j^\sigma}\|_\Sigma^2/N_{t,\delta})$. Therefore, combining the last display with (5.9) yields

$$\begin{aligned}
 & \mathbb{1}(\widehat{S}_t \neq \mathcal{S}^*) \mathbb{P}_{\widehat{\Pi}_t | \mathcal{H}_t}(\text{Alt}^+(\widehat{S}_t)) \geq \max_{i \notin \widehat{S}_t, \bar{d} \in [d]^{|\widehat{S}_t|}} \mathbb{1}(\forall j \in \widehat{S}_t, \mu_i(\bar{d}_j^\sigma) \geq \mu_j(\bar{d}_j^\sigma)) \\
 & \cdot \mathbb{P}_Y \left(Y > \sqrt{\frac{N_{t,i}}{c(t-1, \delta)}} (\mu_i - \hat{\mu}_{t,i}) \right) \cdot \prod_{j \in \widehat{S}_t} \mathbb{P}_{X \sim \mathcal{N}(0,1)} \left(X > \sqrt{\frac{N_{t,j}}{c(t-1, \delta)}} \|\mu_j - \hat{\mu}_{t,j}\|_{\Sigma^{-1}} \right),
 \end{aligned}$$

which completes the proof. \square

Special Case: PSI with Independent Marginals In this section, we specialize Lemma 5.5.1 to the case where Σ is diagonal. In this case, $\Sigma = \text{diag}(\{\sigma_c^2\}_{c \in [d]})$ and

$$(*) = \mathbb{P}_{Y \sim \mathcal{N}(0, \Sigma)} \left(Y \geq \sqrt{\left(\frac{1}{N_{t,i}} + \frac{1}{N_{t,j}} \right)^{-1} \frac{1}{c(t-1, \delta)}} ((\hat{\mu}_{t,i} - \hat{\mu}_{t,j}) - (\mu_i - \mu_j)) \right) =$$

$$\prod_{c \in [d]} \mathbb{P}_{X \sim \mathcal{N}(0,1)} \left(X > \underbrace{\sqrt{\left(\frac{1}{N_{t,i}} + \frac{1}{N_{t,j}} \right)^{-1} \frac{1}{\sigma_c^2 c(t-1, \delta)}} ((\hat{\mu}_{t,i} - \hat{\mu}_{t,j}) - (\mu_i - \mu_j)) (c)}_{Z_{i,j}(c)} \right),$$

which we may rewrite with Mills ratio as $R(x) := \frac{\mathbb{P}(X > x)}{f_X(x)}$ so

$$\mathbb{P}(X > x) = R(x) \exp\left(-\frac{1}{2}x^2\right) \frac{1}{\sqrt{2\pi}} = \tilde{R}(x) \exp\left(-\frac{1}{2}x^2\right),$$

and

$$(*) = \exp\left(-\frac{1}{2} \sum_{c \in [d]} Z_{i,j}(c)\right) \prod_{c \in [d]} \tilde{R}(Z_{i,j}(c))$$

$$= \exp\left(-\left(\frac{1}{N_{t,i}} + \frac{1}{N_{t,j}}\right)^{-1} \frac{1}{c(t-1, \delta)} \sum_{c \in [d]} \frac{((\hat{\mu}_{t,i} - \hat{\mu}_{t,j}) - (\mu_i - \mu_j)) (c)^2}{2\sigma_c^2}\right) \prod_{c \in [d]} \tilde{R}(Z_{i,j}(c))$$

and by Cauchy-Schwarz inequality,

$$|((\hat{\mu}_{t,i} - \hat{\mu}_{t,j}) - (\mu_i - \mu_j)) (c)| \leq \sqrt{\frac{1}{N_{t,i}} + \frac{1}{N_{t,j}}} \sqrt{N_{t,i}(\hat{\mu}_{t,i} - \mu_i)(c)^2 + N_{t,j}(\mu_j - \hat{\mu}_{t,j})(c)^2}.$$

Therefore

$$(*) \geq \exp\left(-\frac{1}{c(t-1, \delta)} \sum_{k \in \{i,j\}} \frac{N_{t,k}}{2} \|\hat{\mu}_{t,k} - \mu_k\|_{\Sigma^{-1}}^2\right) \prod_{c \in [d]} \tilde{R}(Z_{i,j}(c)), \quad (5.10)$$

and using the following lemma,

Lemma 5.5.2. *The Mills ratio R is decreasing, log-convex and satisfies for all $(x_1, \dots, x_p) \in \mathbb{R}^p$,*

$$R\left(\frac{1}{p} \sum_{i=1}^p x_i\right) \leq \prod_{i=1}^p R(x_i).$$

we obtain

$$\begin{aligned}
 \prod_{c \in [d]} \tilde{R}(Z_{i,j}(c)) &\geq \tilde{R} \left(\frac{1}{d} \sum_{c \in [d]} Z_{i,j}(c) \right)^d \\
 &\stackrel{(a)}{\geq} \tilde{R} \left(\frac{1}{d} \sum_{c \in [d]} \sqrt{\frac{1}{\sigma_c^2 c(t-1, \delta)}} \sqrt{N_{t,i}(\hat{\mu}_{t,i} - \mu_i)(c)^2 + N_{t,j}(\mu_j - \hat{\mu}_{t,j})(c)^2} \right)^d \\
 &\stackrel{(b)}{\geq} \tilde{R} \left(\frac{\sqrt{2}}{\sqrt{c(t-1, \delta)d}} \sqrt{\sum_{k \in \{i,j\}} \frac{N_{t,k}}{2} \|\hat{\mu}_{t,k} - \mu_k\|_{\Sigma^{-1}}^2} \right)^d,
 \end{aligned}$$

where (a) follows since \tilde{R} is decreasing (Lemma 5.5.2) and (b) follows from this monotonicity and Cauchy-Schwarz. Combining the last display with (5.10), we prove that on the event $\mathcal{E}_{t,\delta/2}$, we have

$$(*) \geq \exp \left(-\frac{\beta(t-1, \delta/2)}{c(t-1, \delta)} \right) \tilde{R} \left(\sqrt{\frac{2\beta(t-1, \delta/2)}{c(t-1, \delta)d}} \right)^d$$

Choosing

$$c(t, \delta) := \frac{\beta(t, \delta/2)}{\log(1/\delta)}$$

yields,

$$(*) \geq \delta \tilde{R} \left(\sqrt{\frac{2 \log(1/\delta)}{d}} \right)^d = r(\delta, d) \quad \text{with} \quad r(\delta, n) = \tilde{R} \left(\sqrt{\frac{2 \log(1/\delta)}{n}} \right)^n$$

and proceeding identically, we prove that under the event $\mathcal{E}_{t,\delta/2}$,

$$\begin{aligned}
 &\mathbb{P}_{Y \sim \mathcal{N}(0_d, \Sigma)} \left(Y > \sqrt{\frac{N_{t,i}}{c(t-1, \delta)}} (\mu_i - \hat{\mu}_{t,i}) \right) \prod_{j \in \hat{S}_t} \mathbb{P}_{X \sim \mathcal{N}(0,1)} \left(X > \sqrt{\frac{N_{t,j}}{c(t-1, \delta)}} \|\mu_j - \hat{\mu}_{t,j}\|_{\Sigma^{-1}} \right) \\
 &\geq \delta \tilde{R} \left(\sqrt{\frac{2 \log(1/\delta)}{d + |\hat{S}_t|}} \right)^{d + |\hat{S}_t|}.
 \end{aligned}$$

All these combined, we proved that

$$\begin{aligned}
 &\mathbb{1}(\mathcal{E}_{t,\delta/2}) \mathbb{1}(\hat{S}_t \neq \mathcal{S}^*) \mathbb{P}_{\hat{\Pi}_t | \mathcal{H}_t}(\text{Alt}(\hat{S}_t)) \geq \delta \max \left\{ \mathbb{1}(\mu \in \text{Alt}^-(\hat{S}_t)) r(\delta, d), \right. \\
 &\quad \left. \mathbb{1}(\mu \in \text{Alt}^+(\hat{S}_t)) r(\delta, d + |\hat{S}_t|) \right\} \\
 &\geq \delta \min\{r(\delta, d), r(\delta, d + |\hat{S}_t|)\}
 \end{aligned}$$

which follows since $\hat{S}_t \neq \mathcal{S}^* \implies \mu \in \text{Alt}(\hat{S}_t) = \text{Alt}^+(\hat{S}_t) \cup \text{Alt}^-(\hat{S}_t)$.

The conclusion is immediate by taking $M(t, \delta) = \left\lceil \frac{\log(2t^s \zeta(s)/\delta)}{\delta q(t, \delta)} \right\rceil$ where

$$q(t, \delta) = \min \left\{ r(\delta, d), r(\delta, d + |\widehat{S}_t|) \right\}.$$

Moreover, it is known (Birnbaum 1942) that for $x \geq 0$,

$$R(x) \geq \frac{2}{x + \sqrt{x^2 + 4}}$$

and $R(x) \sim \frac{1}{x}, x \rightarrow +\infty$. The Mills ratio of the standard normal is implemented in common scientific libraries. To compute $M(t, \delta)$, we just need to compute at most the values $(r(\delta, d + k))_{k \in [K] \cup \{0\}}$ at the initialisation of the algorithm. Using the lower bound on the Mills ratio, one can easily check that

$$\frac{\log(M(t, \delta))c(t, \delta)}{\log(1/\delta)} \underset{\delta \rightarrow 0}{\leq} 1.$$

Proof of Lemma 5.5.2. The monotonicity is well-known and simple to prove by noting that

$$\begin{aligned} R(x) &:= \exp(x^2/2) \int_x^\infty \exp(-t^2/2) dt \\ &= \exp(x^2/2) \int_0^\infty \exp(-(t+x)^2/2) dt \\ &= \int_0^\infty \exp(-tx - t^2/2) dt, \end{aligned}$$

from which we have $R(x + \alpha^2) = \int_0^\infty \exp(-t\alpha^2) \exp(-tx - t^2/2) dt < R(x)$. The log-convexity is proven in Theorem 2.5 of Baricz 2008. From the log-convexity and using Jensen inequality, we have

$$\begin{aligned} \log R\left(\frac{1}{p} \sum_{i=1}^p x_i\right) &\leq \frac{1}{p} \sum_{i=1}^p \log R(x_i) \\ &= \frac{1}{p} \log \left(\prod_{i=1}^p R(x_i) \right) \end{aligned}$$

which by monotonicity of the log proves the claimed statement. □

PSI with a full covariance matrix. We specialize Lemma 5.5.1 to the case where the covariance Σ is non-diagonal. To derive this result, we will use the following lemma.

Lemma 5.5.3. *Let Σ a covariance matrix, V diagonal matrix such that $(V - \Sigma^{-1})$ is psd, and define $d_\Sigma := \|1_d\|_{(V^{-1/2}\Sigma^{-1}V^{-1/2})}^2$. Then, for all $x \in \mathbb{R}^d$, it holds that*

$$\mathbb{P}_{X \sim \mathcal{N}(0_d, \Sigma)}(X \geq x) \geq (2\pi)^{-d/2} \det(V\Sigma)^{-1/2} \exp\left(-\frac{1}{2}x^\top \Sigma^{-1}x\right) \prod_{c \in [d]} R(e_c^\top V^{-1/2} \Sigma^{-1} x),$$

and in particular, $\prod_{c \in [d]} R(e_c^\top V^{-1/2} \Sigma^{-1} x) \geq R(\|x\|_{\Sigma^{-1}} \sqrt{d_\Sigma}/d)^d$.

Applying this result with $V = \bar{\sigma}I_d$ where $\bar{\sigma} = \|\Sigma^{-1}\|$, we have

$$\mathbb{P}_{Y \sim \mathcal{N}(0_d, \Sigma)} \left(Y > \sqrt{\frac{N_{t,i}}{c(t-1, \delta)}} (\mu_i - \hat{\mu}_{t,i}) \right) \geq \frac{1}{\sqrt{(2\pi)^d \det(\bar{\sigma}\Sigma)}} \exp \left(-\frac{1}{2} \|\mu_i - \hat{\mu}_{t,i}\|_{\Sigma^{-1}}^2 \right) \cdot \prod_{c \in [d]} R(e_c^\top V^{-1/2} \Sigma^{-1} u_t),$$

with $u_t := \sqrt{\frac{N_{t,i}}{c(t-1, \delta)}} (\mu_i - \hat{\mu}_{t,i})$, therefore, letting

$$(**) := \mathbb{P}_{Y \sim \mathcal{N}(0_d, \Sigma)} \left(Y > \sqrt{\frac{N_{t,i}}{c(t-1, \delta)}} (\mu_i - \hat{\mu}_{t,i}) \right) \prod_{j \in \hat{S}_t} \mathbb{P}_{X \sim \mathcal{N}(0,1)} \left(X > \sqrt{\frac{N_{t,j}}{c(t-1, \delta)}} \|\mu_j - \hat{\mu}_{t,j}\|_{\Sigma^{-1}} \right)$$

we have

$$(**) \geq \det(\bar{\sigma}\Sigma)^{-1/2} \exp \left(-\sum_{k \in \{i\} \cup \hat{S}_t} \frac{1}{2} \|\mu_k - \hat{\mu}_{t,k}\|_{\Sigma^{-1}}^2 \right) \cdot \prod_{j \in \hat{S}_t} \tilde{R} \left(\sqrt{\frac{N_{t,j} \|\mu_j - \hat{\mu}_{t,j}\|_{\Sigma^{-1}}^2}{c(t-1, \delta)}} \right) \prod_{c \in [d]} \tilde{R}(e_c^\top V^{-1/2} \Sigma^{-1} u_t),$$

then note that by Lemma 5.5.2, we have

$$\prod_{j \in \hat{S}_t} \tilde{R} \left(\sqrt{\frac{N_{t,j}}{c(t-1, \delta)}} \|\mu_j - \hat{\mu}_{t,j}\|_{\Sigma^{-1}} \right) \prod_{c \in [d]} \tilde{R}(e_c^\top V^{-1/2} \Sigma^{-1} u_t / d) \geq \tilde{R} \left(\frac{1_d^\top V^{-1/2} \Sigma^{-1} u_t + 1_{|\hat{S}_t|}^\top h_t}{d + |\hat{S}_t|} \right)^{d+|\hat{S}_t|}$$

with $h_t := \left(\sqrt{\frac{N_{t,j}}{c(t-1, \delta)}} \|\mu_j - \hat{\mu}_{t,j}\|_{\Sigma^{-1}} \right)_{j \in |\hat{S}_t|}$ and by Cauchy-Swcharz inequality, replacing $V = \bar{\sigma}I_d$,

$$\begin{aligned} 1_d^\top V^{-1/2} \Sigma^{-1} u_t + 1_{|\hat{S}_t|}^\top h_t &\leq \sqrt{\|1_d\|_{(\bar{\sigma}\Sigma)^{-1}}^2 + \|1_{|\hat{S}_t|}\|^2} \sqrt{\|u_t\|_{\Sigma^{-1}}^2 + \|h_t\|^2} \\ &= \sqrt{d_\Sigma + |\hat{S}_t|} \left(\sqrt{\frac{1}{c(t-1, \delta)} \sum_{k \in \{i\} \cup \hat{S}_t} N_{k,t} \|\mu_k - \hat{\mu}_{t,k}\|_{\Sigma^{-1}}^2} \right) \\ &\leq \sqrt{d_\Sigma + |\hat{S}_t|} \sqrt{\frac{2\beta(t-1, \delta/2)}{c(t-1, \delta)}} \end{aligned}$$

where the last inequality follows on the event $\mathcal{E}_{t, \delta/2}$. Combining these displays with $c(t, \delta) := \frac{\beta(t, \delta/2)}{\log(1/\delta)}$, we have on the event $\mathcal{E}_{t, \delta/2}$,

$$\begin{aligned} (**) &\geq \det(\bar{\sigma}\Sigma)^{-1/2} \delta \tilde{R} \left(\frac{\sqrt{d_\Sigma + |\hat{S}_t|}}{d + |\hat{S}_t|} \sqrt{2 \log(1/\delta)} \right) \\ &= \det(\bar{\sigma}\Sigma)^{-1/2} \delta r \left(\delta^{\frac{d_\Sigma + |\hat{S}_t|}{d + |\hat{S}_t|}}, d + |\hat{S}_t| \right). \end{aligned}$$

Following the same reasoning, it is simple to show that under the event $\mathcal{E}_{t,\delta/2}$, we have

$$\begin{aligned} & \mathbb{P}_{Y \sim \mathcal{N}(0_d, \Sigma)} \left(Y \geq \sqrt{\left(\frac{1}{N_{t,i}} + \frac{1}{N_{t,j}} \right)^{-1} \frac{1}{c(t-1, \delta)}} \left((\hat{\mu}_{t,i} - \hat{\mu}_{t,j}) - (\mu_i - \mu_j) \right) \right) \\ & \geq \det(\bar{\sigma}\Sigma)^{-1/2} \delta \tilde{R}(\sqrt{2 \log(1/\delta)} d_\Sigma / d) = \det(\bar{\sigma}\Sigma)^{-1/2} \delta r(\delta^{\frac{d_\Sigma}{d}}, d). \end{aligned}$$

Combining these, we have

$$\begin{aligned} \mathbb{1}(\mathcal{E}_{t,\delta/2}) \mathbb{1}(\hat{S}_t \neq \mathcal{S}^*) \mathbb{P}_{\hat{\Pi}_t | \mathcal{H}_t}(\text{Alt}(\hat{S}_t)) & \geq \delta \det(\bar{\sigma}\Sigma)^{-1/2} \max \left\{ \mathbb{1} \left(\mu \in \text{Alt}^-(\hat{S}_t) \right) r(\delta^{\frac{d_\Sigma}{d}}, d), \right. \\ & \quad \left. \mathbb{1} \left(\mu \in \text{Alt}^+(\hat{S}_t) \right) r(\delta^{\frac{d_\Sigma + |\hat{S}_t|}{d + |\hat{S}_t|}}, d + |\hat{S}_t|) \right\} \\ & \geq \delta \det(\bar{\sigma}\Sigma)^{-1/2} \min \left\{ r(\delta^{\frac{d_\Sigma}{d}}, d), r(\delta^{\frac{d_\Sigma + |\hat{S}_t|}{d + |\hat{S}_t|}}, d + |\hat{S}_t|) \right\} \end{aligned}$$

which follows similarly to the Σ diagonal case, since $\hat{S}_t \neq \mathcal{S}^* \implies \mu \in \text{Alt}(\hat{S}_t) = \text{Alt}^+(\hat{S}_t) \cup \text{Alt}^-(\hat{S}_t)$. Note that when $\Sigma = \sigma I_d$, we recover the results of section 5.5.1.

Proof of Lemma 5.5.3. Σ is a $d \times d$ covariance matrix and f_Σ is the density function of $\mathcal{N}(0_d, \Sigma)$. We let \mathcal{R}_Σ denote the Mills ratio of the distribution $\mathcal{N}(0_d, \Sigma)$, which is defined for a vector $x \in \mathbb{R}^d$ as

$$\mathcal{R}_\Sigma(x) := \frac{\mathbb{P}(\mathcal{N}(0_d, \Sigma) \geq x)}{f_\Sigma(x)} \quad (5.11)$$

where for two vectors x, y , the notation $x \leq y$ should be understood component-wise. Expanding (5.11) gives

$$\mathcal{R}_\Sigma(x) = \exp(\|x\|_{\Sigma^{-1}}^2 / 2) \int_{(x, \infty)} \exp(-\|u\|_{\Sigma^{-1}}^2 / 2) du,$$

with $(x, \infty) = (x(1), \infty) \times \dots \times (x(d), \infty)$. By a simple translation, it follows that

$$\mathcal{R}_\Sigma(x) = \int_{(0_d, \infty)} \exp(-x^\top \Sigma^{-1} u - \|u\|_{\Sigma^{-1}}^2 / 2) du.$$

Since $(V - \Sigma^{-1})$ is psd, we have

$$\begin{aligned} \mathcal{R}_\Sigma(x) & \geq \int_{(0_d, \infty)} \exp(-x^\top \Sigma^{-1} u - u^\top V u / 2) du \\ & = \det(V)^{-1/2} \int_{(0_d, \infty)} \exp(-x^\top \Sigma^{-1} V^{-1/2} u - u^\top u / 2) du \end{aligned}$$

which follows by change of variable. Thus, letting $z = V^{-1/2} \Sigma^{-1} x$, we have

$$\begin{aligned} \mathcal{R}_\Sigma(x) & \geq \det(V)^{-1/2} \int_{(0_d, \infty)} \exp(-u^\top z - u^\top u / 2) du \\ & = \det(V)^{-1/2} \prod_{c \in [d]} R(z(c)) \end{aligned}$$

then using Lemma 5.5.2, it follows that

$$\begin{aligned}\mathcal{R}_\Sigma(x) &\geq \det(V)^{-1/2} R \left(\mathbb{1}_d^\top V^{-1/2} \Sigma^{-1} x / d \right)^d \\ &\geq \det(V)^{-1/2} R \left(\| \mathbb{1}_d \|_{V^{-1/2} \Sigma^{-1} V^{-1/2}} \| x \|_{\Sigma^{-1}} / d \right)^d.\end{aligned}$$

□

5.5.2 Sample complexity

We recall the following result on the sampling rule, proven in [Kone, Jourdan, et al. 2025](#).

Theorem 5.5.4 (Theorem 3 of [Kone, Jourdan, et al. 2025](#)). *There exists events $(\Xi_t)_{t \geq 1}$ and $T_0 \in \mathbb{N}$ such that for all $t \geq T_0$,*

$$2GLR(t) \geq (t-1)T^*(\nu)^{-1} - f(t), \quad (5.12)$$

with $f(t) =_\infty o(t)$ and $\mathbb{P}_\nu(\Xi_t) \geq 1 - \mathcal{O}(1/t^2)$. In particular for $t \geq T_0$, $\widehat{S}_t = \mathcal{S}^*$.

The next lemma is proven in section 5.3.1

Lemma 5.5.5 (Posterior probability bound via GLR). *At each round t , it holds that*

$$\mathbb{P}_{\widehat{\Pi}_t | \mathcal{H}_t}(\text{Alt}(\widehat{S}_t)) \leq \alpha_t \exp\left(-\frac{GLR(t)}{c(t-1, \delta)}\right),$$

where $\alpha_t := \frac{1}{2}n(p_t)$ with $p_t := |\widehat{S}_t|$ and $n(p_t) = p_t(p_t - 1) + (K - p_t)d^{p_t}$.

Theorem 5.3.1. *For any budget M and inflation c such that $\limsup_{\delta \rightarrow 0} \frac{c(t, \delta) \log M(t, \delta)}{\log(1/\delta)} \leq 1$, PSIPS satisfies that*

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_\nu[\tau_\delta^{\text{PS}}]}{\log \frac{1}{\delta}} \leq T^*(\nu) \quad \text{and} \quad \mathbb{P}\left(\limsup_{\delta \rightarrow 0} \frac{\tau_\delta^{\text{PS}}}{\log \frac{1}{\delta}} \leq T^*(\nu)\right) = 1, \quad \forall \nu \in \mathcal{D}^K.$$

Proof. In this section, $(\Xi_t)_{t \geq 1}$ denotes the sequence of events Theorem 5.5.4 and $T_0 \in \mathbb{N}$.

We have $\tau_\delta^{\text{PS}} = 1 + \sum_{t \geq 1} \mathbb{1}(\tau_\delta^{\text{PS}} > t)$ then

$$\begin{aligned}\mathbb{E}_\nu[\tau_\delta^{\text{PS}}] &= 1 + \mathbb{E}\left[\sum_{t \geq 1} \mathbb{P}_\nu(\tau_\delta^{\text{PS}} > t | \mathcal{H}_t)\right] \\ &= 1 + \mathbb{E}_\nu\left[\sum_{t \geq 1} \mathbb{P}_\nu\left(\exists m \in [M(t, \delta)] : \theta_t^m \in \text{Alt}(\widehat{S}_t) | \mathcal{H}_t\right)\right] \\ &\leq \mathbb{E}_\nu\left[\sum_{t \geq T_0} \mathbb{1}(\Xi_t) M(t, \delta) \mathbb{P}_\nu\left(\theta_t \in \text{Alt}(\widehat{S}_t) | \mathcal{H}_t\right)\right] + \left[\sum_{t \geq 1} \mathbb{P}_\nu(\Xi_t^c)\right] + (T_0 + 1),\end{aligned}$$

where $\theta_t, \theta_t^1, \dots, \dots, \theta_t^m$ are *i.i.d.* given \mathcal{H}_t . Moreover, by Lemma 5.5.5, we have

$$\mathbb{P}_\nu \left(\theta_t \in \text{Alt}(\widehat{S}_t) | \mathcal{H}_t \right) \leq \alpha_t \exp \left(-\frac{\text{GLR}(t)}{c(t-1, \delta)} \right)$$

therefore,

$$\mathbb{E}_\nu[\tau_\delta^{\text{PS}}] \leq \underbrace{\mathbb{E}_\nu \left[\sum_{t \geq t_3} \mathbb{1}(\Xi_t) \alpha_0 M(t, \delta) \exp \left(-\frac{\text{GLR}(t)}{c(t-1, \delta)} \right) \right]}_{L_1(\delta)} + \underbrace{\left[\sum_{t \geq 1} \mathbb{P}_\nu(\Xi_t^c) \right]}_{L_2} + T_0,$$

then, since $\mathbb{P}_\nu(\Xi_t^c) \leq \mathcal{O}(1/t^2)$ we immediately have

$$L_2 \leq 5\pi^2/6. \quad (5.13)$$

It remains to bound $L_1(\delta)$. Using the saddle-point convergence property on the event Ξ_t ensures that (Theorem 5.5.4) for $t \geq T_0$,

$$\text{GLR}(t) \geq tT^*(\nu)^{-1} - f(t), \quad (5.14)$$

which further results in

$$\begin{aligned} L_1(\delta) &\leq \sum_{t \geq t_3} \alpha_0 M(t, \delta) \exp \left(-\frac{t}{T^*(\nu)c(t, \delta)} + f(t)/c(t, \delta) \right) \\ &= \sum_{t \geq t_3} \exp \left(-\frac{t}{T^*(\nu)c(t, \delta)} + f(t)/c(t, \delta) + \log(\alpha_0 M(t, \delta)) \right). \end{aligned}$$

To bound the above quantity, let us introduce

$$T(\delta) := \sup \left\{ t \mid \frac{t}{T^*(\nu)c(t, \delta)} - \log(M(t, \delta)) - f(t)/c(t, \delta) - \log(\alpha_0) \leq \log(t \log(t)) \right\}, \quad (5.15)$$

then it follows that

$$L_1(\delta) \leq T(\delta) + \sum_{t \geq t_3} (t \log(t))^{-1}. \quad (5.16)$$

Further observe that $T(\delta)$ can be rewritten as

$$T(\delta) := \sup \{ t \mid t \leq T^*(\nu)c(t, \delta) (\log(M(t, \delta)) + f(t)/c(t, \delta) + \log(\alpha_0) + \log(t \log(t))) \}.$$

Bounding $T(\delta)$. Since f is sub-linear in t , there exists $\varepsilon \in (0, 1)$ such that $f(t) = o_{t \rightarrow \infty}(t^\varepsilon)$ then $f(\log(1/\delta)^{1/\varepsilon}) = o(\log(1/\delta))$. Further observe that for $t_\delta = \log(1/\delta)^{1/\alpha}$,

$$\underbrace{c(t_\delta, \delta) \log(M(t_\delta, \delta)) + f(\log(1/\delta)^{1/\varepsilon}) + c(t_\delta, \delta) (\log(t_\delta \log(t_\delta)) + \alpha_0)}_{b(t_\delta)} \underset{\delta \rightarrow 0}{\sim} \log(1/\delta).$$

Let $\delta_{\min} \in (0, 1)$ be defined as

$$\delta_{\min} := \inf \{ \delta \in (0, 1) \mid b(t_\delta) > \log(1/\delta)^{1/\varepsilon} T^*(\nu)^{-1} \}$$

which is well defined as $b(t_\delta) \underset{\delta \rightarrow 0}{\sim} \log(1/\delta)$ and $\varepsilon \in (0, 1)$. Letting $T_{\max} = \log(1/\delta_{\min})^{1/\varepsilon}$, further note that for all $t \geq T_{\max}$, there is $(0, 1) \ni \delta' \leq \delta_{\min}$ $t'_\delta = t$ and $b(t_{\delta'}) < t_{\delta'}$. Therefore, for all $\delta \leq \delta_{\min}$

$$T(\delta) \leq \log(1/\delta)^{1/\varepsilon} \quad (5.17)$$

and further noting that by definition $T(\delta) \leq T^*(\nu)b(T(\delta))$ and b is increasing, it follows that

$$T(\delta) \leq T^*(\nu)b(\log(1/\delta)^{1/\varepsilon}). \quad (5.18)$$

Combining (5.13), (5.16) and (5.18), it follows that for $\delta \leq \delta_{\min}$,

$$\mathbb{E}_\nu[\tau_\delta^{\text{PS}}] \leq T^*(\nu)b(\log(1/\delta)^{1/\varepsilon}) + 5\pi^2/6 + (T_0 + 1) + \sum_{t \geq T_0} (t \log(t))^{-1}. \quad (5.19)$$

Finally, noting that $b(\log(1/\delta)^{1/\varepsilon}) \underset{\delta \rightarrow 0}{\sim} \log(1/\delta)$, we have proved that

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}[\tau_\delta^{\text{PS}}]}{\log \frac{1}{\delta}} \leq T^*(\nu). \quad (5.20)$$

Almost-sure upper bound. The almost-sure bound on the sample complexity is derived from an application of Borel-Cantelli's lemma. Indeed, introducing for any T the event

$$\mathcal{E}(T) := \left\{ \sum_{t=1}^T \mathbb{1}(\tau_\delta^{\text{PS}} > t) - \sum_{t=1}^T \mathbb{P}(\tau_\delta > t | \mathcal{H}_{t-1}) \leq \sqrt{2T \log T^2} \right\},$$

we have when $\mathcal{E}(T)$ holds

$$\begin{aligned} \min \{ \tau_\delta^{\text{PS}}, T \} &\leq 1 + \sum_{t=1}^T \mathbb{1}(\tau_\delta^{\text{PS}} > t) \\ &\leq 1 + \sum_{t=1}^T \mathbb{P}(\tau_\delta^{\text{PS}} > t | \mathcal{H}_{t-1}) + \sqrt{2T \log T^2} \\ &\leq T(\delta) + \sqrt{2T \log T^2} + T_0, \end{aligned}$$

where the last inequality follows from the derivations above on the event Ξ_T and holds with probability larger than $1 - \mathcal{O}(1/T^2)$. Next, assume $\mathcal{E}(T)$ holds. Then, if

$$T \geq T'(\delta) := \sup \left\{ t \geq 1 : t < T_0 + T(\delta) + \sqrt{2t \log t^2} \right\},$$

and $\mathcal{E}(T) \cap \Xi(T)$ holds then the stopping rule would trigger before episode T . Thanks to Azuma-Hoeffding $\mathcal{E}(T)$ holds with probability at least $1 - \mathcal{O}(1/T^2)$. Thus $\Xi(T)^c$ holds with probability at most $\mathcal{O}(1/T^2)$, similarly for $\mathcal{E}(T)^c$ and, as $\sum_T \frac{1}{T^2} < \infty$, we invoke Borel-Cantelli's lemma to justify that with probability 1, there exists \tilde{T} (possibly random) such that for $T \geq \tilde{T}$, $\Xi(T) \cap \mathcal{E}(T)$ holds.

From the calculations above, we observe that the leading term in $T'(\delta)$ is at most of order $\log(1/\delta)$. Let δ'_{\min} be such that $\forall \delta \leq \delta'_{\min}$,

$$\lceil \log(1/\delta)^{3/2} \rceil > T'(\delta),$$

which is well-defined. Let $\delta \leq \delta'_{\min}$ such that additionally $\lceil \log(1/\delta)^{3/2} \rceil > \tilde{T}$. We then have

$$\begin{aligned} \tau_{\delta}^{\text{PS}} &\leq T_0 + T(\delta) + \sqrt{2 \lceil \log(1/\delta)^{3/2} \rceil \log \lceil \log(1/\delta)^{3/2} \rceil^2} \\ &\leq T_0 + T(\delta) + \log(1/\delta)^{3/4} \sqrt{8 \log(1 + \log(1/\delta)^{3/2})}. \end{aligned}$$

Finally recalling that

$$\limsup_{\delta \rightarrow 0} \frac{T(\delta)}{\log \frac{1}{\delta}} \leq T^*(\nu).$$

We conclude by noting that

$$\begin{aligned} \limsup_{\delta \rightarrow 0} \frac{\tau_{\delta}^{\text{PS}}}{\log \frac{1}{\delta}} &\leq \limsup_{\delta \rightarrow 0} \frac{T(\delta)}{\log \frac{1}{\delta}} + \lim_{\delta \rightarrow 0} \frac{T_0 + \log(1/\delta)^{3/4} \sqrt{8 \log(1 + \log(1/\delta)^{3/2})}}{\log \frac{1}{\delta}} \\ &\leq T^*(\nu). \end{aligned}$$

□

5.5.3 Posterior convergence

We prove that the posterior contraction rate of PSIPS is unimprovable.

In BAI, Russo 2016 and Z. Li et al. 2024 proved a similar result for truncated Gaussian (restricted to a bounded domain). We prove it more generally in the unbounded setting for Gaussian distributions, thanks to the novel Lemma 5.5.7.

Theorem 5.3.2. *Let $\tilde{\Pi}_t := \bigotimes_{k=1}^K \mathcal{N}(\hat{\mu}_{t,k}, \Sigma/N_{t,k})$ be the posterior distribution under a flat Gaussian prior (without inflation). For any bandit instance $\nu \in \mathcal{D}^K$, it holds with probability 1 that*

$$\limsup_{t \rightarrow +\infty} -\frac{1}{t} \log \mathbb{P}_{\tilde{\Pi}_t | \mathcal{H}_t}(\text{Alt}(\mathcal{S}^*)) \leq T^*(\nu)^{-1},$$

for any algorithm, and PSIPS almost surely satisfies that

$$\liminf_{t \rightarrow +\infty} -\frac{1}{t} \log \mathbb{P}_{\tilde{\Pi}_t | \mathcal{H}_t}(\text{Alt}(\mathcal{S}^*)) \geq T^*(\nu)^{-1}.$$

Proof. We first prove the upper bound. Similarly to Lemma 5.5.5, we can derive (using Lemma 5.3.4 and Lemma 5.3.5),

$$\mathbb{P}_{\tilde{\Pi}_t | \mathcal{H}_t}(\text{Alt}(\mathcal{S}^*)) \leq \alpha_0 \exp \left(- \inf_{\lambda \in \text{Alt}(\mathcal{S}^*)} \sum_{k=1}^K \frac{N_{t,k}}{2} \|\hat{\mu}_{t,k} - \lambda_k\|_{\Sigma^{-1}}^2 \right).$$

From Theorem 5.5.4, there exists events $(\Xi_t)_{t \geq 1}$ and $T_0 \in \mathbb{N}$ such that for $t \geq T_0$, $\widehat{S}_t = S^*$, and

$$\text{GLR}(t) \geq (t-1)T^*(\nu)^{-1} - f(t), \quad (5.21)$$

with $f(t) = o(t)$ and $\mathbb{P}_\nu(\Xi_t) \geq 1 - \mathcal{O}(1/t^2)$. Since $\sum_{t \geq 1} \mathbb{P}(\Xi_t^c) < \infty$, by Borel-Cantelli's lemma, with probability 1, there exists a finite time \widetilde{T}_0 possibly random such that for $t \geq \widetilde{T}_0$, Ξ_t holds. So for $t \geq \max(T_0, \widetilde{T}_0)$, we have $\widehat{S}_t = S^*$ and

$$\begin{aligned} \inf_{\lambda \in \text{Alt}(S^*)} \frac{N_{t,k}}{2} \|\lambda_k - \widehat{\mu}_{t,k}\|_{\Sigma^{-1}}^2 = \text{GLR}(t) &\geq (t-1)T^*(\nu)^{-1} - f(t), \\ &\geq tT^*(\nu)^{-1} - o(t) \end{aligned}$$

then

$$\mathbb{P}_{\widetilde{\Pi}_t | \mathcal{H}_t}(\text{Alt}(S^*)) \leq \alpha_0 \exp(-tT^*(\nu)^{-1} + o(t)),$$

so that

$$-\frac{1}{t} \log(\mathbb{P}_{\widetilde{\Pi}_t | \mathcal{H}_t}(\text{Alt}(S^*))) \geq T^*(\nu)^{-1} - (1/t) \log(\alpha_0) - o(t)/t.$$

Put together, the above displays show that with probability 1,

$$\liminf_{t \rightarrow \infty} -\frac{1}{t} \log(\mathbb{P}_{\widetilde{\Pi}_t | \mathcal{H}_t}(\text{Alt}(S^*))) \geq T^*(\nu)^{-1}.$$

The proof of the lower bound uses Lemma 5.5.6.

Let \mathbb{B}_ε be the ball centered on μ and with radius ε as in Lemma 5.5.7. We have $-\mathbb{P}(\text{Alt}(S^*)) \leq -\mathbb{P}(\text{Alt}(S^*) \cap \mathbb{B}_\varepsilon)$, then $O = \text{Alt}(S^*) \cap \mathbb{B}_\varepsilon$ is countably convex (Lemma 5.3.4) and \mathbb{B}_ε is convex) and bounded. Applying Lemma 5.5.6, we obtain

$$\begin{aligned} \limsup_{t \rightarrow \infty} -\frac{1}{t} \log \mathbb{P}_{\widetilde{\Pi}_t | \mathcal{H}_t}(\text{Alt}(S^*)) &\leq \limsup_{t \rightarrow \infty} -\frac{1}{t} \log \mathbb{P}_{\widetilde{\Pi}_t | \mathcal{H}_t}(\text{Alt}(S^*) \cap \mathbb{B}_\varepsilon) \\ &\leq \frac{1}{2} \sup_{w \in \Delta_K} \inf_{\lambda \in \text{Alt}(S^*) \cap \mathbb{B}_\varepsilon} \left[\sum_{k=1}^K \frac{w_k}{2} \|\mu_k - \lambda_k\|_{\Sigma^{-1}}^2 \right] \\ &= \frac{1}{2} \sup_{w \in \Delta} \inf_{\lambda \in \text{Alt}(S^*)} \left[\sum_{k=1}^K \frac{w_k}{2} \|\mu_k - \lambda_k\|_{\Sigma^{-1}}^2 \right] \quad (\text{Lemma 5.5.7}) \\ &= T^*(\nu)^{-1}, \end{aligned}$$

which concludes the proof. \square

Lemma 5.5.6. *Let $O \subset \mathcal{I}^K$ be a countably convex bounded set with a non-empty interior. With probability 1, it holds that*

$$\limsup_{t \rightarrow \infty} -\frac{1}{t} \log \mathbb{P}_{\widetilde{\Pi}_t | \mathcal{H}_t}(O) \leq \frac{1}{2} \sup_{w \in \Delta_K} \inf_{\lambda \in O} \left[\sum_{k=1}^K \|\mu_k - \lambda_k\|_{\Sigma^{-1}}^2 \right].$$

Proof. Since O is countably there exists convex sets $C_1 \dots C_n$ such that $O = \cup_{i \in [n]} C_i$.

We have

$$\mathbb{P}_{\tilde{\Pi}_t | \mathcal{H}_t}(O) = (2\pi)^{-dh/2} D_t^{-1/2} \int_O \exp \left(- \sum_{k=1}^K \frac{N_{t,k}}{2} \|\hat{\mu}_{t,k} - \lambda_k\|_{\Sigma^{-1}}^2 \right) d\lambda,$$

where $D_t = \prod_{k=1}^K \det(\Sigma/N_{t,k})$.

Let $\gamma > 0$ (to be defined) and

$$\tilde{\lambda}_t \in \operatorname{argmin}_{\lambda \in O} \left[\sum_{k=1}^K \frac{N_{t,k}}{2} \|\lambda_k - \hat{\mu}_{t,k}\|_{\Sigma^{-1}}^2 \right].$$

Since O is a union of convex sets, there exists a convex set $\mathcal{C} \subset O$ such that $\tilde{\lambda}_t \in \mathcal{C}$. Then letting $\mathcal{N}_\gamma := \{(1-\gamma)\tilde{\lambda}_t + \gamma\lambda, \lambda \in \mathcal{C}\} = (1-\gamma)\tilde{\lambda}_t + \gamma\mathcal{C}$ we have $\mathcal{N}_\gamma \subset \mathcal{C} \subset O$ and

$$\begin{aligned} I_t &\geq \int_{\mathcal{N}_\gamma} \exp \left(- \sum_{k=1}^K \frac{N_{t,k}}{2} \|\hat{\mu}_{t,k} - \lambda_k\|_{\Sigma^{-1}}^2 \right) d\lambda \\ &= \int_{\gamma\mathcal{C}} \exp \left(- \sum_{k=1}^K \frac{N_{t,k}}{2} \|\hat{\mu}_{t,k} - (1-\gamma)\tilde{\lambda}_{t,k} - \lambda_k\|_{\Sigma^{-1}}^2 \right) d\lambda \\ &= \int_{\mathcal{C}} \gamma^{Kd} \exp \left(- \sum_{k=1}^K \frac{N_{t,k}}{2} \|(1-\gamma)(\hat{\mu}_{t,k} - \tilde{\lambda}_{t,k}) + \gamma(\hat{\mu}_{t,k} - \lambda_k)\|_{\Sigma^{-1}}^2 \right) d\lambda \\ &\geq \int_{\mathcal{C}} \gamma^{Kd} \exp \left(- \sum_{k=1}^K \frac{N_{t,k}}{2} \left[(1-\gamma)\|\hat{\mu}_{t,k} - \tilde{\lambda}_{t,k}\|_{\Sigma^{-1}}^2 + \gamma\|\hat{\mu}_{t,k} - \lambda_k\|_{\Sigma^{-1}}^2 \right] \right) d\lambda \end{aligned}$$

which follows thanks to the convexity of the squared norm. Therefore

$$\begin{aligned} \log I_t &\geq Kd \log(\gamma) - (1-\gamma) \sum_{k=1}^K \frac{N_{t,k}}{2} \|\hat{\mu}_{t,k} - \tilde{\lambda}_{t,k}\|_{\Sigma^{-1}}^2 - \\ &\quad \frac{\gamma \mathbb{E}_{\lambda \sim \text{unif}(\mathcal{C})} \left[\sum_{k=1}^K \frac{N_{t,k}}{2} \|\lambda_k - \hat{\mu}_{t,k}\|_{\Sigma^{-1}}^2 \right]}{2} + \log(\text{vol}(\mathcal{C}_*)) \end{aligned} \tag{5.22}$$

where \mathcal{C}_* is the set of minimum positive volume among $\mathcal{C}_1, \dots, \mathcal{C}_n$. We introduce the event

$$\Xi_{1,t} := \left\{ \forall s \leq t, \sum_{k=1}^K \frac{N_{s,k}}{2} \|\mu_k - \hat{\mu}_{s,k}\|_{\Sigma^{-1}}^2 \leq \beta(t, 1/t^2) =: f_1(t) \right\},$$

where $\beta(t, \delta)$ is defined as in Lemma 5.2.1, and $\mathbb{P}_\nu(\Xi_{1,t}) \geq 1 - 1/t^2$.

Moreover, $f_1(t)$ is logarithmic in t . Since $\sum_{t \geq 1} \mathbb{P}(\Xi_t^c) < \infty$, by Borel-Cantelli's lemma, with probability 1, there exists a finite time \tilde{T} possibly random such that for $t \geq \tilde{T}$, $\Xi_{1,t}$ holds.

Taking $\gamma = 1/t$ in Equation (5.22), and for $t \geq \tilde{T}$, we get (after simplification)

$$\log I_t \geq -dh \log(t) - \sum_{k=1}^K \frac{N_{t,k}}{2} \|\hat{\mu}_{t,k} - \tilde{\lambda}_{t,k}\|_{\Sigma^{-1}}^2 - \mathcal{O}(\sqrt{tL(O)} + \sqrt{f_1(t)})^2/t + \log(\text{vol}(\mathcal{C}_*)),$$

where $L(O) = \max_{\lambda \in O} \sum_{k=1}^K \|\lambda_k - \mu_k\|_{\Sigma^{-1}}^2$. Now, by convexity, we have for any $\lambda \in O$,

$$\begin{aligned} \sum_{k=1}^K N_{t,k} \|\hat{\mu}_{t,k} - \lambda_k\|_{\Sigma^{-1}}^2 &\leq \sum_{k=1}^K N_{t,k} \|\mu_k - \lambda_k\|_{\Sigma^{-1}}^2 + 2 \sum_{k=1}^K N_{t,k} \|\hat{\mu}_{t,k} - \lambda_k\|_{\Sigma^{-1}} \|\hat{\mu}_{t,k} - \mu_k\|_{\Sigma^{-1}} \\ &\leq \sum_{k=1}^K N_{t,k} \|\mu_k - \lambda_k\|_{\Sigma^{-1}}^2 + 4f_1(t) + 4\sqrt{f_1(t)}\sqrt{tL(O)}, \end{aligned}$$

thus we have

$$\sum_{k=1}^K N_{t,k} \|\hat{\mu}_{t,k} - \tilde{\lambda}_{t,k}\|_{\Sigma^{-1}}^2 \leq \inf_{\lambda \in O} \sum_{k=1}^K N_{t,k} \|\mu_k - \lambda_k\|_{\Sigma^{-1}}^2 + 4f_1(t) + 4\sqrt{f_1(t)}\sqrt{tL(O)}.$$

In all cases, we have proved that

$$-\log \mathbb{P}_{\tilde{\Pi}_t | \mathcal{H}_t}(O) \leq \frac{1}{2} \inf_{\lambda \in O} \sum_{k=1}^K N_{t,k} \|\mu_{t,k} - \lambda_k\|_{\Sigma^{-1}}^2 + o(t) + D_t.$$

Next, due to the initialization, we have

$$D_t = \prod_{k=1}^K \det(\Sigma/N_{t,k}) \leq \det(\Sigma)^K.$$

Combining the displays above yields,

$$-\log \mathbb{P}_{\tilde{\Pi}_t | \mathcal{H}_t}(O) \leq \inf_{\lambda \in O} \sum_{k=1}^K \frac{N_{t,k}}{2} \|\mu_k - \lambda_k\|_{\Sigma^{-1}}^2 + o(t),$$

which yields the claimed result by dividing by t and taking the limit. \square

The result below shows that the best response always exists as an inf over a compact subset of the alternative. Although its closed form is unknown in PSI, we show that it belongs to a ball centered at θ and whose radius also depends on θ .

Lemma 5.5.7. *Let $w \in \mathbb{R}_+^K$. For any $\theta := (\theta_1, \dots, \theta_K) \in \mathcal{I}^K$, the following statement holds*

$$\inf_{\lambda \in \text{Alt}(S^*(\theta))} \left[\sum_{k=1}^K \frac{w_k}{2} \|\theta_k - \lambda_k\|_{\Sigma^{-1}}^2 \right] = \inf_{\lambda \in \text{Alt}(S^*(\theta)) \cap \{\lambda \mid \max_i \|\theta_i - \lambda_i\|_{\Sigma^{-1}} < \varepsilon\}} \left[\sum_{k=1}^K \frac{w_k}{2} \|\theta_k - \lambda_k\|_{\Sigma^{-1}}^2 \right], \quad (5.23)$$

where $\varepsilon(\theta) := \max \{2 \max_{i \notin S} \max_{j \in S} \|\theta_i - \theta_j\|_{\Sigma^{-1}}, \max_{i,j \in S^2} \|\theta_i - \theta_j\|_{\Sigma^{-1}}\}$, and $S = S^*(\theta)$.

Proof. The idea of the proof is to show that when there is a ball \mathbb{B} such that if an alternative parameter of θ does not belong to \mathbb{B} , then there is a parameter in $\text{Alt}(\mathcal{S}^*(\theta)) \cap \mathbb{B}$ for which the transportation cost will be smaller. Let $\theta := (\theta_1, \dots, \theta_K)^\top$, where θ_i denotes the vector mean of arm i and we let $\lambda := (\lambda_1, \dots, \lambda_K)^\top$ where $\lambda_i \in \mathbb{R}^d$. To ease notation, we let $S = \mathcal{S}^*(\theta)$. Introducing

$$\text{Alt}^-(S) := \bigcup_{i,j \in S^2: i \neq j} W_{i,j} \quad \text{and} \quad \text{Alt}^+(S) := \bigcup_{i \in S^c} V_i, \quad (5.24)$$

where $\mathcal{I}^K := \mathbb{R}^{K \times d}$, we define

$$\begin{aligned} W_{i,j} &:= \{ \lambda := (\lambda_1 \dots \lambda_K)^\top \in \mathcal{I}^K \mid \lambda_i \leq \lambda_j \} \quad \text{and} \\ V_i &:= \{ \lambda := (\lambda_1 \dots \lambda_K)^\top \in \mathcal{I}^K \mid \exists i \in (S)^c : \forall j \in S, \lambda_i \not\leq \lambda_j \}. \end{aligned}$$

Let us introduce

$$D(w, \lambda; \theta) := \sum_{i=1}^K w_i \|\lambda_i - \theta_i\|_{\Sigma^{-1}}^2.$$

Let $(i, j) \in S^2$ be fixed and $\lambda \in W_{i,j}$. Let $\alpha_{i,j} = \|\theta_i - \theta_j\|_{\Sigma^{-1}}$. If $\|\lambda_i - \theta_i\|_{\Sigma^{-1}} > \alpha_{i,j}$ then define the instance $\tilde{\lambda}$ as

$$\tilde{\lambda}_k := \begin{cases} \theta_k & \text{if } k \notin \{i, j\}, \\ \theta_j & \text{else,} \end{cases}$$

which satisfies $\tilde{\lambda}_i \leq \tilde{\lambda}_j$, so $\tilde{\lambda} \in W_{i,j}$, and

$$\begin{aligned} D(w, \tilde{\lambda}; \theta) &= w_i \|\theta_i - \theta_j\|_{\Sigma^{-1}}^2 + \sum_{k \notin \{i,j\}} w_k \|\lambda_k - \theta_k\|_{\Sigma^{-1}}^2, \\ &= w_i \|\theta_i - \theta_j\|_{\Sigma^{-1}}^2 < w_i \|\lambda_i - \theta_i\|_{\Sigma^{-1}}^2 \\ &< D(w, \lambda, \theta), \end{aligned}$$

and further observe that $\max_{k \in [K]} \|\tilde{\lambda}_k - \theta_k\|_{\Sigma^{-1}} \leq \alpha_{i,j}$. We proceed similarly if $\|\lambda_j - \theta_j\|_{\Sigma^{-1}} > \alpha_{i,j}$, by defining

$$\tilde{\lambda}_k := \begin{cases} \theta_k & \text{if } k \notin \{i, j\}, \\ \theta_i & \text{else,} \end{cases}$$

and the same conclusion follows. So we have proved that there exists $\tilde{\lambda} \in W_{i,j}$ with $\max_{k \in [K]} \|\tilde{\lambda}_k - \theta_k\|_{\Sigma^{-1}} \leq \alpha_{i,j}$ and for which the transportation cost is not larger than that of λ . We prove a similar property for $\text{Alt}^+(S)$.

Now, fix $i \notin S$ and take $\lambda \in V_i$. Let $b_i = \max_{k \in S} \|\theta_i - \theta_k\|_{\Sigma^{-1}}$. If $\|\lambda_i - \theta_i\|_{\Sigma^{-1}} > b_i$ then it suffices to define $\tilde{\lambda}$ as

$$\tilde{\lambda}_p := \begin{cases} \theta_p & \text{if } p \neq i, \\ \theta_{\tilde{i}} & \text{else,} \end{cases}$$

for $\tilde{i} \in S$, to ensure that $\tilde{\lambda} \in V_i$ and

$$\begin{aligned} D(w, \tilde{\lambda}; \theta) &= w_i \|\theta_{\tilde{i}} - \theta_i\|_{\Sigma^{-1}}^2 + \sum_{k \neq i} w_k \|\lambda_k - \theta_k\|_{\Sigma^{-1}}^2 \\ &\leq w_i b_i^2 \\ &< w_i \|\lambda_i - \theta_i\|_{\Sigma^{-1}}^2 \\ &< D(w, \lambda; \theta), \end{aligned}$$

where the second line also uses the fact that $\tilde{i} \in S$. Assume $\|\lambda_i - \theta_i\|_{\Sigma^{-1}} < b_i$ and that $H_i := \{k \in S : \|\lambda_k - \theta_k\|_{\Sigma^{-1}} > 2b_i\}$ is non-empty. Let us define the instance $\tilde{\lambda}$:

$$\tilde{\lambda}_k := \begin{cases} \lambda_i & \text{if } k \in H_i, \\ \lambda_k & \text{if } k \in (S \cup \{\tilde{i}\}) \setminus H_i, \\ \theta_k & \text{else.} \end{cases}$$

Since $\lambda \in V_i$, $\tilde{\lambda}$ as defined above satisfies $\tilde{\lambda}_i \neq \tilde{\lambda}_k, k \in S$, that is $\tilde{\lambda} \in V_i$ and

$$\begin{aligned} D(w, \tilde{\lambda}; \theta) &= \sum_{k \in (S \cup \{\tilde{i}\})} w_k \|\tilde{\lambda}_k - \theta_k\|_{\Sigma^{-1}}^2 \\ &= \sum_{k \in (S \cup \{\tilde{i}\}) \setminus H_i} w_k \|\lambda_k - \mu_k\|_{\Sigma^{-1}}^2 + \sum_{k \in H_i} w_k \|\lambda_i - \theta_k\|_{\Sigma^{-1}}^2 \\ &< \sum_{k \in (S \cup \{\tilde{i}\}) \setminus H_i} w_k \|\lambda_k - \theta_k\|_{\Sigma^{-1}}^2 + \sum_{k \in H_i} w_k 4b_i^2, \end{aligned}$$

which follows since $\|\lambda_i - \theta_k\|_{\Sigma^{-1}} \leq \|\lambda_i - \theta_i\|_{\Sigma^{-1}} + \|\theta_i - \theta_k\|_{\Sigma^{-1}} \leq 2b_i$. Recalling that for $k \in H_i$,

$$4b_i < \|\lambda_k - \theta_k\|_{\Sigma^{-1}}^2,$$

it follows that

$$D(w, \tilde{\lambda}; \theta) < D(w, \lambda; \theta).$$

Combining these results, we have proved that for all $\lambda \in V_i$, there exists $\tilde{\lambda} \in V_i$ whose transportation cost is not larger than that of λ and which additionally satisfies: $\max_{k \in [K]} \|\tilde{\lambda}_k - \theta_k\|_{\Sigma^{-1}} \leq 2b_i$.

To conclude, let us define $\varepsilon := \max\{2 \max_{i \notin S} b_i, \max_{i,j \in S^2} \alpha_{i,j}\}$. Using what precedes, we have proved that

$$\inf_{\lambda \in \text{Alt}(S) \cap \{\lambda : \max_k \|\lambda_k - \theta_k\|_{\Sigma^{-1}} < \varepsilon\}} D(w, \lambda; \theta) \leq \inf_{\lambda \in \text{Alt}(S)} D(w, \lambda; \theta),$$

which proves the claimed result as we have $\text{Alt}(S) \cap \{\lambda : \max_k \|\lambda_k - \theta_k\|_{\Sigma^{-1}} < \varepsilon\} \subset \text{Alt}(S)$. \square

Chapter 6

Pareto Set Identification with Linear Constraints

This chapter studies the problem of *Pareto Set Identification* (PSI) under feasibility constraints in a multi-objective bandit framework. Given a K -armed bandit with unknown mean vectors $\mu_1, \dots, \mu_K \in \mathbb{R}^d$, we aim to identify, with high confidence, the arms that are both *feasible*, *i.e.*, satisfying a known set of linear constraints, and *non-dominated* with respect to all objectives. This constrained variant of PSI is particularly relevant in applications such as clinical trials, where treatment candidates must satisfy safety or efficacy thresholds before being compared for optimality.

To the best of our knowledge, we propose the first theoretically grounded algorithms for constrained PSI in the fixed-confidence setting. Our main algorithm, e-CAPE, jointly handles feasibility assessment and Pareto set identification by adaptively allocating samples based on upper and lower confidence bounds. We establish information-theoretic lower bounds showing that e-CAPE is near-optimal in sample complexity.

This chapter is based on joint work with Émilie Kaufmann and Laura Richert, published in the proceedings of *ICML 2025*.

6.1	Introduction	156
6.2	On the complexity of constrained PSI	159
6.2.1	Constrained PSI without explainability (cPSI)	159
6.2.2	Constrained PSI with Explainability (e-cPSI)	160
6.3	Constrained Adaptive Pareto Exploration	161
6.4	Main theoretical results	165
6.4.1	Sample complexity	167
6.5	Numerical study and discussion	170
6.6	Additional proofs	174
6.6.1	Stopping time	174
6.6.2	Sample complexity: Proof of Theorem 6.4.1	182
6.6.3	Technical results	185
6.6.4	Lower bound of e-cPSI	186
6.6.5	Complexity of Best response for cPSI	192

6.1 Introduction

In previous chapters, we studied the *Pareto Set Identification* (PSI) problem in the fixed-confidence and fixed-budget settings. So far, no constraints have been imposed on admissible arms, allowing for arms that perform extremely poorly on one objective to still belong to the Pareto set if they excel on another. In many real-world applications, however, such unconstrained trade-offs are unacceptable.

A prime example arises in clinical trials, particularly in dose-finding or early-phase vaccine studies, where multiple biological markers, such as antibody titers, safety indicators, or biomarkers of immune response, must simultaneously meet acceptable thresholds. Here, the experimenter seeks to identify treatments that not only achieve a desirable balance between efficacy and safety but also satisfy explicit feasibility constraints, such as a minimum efficacy level or a maximum toxicity bound (Jennison & Turnbull 1993; Munro et al. 2021).

Similarly, in multi-objective recommender systems or material design optimization, one might be interested in finding the Pareto set of items that satisfy some minimal performance metrics or dimensional constraints (Afshari et al. 2019; Kumar et al. 2021).

This motivated us to introduce the *constrained Pareto Set Identification* (cPSI). Formally, given K arms with multivariate means $\mu_k \in \mathbb{R}^d$, we are provided with a set of known linear constraints defining a polyhedron

$$P := \{x \in \mathbb{R}^d : Ax \leq b\}, \quad A \in \mathbb{R}^{m \times d}, \quad b \in \mathbb{R}^m,$$

encoding m feasibility requirements. Any arm k such that $\mu_k \in P$ is deemed *feasible*. The learner’s objective is to identify the Pareto set of these feasible arms, denoted $\mathcal{S}_{\text{feas}}$.

We study this problem in the *fixed-confidence* setting, where the learner must output the correct constrained Pareto set with probability at least $1 - \delta$ while minimizing the expected number of samples. This formulation naturally extends the classical fixed-confidence PSI framework introduced in Chapter 3 to handle feasibility constraints.

From PSI to Constrained PSI. A straightforward approach might proceed in two stages: first, identify the feasible arms, and then apply a PSI algorithm on the remaining subset. However, as we demonstrate in Section 6.3, such a decoupled strategy can be arbitrarily inefficient, since feasibility and dominance information are statistically intertwined. This observation motivates algorithms that integrate both tasks adaptively.

Explainable Constrained PSI. Beyond identifying the constrained Pareto set, certain applications—again notably in clinical contexts—require explicit *justifications* for excluding arms. We therefore introduce an extended variant, the *explainable constrained PSI* (e-cPSI), in which the learner must classify every non-optimal arm as either (i) infeasible or (ii) dominated by a feasible arm. This provides interpretable reasoning for every decision, a property desirable in safety-critical domains.

Related work. As illustrated throughout Chapters 2 to 4, a growing body of work has addressed Pareto set identification in the fixed-budget, fixed-confidence, and linear settings (Zuluaga, Sargent, et al. 2013; Auer et al. 2016; Zuluaga, Krause, et al. 2016; Kone, Kaufmann, et al. 2023; Karagözlü et al. 2024; Kone, Kaufmann, et al. 2025a). Despite these advances, none of these studies consider the additional challenge of handling constraints on the Pareto set.

Feasibility detection in multi-objective bandit problems was independently studied by Katz-Samuels & Scott 2018, who considered identifying arms whose mean vectors lie within a known polyhedral constraint set $P \subset \mathbb{R}^d$. However, their goal was to determine feasibility rather than to compare arms with respect to one another.

A related constrained BAI setting was analyzed by Katz-Samuels & Scott 2019, in which the learner must identify the arm maximizing a weighted combination $w^\top \mu_k$ among feasible arms. While this formulation assumes a known preference vector w to rank the arms, it fundamentally differs from the constrained PSI setting studied here, where no such scalarization is available and dominance is defined in the original multi-objective space.

The recent work of Faizal & Nair 2022 further investigated constrained BAI in dimension $d = 2$, which can be viewed as a particular instance of Katz-Samuels & Scott 2019. The authors proposed specialized algorithms and a lower bound for this setting, but extending these results to general multi-objective problems and Pareto identification remains an open challenge.

Finally, beyond pure exploration, several works have examined *constrained regret minimization*, either in single-objective problems (Lattimore & Szepesvari 2020) or in multi-objective settings (Drugan & Nowe 2013; Amani et al. 2019; D. Li et al. 2024), where the learner must select arms satisfying feasibility constraints while maximizing cumulative reward.

In contrast, our focus here is on the fixed-confidence identification of the *constrained Pareto set*, for which no existing framework or algorithm currently provides theoretical guarantees.

Contributions. This chapter makes the following contributions. We first formalize the *constrained Pareto Set Identification* (cPSI) problem and its explainable variant (e-cPSI), extending the standard PSI framework introduced in Chapters 3–4 to settings with feasibility constraints and interpretability requirements.

We then derive information-theoretic lower bounds characterizing the optimal sample complexity of both problems, thereby establishing the fundamental limits of constrained multi-objective exploration. For cPSI, we propose a gradient-ascent procedure that matches the lower bound asymptotically with a computational complexity polynomial in the number of arms.

Yet such an approach is computationally expensive for the more challenging e-cPSI, and we devote the rest of our effort to proposing and analyzing an efficient algorithm for this setting.

We introduce e-cAPE, an adaptive algorithm that jointly handles feasibility detection and Pareto set identification through coupled upper and lower confidence bounds.

We provide upper bounds on its sample complexity and show that it can be significantly smaller than that of two-stage or racing-based baselines, while remaining near-optimal for a broad class of instances.

Learning model. We denote by K the number of arms, by $d \geq 1$ the number of objectives, and by $[K] := \{1, \dots, K\}$ the index set of arms. Let \mathcal{D} be a family of distributions over \mathbb{R}^d . Given a multivariate bandit instance $\nu := (\nu_1, \dots, \nu_K) \in \mathcal{D}^K$, we denote its mean vectors by $\mu := (\mu_1, \dots, \mu_K) \in \mathcal{I}^K$, where $\mathcal{I} \subset \mathbb{R}^d$ and $\mu_k := \mathbb{E}_{X \sim \nu_k}[X]$. In the sequel, μ and ν will be used interchangeably.

We assume that \mathcal{D} is either a family of marginally subgaussian distributions,¹ used for algorithmic analysis, or a parametric family of multivariate Gaussian distributions, used when deriving lower bounds. At each round $t = 1, 2, \dots$, the learner selects an arm $A_t \in [K]$ (based on past observations) and observes an outcome $Z_t \sim \nu_{A_t}$, independent of previous samples. We denote by $\mathcal{H}_t := \sigma(A_1, Z_1, \dots, A_t, Z_t)$ the σ -algebra generated by the history up to time t , and by \mathbb{P}_ν and \mathbb{E}_μ the probability and expectation under instance ν .

In this chapter, we consider the case where the learner is additionally given a convex polyhedron

$$P := \{x \in \mathbb{R}^d : Ax \leq b\},$$

where $A \in \mathbb{R}^{m \times d}$, $b \in \mathbb{R}^m$, and m is the number of linear constraints. Given P , we study two closely related pure exploration problems:

Constrained PSI (cPSI). Identify $\mathcal{S}_{\text{feas}}$, the Pareto set of arms whose mean vectors lie within P .

Constrained PSI with explainability (e-cPSI). Identify $\mathcal{S}_{\text{feas}}$ and, in addition, provide a justification for rejecting every arm outside $\mathcal{S}_{\text{feas}}$. Each non-optimal arm must be classified as either *infeasible* (violating the constraints) or *dominated by a feasible arm*. When both conditions hold, both explanations are admissible.

While cPSI may suffice in some applications, e-cPSI is particularly relevant in domains such as clinical trials, where investigators must provide explicit reasons for excluding candidate treatments—distinguishing, for example, between infeasibility (e.g., excessive toxicity) and domination by a safer or more effective alternative. This need for interpretability makes e-cPSI our primary focus in this chapter.

We note that [Katz-Samuels & Scott 2019](#) considered a related “ δ -PAC-Explanatory” framework for identifying a feasible arm maximizing a linear combination of objectives, though their approach assumes a fixed preference vector and thus does not extend to the Pareto setting considered here.

¹Let $(X_i^c)_{c \leq d}$ denote a realization of arm i . It is marginally σ -subgaussian if for all $c \leq d$ and $\lambda \in \mathbb{R}$, $\mathbb{E}[e^{\lambda(X_i^c - \mu_i^c)}] \leq e^{\lambda^2 \sigma^2 / 2}$.

Additional notation. For any $x \in \mathbb{R}^d$, let x^c denote its c -th component, and define the distance to a set $\mathcal{X} \subset \mathbb{R}^d$ as $\text{dist}(x, \mathcal{X}) := \inf\{\|x - y\|_2 : y \in \mathcal{X}\}$. We denote by $\text{cl}(P)$ the closure of P , by $\text{int}(P)$ its interior, and by $\partial P := \text{cl}(P) \setminus \text{int}(P)$ its boundary. For $a, b \in \mathbb{R}$, we write $(a)_+ := \max\{a, 0\}$, $a \wedge b := \min\{a, b\}$, and $a \vee b := \max\{a, b\}$.

6.2 On the complexity of constrained PSI

We now study the intrinsic difficulty of the constrained Pareto set identification problems introduced above. In particular, we establish information-theoretic lower bounds on the expected sample complexity of any δ -correct algorithm for both cPSI and e-cPSI. We begin by introducing basic definitions that formalize these two settings.

Given a polyhedron P and a bandit instance $\mu \in \mathcal{I}^K$, we define the *feasible set*

$$F(\mu) := \{k \in [K] : \mu_k \in P\},$$

i.e., the set of arms satisfying all linear constraints.

For two arms i, j , we recall that arm i is (strictly) *dominated* by j (written $i \prec j$ or $\mu_i \prec \mu_j$) if $\mu_i^c < \mu_j^c$ for all $c \in [d]$. For any subset $S \subseteq [K]$ and parameter $\lambda \in \mathcal{I}^K$, let $\text{Par}(S, \lambda)$ denote the (strict) Pareto set of $\{\lambda_i : i \in S\}$ —the arms in S not dominated by any other in S under the instance λ . With this notation, the target Pareto set of feasible arms is

$$\mathcal{S}_{\text{feas}}(\mu) := \text{Par}(F(\mu), \mu),$$

and we further define the set of arms dominated by a feasible one as

$$\text{SubOpt}(\mu) := \{i \in [K] : \exists j \in F(\mu) \text{ such that } \mu_i \prec \mu_j\}.$$

When clear from context, we simply write F , $\mathcal{S}_{\text{feas}}$, and SubOpt for $F(\mu)$, $\mathcal{S}_{\text{feas}}(\mu)$, and $\text{SubOpt}(\mu)$, respectively.

6.2.1 Constrained PSI without explainability (cPSI)

In the cPSI problem, an algorithm with stopping time τ outputs a recommendation $R_\tau = O_\tau$, which is its estimate of $\mathcal{S}_{\text{feas}}(\mu)$.

Definition 6.2.1 (Correctness for cPSI). An algorithm is said to be δ -correct for cPSI on \mathcal{D}^K if, for any instance $\nu \in \mathcal{D}^K$ with mean parameter μ , it satisfies

$$\mathbb{P}_\nu(\tau < \infty, O_\tau \neq \mathcal{S}_{\text{feas}}(\mu)) \leq \delta.$$

We are interested in δ -correct algorithms with small expected sample complexity $\mathbb{E}_\mu[\tau]$. To correctly identify $\mathcal{S}_{\text{feas}}(\mu)$, an algorithm must be able to distinguish μ from any alternative instance λ for which the optimal set differs. We denote the set of such alternatives by

$$\text{Alt}(S) := \{\lambda \in \mathcal{I}^K : \mathcal{S}_{\text{feas}}(\lambda) \neq S\}.$$

Proposition 6.2.2 (Lower bound for cPSI). *Let \mathcal{D} be the class of multivariate Gaussian distributions with identity covariance, and let $\nu \in \mathcal{D}^K$ have mean vectors $\mu := (\mu_1, \dots, \mu_K) \in (\mathbb{R}^d)^K$. For any δ -correct algorithm for cPSI on \mathcal{D} with stopping time τ , it holds that*

$$\mathbb{E}_\mu[\tau] \geq T^*(\mu) \log \left(\frac{1}{2.4\delta} \right), \quad \text{where } T^*(\mu)^{-1} := \max_{w \in \Delta_K} \inf_{\lambda \in \text{Alt}(\mathcal{S}_{\text{feas}})} \sum_{k=1}^K \frac{1}{2} w_k \|\mu_k - \lambda_k\|^2,$$

with $\text{Alt}(S) := \{ \lambda \in \mathcal{I}^K : \mathcal{S}_{\text{feas}}(\lambda) \neq S \}$.

(6.1)

The proof follows the general information-theoretic approach of [Garivier & Kaufmann 2016](#), originally developed for best-arm identification. Algorithms that approach such bounds typically rely on iterative saddle-point or game-theoretic methods ([Ménard 2019](#); [Degenne, Ménard, et al. 2020](#); [P.-A. Wang et al. 2021](#)). However, they require efficient solvers for the inner optimization problem in (6.1), which is generally non-convex.

We show in Section 6.6.5 that subgradients of this inner inf function in Eq. 6.1 can be computed in polynomial time.

Then, applying the techniques in [Ménard 2019](#) yields a computationally efficient algorithm for cPSI called Game-cPSI, whose analysis is given in [Kone, Kaufmann, et al. 2025b](#).

6.2.2 Constrained PSI with Explainability (e-cPSI)

We now turn to the more challenging e-cPSI setting, in which the learner must not only identify the feasible Pareto set but also provide a justification for rejecting each non-optimal arm. At stopping time τ , an algorithm outputs a partition

$$R_\tau = (O_\tau, S_\tau, I_\tau)$$

of the arms such that

$$i) O_\tau = \mathcal{S}_{\text{feas}}, \quad ii) S_\tau \subset \text{SubOpt}, \quad iii) I_\tau \subset \mathbb{F}^c.$$

Since some arms may be both infeasible and dominated by a feasible one, multiple correct classifications may exist. Accordingly, we define the set of valid *explanatory partitions* as

$$\mathcal{M}(P, \mu) := \left\{ (S, I) \mid S \subset \text{SubOpt}(\mu), I \subset \mathbb{F}(\mu)^c, S \cap I = \emptyset, S \cup I = (\mathcal{S}_{\text{feas}}(\mu))^c \right\}.$$

When P and μ are clear from the context, we simply write \mathcal{M} . This setting thus corresponds to a pure exploration problem with multiple correct answers, a framework also studied by [Degenne & W. Koolen 2019](#) for single-objective bandits.

Definition 6.2.3 (δ -correctness for e-cPSI). *An algorithm is δ -correct for e-cPSI on \mathcal{D}^K if, for any instance $\nu \in \mathcal{D}^K$ with mean parameter μ , it outputs a partition (O_τ, S_τ, I_τ) of $[K]$ such that*

$$\mathbb{P}_\nu(\tau < \infty \text{ and } \neg(O_\tau = \mathcal{S}_{\text{feas}}, (S_\tau, I_\tau) \in \mathcal{M})) \leq \delta.$$

We again focus on δ -correct algorithms with small expected sample complexity $\mathbb{E}_\mu[\tau]$. To derive a lower bound, we first define the set of alternative instances where a given partition (S, I) would not be a valid answer:

$$\text{Alt}(S, I) := \{ \lambda \in \mathcal{I}^K : (S, I) \notin \mathcal{M}(P, \lambda) \}.$$

Proposition 6.2.4 (Lower bound for e-cPSI). *Let \mathcal{D} be the class of multivariate normal distributions with identity covariance, and let $\nu \in \mathcal{D}^K$ have mean parameters $\mu \in \mathcal{I}^K$. For any δ -correct algorithm for e-cPSI on \mathcal{D} with stopping time τ ,*

$$\liminf_{\delta \rightarrow 0} \frac{\mathbb{E}_\mu[\tau]}{\log(1/\delta)} \geq T_{\mathcal{M}}^*(\mu) := \min_{(S, I) \in \mathcal{M}} T(\mu, S, I), \quad \text{where}$$

$$T(\mu, S, I)^{-1} := \max_{w \in \Delta_K} \inf_{\lambda \in \text{Alt}(S, I)} \sum_{k=1}^K \frac{w_k}{2} \|\mu_k - \lambda_k\|^2. \quad (6.2)$$

This asymptotic lower bound extends the game-theoretic technique introduced by [Degenne & W. Koolen 2019](#) for pure exploration with multiple correct answers. For any valid partition $(S, I) \in \mathcal{M}$, the quantity $T(\mu, S, I)$ captures the information-theoretic cost required to identify (S, I) as the correct explanation. The bound implies that the minimal achievable sample complexity of any δ -correct algorithm is governed by the easiest valid partition to identify.

Even if μ were known, evaluating $T_{\mathcal{M}}^*(\mu)$ remains computationally challenging, as it involves solving the inner inf problem in (6.2) for every $(S, I) \in \mathcal{M}$. The set \mathcal{M} may have size up to 2^n with $n = |\text{SubOpt} \cap \text{F}^c|$, which can grow exponentially with K . Consequently, directly implementing the optimal strategy is infeasible for large-scale problems.

[Degenne & W. Koolen 2019](#) proposed an optimal algorithm for similar single-objective problems, but it requires enumerating all possible answers and solving the inner optimization at each step. In our setting, such enumeration over the power set of $\text{SubOpt} \cap \text{F}^c$ would be computationally prohibitive. Instead, in the next section, we introduce a practical alternative: an adaptive confidence-bound algorithm, e-cAPE, which efficiently balances exploration between feasibility testing and Pareto set identification, while achieving near-optimal sample complexity in many instances.

6.3 Constrained Adaptive Pareto Exploration

We now introduce the main algorithm of this chapter, the *Explainable Constrained Adaptive Pareto Exploration* (e-cAPE), which efficiently balances the detection of feasible arms and the identification of the constrained Pareto set. Similar to the APE algorithm presented in Chapter 3, e-cAPE relies on confidence regions around pairwise dominance and feasibility quantities. Throughout this section, all arms are assumed to have marginally σ -subgaussian rewards.

Preliminaries. As in previous chapters, we characterize pairwise dominance using the operators

$$M(i, j) := \max_{c \leq d} [\mu_i^c - \mu_j^c], \quad m(i, j) := \min_{c \leq d} [\mu_j^c - \mu_i^c], \quad (6.3)$$

introduced in Chapter 3. These quantities measure how far arm i is from being dominated by j : note that $M(i, j) > 0$ if and only if $\mu_i \not\preceq \mu_j$. They play a central role in constructing empirical dominance tests.

In the constrained setting, we must also assess whether each arm satisfies the feasibility constraints. Following [Katz-Samuels & Scott 2018](#), we define the signed feasibility margin:

$$\eta_i := \begin{cases} \text{dist}(\mu_i, P), & \text{if } \mu_i \notin P, \\ \text{dist}(\mu_i, P^c), & \text{if } \mu_i \in P, \end{cases} \quad (6.4)$$

which measures the distance of μ_i to the constraint boundary.

Empirical Estimates and Confidence Bounds. At each round t , we denote by $N_{t,i}$ the number of times arm i has been pulled and by $\hat{\mu}_{t,i}$ its empirical mean. We write $M(i, j; t)$, $m(i, j; t)$, and $\eta_i(t)$ for the empirical counterparts of (6.3) and (6.4). The empirical feasible set is

$$F_t := \{i : \hat{\mu}_{t,i} \in P\},$$

and its empirical Pareto set is given by

$$O_t := \text{Par}(F_t, \hat{\mu}_t).$$

To control uncertainty, we define a sequence of high-probability events:

$$\mathcal{E}_t := \left\{ \forall i \in [K], \|\hat{\mu}_{t,i} - \mu_i\|_\infty \leq \beta_i(t, \delta) \text{ and } \|\hat{\mu}_{t,i} - \mu_i\|_2 \leq U_i(t, \delta) \right\}, \quad \mathcal{E} := \bigcap_{t \geq 1} \mathcal{E}_t,$$

with confidence bonuses of the form

$$\beta_i(t, \delta) = \sqrt{\frac{2\sigma^2 f(t, \delta)}{N_{t,i}}}, \quad U_i(t, \delta) = \sqrt{\frac{2\sigma^2 g(t, \delta)}{N_{t,i}}},$$

where f and g are chosen so that $\mathbb{P}_\nu(\mathcal{E}) \geq 1 - \delta$.

Lemma 6.3.1 (Concentration of empirical quantities). *Under \mathcal{E}_t , for all arms i, j :*

- (i) $|M(i, j) - M(i, j; t)| \leq \beta_i(t, \delta) + \beta_j(t, \delta)$ and $|m(i, j) - m(i, j; t)| \leq \beta_i(t, \delta) + \beta_j(t, \delta)$;
- (ii) for any $\mathcal{X} \subset \mathbb{R}^d$, $|\text{dist}(\hat{\mu}_{t,i}, \mathcal{X}) - \text{dist}(\mu_i, \mathcal{X})| \leq U_i(t, \delta)$.

This leads to the confidence intervals

$$M^\pm(i, j; t) = M(i, j; t) \pm (\beta_i(t, \delta) + \beta_j(t, \delta)), \quad m^\pm(i, j; t) = m(i, j; t) \pm (\beta_i(t, \delta) + \beta_j(t, \delta)).$$

In particular, under \mathcal{E}_t , $M^-(i, j; t) > 0$ guarantees that $\mu_i \not\preceq \mu_j$, while $M^+(i, j; t) < 0$ implies that $\mu_i \prec \mu_j$.

Feasibility Detection. For any arm i , introducing the quantity $\gamma_i(t) := \frac{1}{2\sigma^2} N_{t,i} \eta_i^2(t) - g(t, \delta)$, one can show using Lemma 6.3.1 that when $\gamma_i(t) > 0$, the vectors $\hat{\mu}_{t,i}$ and μ_i can either both belong to P or to P^c , with high confidence. We let $G_t := \{i \in F_t^c \mid \gamma_i(t) < 0\}$ be the subset of F_t^c of empirically infeasible arms that cannot be confidently ruled out as infeasible at time t .

Stopping and Recommendation Rules. An algorithm for e-cPSI should stop as soon as any confidently valid answer has been identified. To introduce the stopping rule, we first define $Z_1^F(t) := \min_{i \in O_t} \gamma_i(t)$ and $Z_1^{\text{PS}}(t) := \min_{i \in O_t} \min_{j \in O_t \setminus \{i\}} [M^-(i, j; t)]$. We further define

$$Z_1(t) := \min\{Z_1^F(t), Z_1^{\text{PS}}(t)\}. \quad (6.5)$$

When $Z_1^F(t)$ is positive, arms in O_t can be proved to be feasible, and having $Z_1^{\text{PS}}(t)$ positive is sufficient to guarantee that arms in O_t are not dominated by each other. Thus, when $Z_1(t) > 0$, arms in O_t are all confidently feasible and non-dominated by each other. However, this alone is not sufficient to ensure that $O_t = \mathcal{S}_{\text{feas}}$. We also define $\xi_i(t) := \max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t)$ and note that $\xi_i(t) > 0$ implies that there exists $j \neq i \in (F_t \cup G_t)$ such that $\mu_i \prec \mu_j$. Therefore, $\mathbb{1}(i \notin F_t) ((\gamma_i(t) \vee \xi_i(t))) > 0$ implies that either arm i is infeasible or it is dominated by an arm of $(F_t \cup G_t)$. We then introduce

$$Z_2(t) := \min_{i \in [K] \setminus O_t} [\xi_i(t) \mathbb{1}(i \in F_t) + (\gamma_i(t) \vee \xi_i(t)) \mathbb{1}(i \notin F_t)], \quad (6.6)$$

for which a positive value guarantees that each arm outside of O_t can be confidently classified as either infeasible or dominated by another arm of $(F_t \cup G_t)$. We will prove that having both $Z_1(t) \geq 0$ and $Z_2(t) \geq 0$ for some t is sufficient to prove that $O_t = \mathcal{S}_{\text{feas}}$ and to identify a correct answer $(S_t, I_t) \in \mathcal{M}$. Thus, we define the stopping time τ of our algorithm as

$$\tau := \inf\{t \geq 1 \mid Z_1(t) \geq 0 \text{ and } Z_2(t) \geq 0\}. \quad (6.7)$$

At stopping, the algorithm recommends the partition (O_τ, S_τ, I_τ) of $[K]$ with $O_\tau = \text{Par}(F_\tau, \hat{\mu}_\tau)$, $S_\tau := (F_\tau \cup G_\tau) \setminus O_\tau$, a set of arms that will be shown to be dominated by some arms in O_τ , and $I_\tau := G_\tau^c \cap F_\tau^c$, a set of arms deemed infeasible.

Lemma 6.3.2 (Correctness). *On the event \mathcal{E} , if e-cAPE outputs (O_τ, S_τ, I_τ) , we have $O_\tau = \mathcal{S}_{\text{feas}}$ and $(S_\tau, I_\tau) \in \mathcal{M}(P, \mu)$.*

The above result holds regardless of the sampling rule that is used, but the sampling rule is crucial to stop early.

Sampling Rule. The challenge in designing an efficient sampling rule for constrained PSI lies in efficiently balancing the information about feasibility and Pareto dominance.

Following the "top-two" philosophy of LUCB (Kalyanakrishnan et al. 2012) and its Pareto extensions from Chapter 3, we select at each round a *leader* arm b_t and a *challenger* c_t . The leader is chosen as the arm that currently limits the stopping condition:

$$b_t := \begin{cases} \text{minimizer of (6.6),} & \text{if } Z_2(t) < 0, \\ \text{minimizer of (6.5),} & \text{otherwise.} \end{cases} \quad (6.8)$$

Intuitively, if $Z_2(t) < 0$, there remains uncertainty about an arm outside O_t being infeasible or dominated; if $Z_1(t) < 0$, some arm in O_t is not yet confidently verified as feasible or mutually non-dominated.

The challenger c_t is the arm most likely to invalidate the current status of b_t :

$$c_t := \operatorname{argmin}_{j \in (F_t \cup G_t) \setminus \{b_t\}} M^-(b_t, j; t).$$

Pulling both b_t and c_t refines the dominance and feasibility estimates most relevant to the current uncertainty, thereby increasing $\min\{Z_1(t), Z_2(t)\}$ and accelerating convergence to the stopping condition.

Algorithm 6.1: e-CAPE: Explainable Constrained Adaptive Pareto Exploration

Require: risk parameter $\delta \in (0, 1)$

```

1 Initialize: Pull each arm once; set  $t \leftarrow K + 1$ 
2 while true do
   // Feasibility and empirical Pareto set estimation
3   Compute the empirical feasible set  $F_t = \{i : \hat{\mu}_{t,i} \in P\}$ ;
4   Compute
      
$$O_t := \{i \in F_t \mid \forall j \in F_t, \hat{\mu}_{t,i} \not\prec \hat{\mu}_{t,j}\}$$

      (empirical non-dominated arms).;
   // Select most ambiguous candidate
5   if  $Z_2(t) < 0$  then
6      $b_t \leftarrow$  minimizer of (6.6);
7   else
8      $b_t \leftarrow$  minimizer of  $\begin{cases} Z_1^F(t), & \text{if } Z_1^{\text{PS}}(t) \geq Z_1^F(t), \\ Z_1^{\text{PS}}(t), & \text{otherwise.} \end{cases}$ ;
   // Find the most conflicting arm
9    $c_t \leftarrow \operatorname{argmin}_{j \in (F_t \cup G_t) \setminus \{b_t\}} M^-(b_t, j; t)$ ;
10  if  $\min\{Z_1(t), Z_2(t)\} \geq 0$  then
11    return  $(O_t, (F_t \cup G_t) \setminus O_t, F_t^c \cap G_t^c)$ ;
   // Sample selected arms
12  Pull arms  $b_t$  and  $c_t$ ;  $t \leftarrow t + 1$ ;

```

6.4 Main theoretical results

We begin by introducing the complexity quantity that drives the analysis of e-CAPE. Beyond the feasibility margins η_i (Definition (6.4)), it depends on Pareto sub-optimality gaps already used for PSI in Chapters 2–3. For any nonempty $S \subset [K]$, let $\text{Par}(S, \mu)$ be the Pareto set in the unconstrained setting and recall the dominance operators $M(i, j) = \max_{c \in [d]} [\mu_i^c - \mu_j^c]$ and $m(i, j) = \min_{c \in [d]} [\mu_j^c - \mu_i^c]$ (so $M(i, j) = -m(j, i)$). For $i \notin \text{Par}(S, \mu)$, the (unconstrained) PSI gap is

$$\Delta_i(S) := \Delta_i^*(S) := \max_{j \in \text{Par}(S, \mu)} m(i, j), \quad (6.9)$$

i.e., the minimal (coordinate-wise) increase that makes i appear Pareto-optimal against $S \setminus \{i\}$. For $i \in \text{Par}(S, \mu)$, we use the standard “two-sided margin”:

$$\Delta_i(S) := \min\{\delta_i^+(S), \delta_i^-(S)\}, \quad (6.10)$$

where

$$\delta_i^+(S) := \min_{j \in \text{Par}(S, \mu) \setminus \{i\}} \min\{M(i, j), M(j, i)\}, \quad \delta_i^-(S) := \min_{j \in S \setminus \text{Par}(S, \mu)} [(M(j, i))_+ + \Delta_j(S)],$$

with the convention $\min_{\emptyset} = +\infty$. Intuitively, δ_i^+ measures the proximity of i to other Pareto arms, whereas δ_i^- captures the margin separating i from sub-optimal arms.

A candidate-answer complexity. To certify a valid explainable answer $(S, I) \in \mathcal{M}$ (cf. Section 6.2), an algorithm must: (i) identify the Pareto set of $\mathcal{S}_{\text{feas}} \cup S$; (ii) certify feasibility of $\mathcal{S}_{\text{feas}}$; and (iii) certify infeasibility of I . We quantify this joint effort by

$$C(\mu, S, I) := \sum_{i \in \mathcal{S}_{\text{feas}}} \frac{1}{\min\{\Delta_i^2(\mathcal{S}_{\text{feas}} \cup S), \eta_i^2\}} + \sum_{i \in S} \frac{1}{\Delta_i^2(\mathcal{S}_{\text{feas}} \cup S)} + \sum_{i \in I} \frac{1}{\eta_i^2},$$

and the “easiest” valid partition by

$$C_{\mathcal{M}}^*(\mu) := \min_{(S, I) \in \mathcal{M}(P, \mu)} C(\mu, S, I).$$

Theorem 6.4.1. *Let $f(t, \delta) = \log\left(\frac{4k_1 K d t^\alpha}{\delta}\right)$ and $g(t, \delta) = 4 \log\left(\frac{4k_1 K 5^d t^\alpha}{\delta}\right)$ with $k_1 > 1 + \frac{1}{\alpha-1}$ and $\alpha > 2$. Over the class \mathcal{D} of marginally σ -subgaussian arms, e-CAPE is δ -correct for any $\delta \in (0, 1)$ and, for any instance μ ,*

$$\mathbb{E}_\mu[\tau_\delta] \leq 256\sigma^2 C_{\mathcal{M}}^*(\mu) \log\left(128\sigma^2 C_{\mathcal{M}}^*(\mu) \left(\frac{4k_1 K d}{\delta}\right)^{1/\alpha}\right) + \Lambda_\alpha + 4H(\mu, F^c),$$

where

$$\Lambda_\alpha \leq \frac{2^{\alpha-1}\delta}{4k_1} \sum_{T \geq 1} (\log T + 1) \frac{f(T, \delta \cdot d/5^d) + f(T, \delta)}{T^{\alpha-1}}, \quad H(\mu, F^c) := \sum_{a \in F^c \cup \mathcal{S}_{\text{feas}}} \frac{32\sigma^2 \log(5^d/d)}{\eta_a^2}.$$

This result shows that the dominant term of the sample complexity scales as $C_{\mathcal{M}}^*(\mu) \log\left(\frac{K C_{\mathcal{M}}^*(\mu)}{\delta}\right)$. While a general, instance-wise comparison to the lower bound in Proposition 6.2.4 is delicate, the theorem below shows that there exist families of instances for which the optimal sample complexity must scale with $C_{\mathcal{M}}^*(\mu)$, implying near-optimality of e-CAPE in a worst-case sense.

Theorem 6.4.2. *There exists a class $\tilde{\mathcal{D}}$ of multivariate Gaussian instances such that any δ -correct algorithm for e-cPSI satisfies, for all μ ,*

$$\liminf_{\delta \rightarrow 0} \frac{\mathbb{E}_{\mu}[\tau_{\delta}]}{\log(1/\delta)} \geq \frac{C_{\mathcal{M}}^*(\mu)}{4}.$$

Comparison to two-stage approaches. Since $(F \setminus \mathcal{S}_{\text{feas}}, F^c)$ is always a valid explanation, we have

$$\begin{aligned} C_{\mathcal{M}}^*(\mu) &\leq C(\mu, F \setminus \mathcal{S}_{\text{feas}}, F^c) = \sum_{i \in \mathcal{S}_{\text{feas}}} \frac{1}{\Delta_i^2(F) \wedge \eta_i^2} + \sum_{i \in F \setminus \mathcal{S}_{\text{feas}}} \frac{1}{\Delta_i^2(F)} + \sum_{i \notin F} \frac{1}{\eta_i^2} \\ &\leq \underbrace{\sum_{i \in F} \frac{1}{\Delta_i^2(F)}}_{H_{\text{PSI}}(\mu)} + \underbrace{\sum_{i \in [K]} \frac{1}{\eta_i^2}}_{H_F(\mu)}. \end{aligned}$$

Here H_F is the feasibility-detection complexity (Katz-Samuels & Scott 2018), and H_{PSI} the Pareto-identification complexity on F , as studied in Chapter 3. Thus the leading term of e-CAPE is *always* no larger than that of a naive two-stage method that first learns F and then identifies $\text{Par}(F, \mu)$: e-CAPE pays a PSI cost only for $\text{SubOpt} \cap F$, and for $\text{SubOpt} \cap F^c$ it pays the cheaper of infeasibility or dominance. Figure 6.1 illustrates an instance where this yields an arbitrarily large advantage.

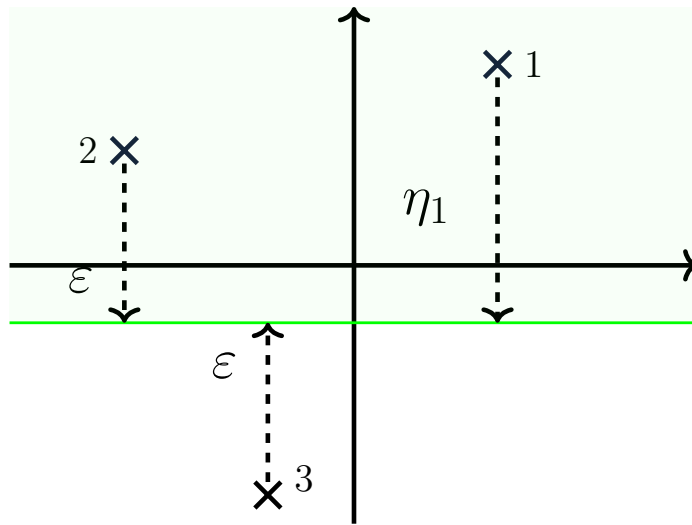


Figure 6.1: A constrained PSI instance. A two-stage approach scales as $1/\eta_1^2 + 2/\epsilon^2 + \sum_{i=1}^2 1/\Delta_i(\{1, 2\})^2$, whereas e-CAPE scales as $1/\eta_1^2 + \sum_{i=1}^3 1/\Delta_i(\{1, 2, 3\})^2$ for $\epsilon \ll 1$.

When all arms are well inside/outside P (i.e., $\eta_i \gg 1$), our bound reduces to the PSI complexities on F . Symmetrically, when PSI is easy (large Δ_i), we match feasibility-only detection rates.

Algorithms based on uniform sampling and rejection for cPSI. Racing-style methods eliminate an arm only when (i) it is confidently infeasible; (ii) it is confidently dominated by an *already* confidently feasible arm; or (iii) it is confidently optimal and does not dominate any other active arm. Condition (ii) is the bottleneck: to discard a dominated arm, one must first certify the feasibility of the dominator. Consider an instance with $F = [K]$, $\mathcal{S}_{\text{feas}} = \{1\}$, and $\text{Par}([K] \setminus \{1\}, \mu) = [K] \setminus \{1\}$. A racing algorithm must pull *every* arm until arm 1 is certified feasible before discarding any $k \neq 1$, yielding a complexity on the order of

$$\frac{K}{\eta_1^2} + \sum_{i \neq 1} \frac{1}{\Delta_i ([K])^2}.$$

In contrast, e-cAPE scales as

$$\frac{1}{\eta_1^2} + \sum_{i \neq 1} \frac{1}{\Delta_i ([K])^2},$$

removing the multiplicative factor K on the feasibility term. Our experiments in the next section corroborate this gap.

6.4.1 Sample complexity

We prove the correctness and refer to the end of the chapter for the sample complexity bound.

Proof. We show that on \mathcal{E} , when cAPE stops and returns a partition (O_τ, S_τ, I_τ) , it is a valid partition. We recall that P is the polyhedron of feasible arms, F is the set of indices of the feasible arms and

$$\mathcal{S}_{\text{feas}} := \{i \mid \mu_i \in P \text{ and } \forall j \text{ s.t. } \mu_j \in P, \mu_i \not\prec \mu_j\}$$

is the Pareto-optimal feasible set; $\text{SubOpt} := \{i \in [K] \mid \exists j \in \mathcal{S}_{\text{feas}} \setminus \{i\} : \mu_i \prec \mu_j\}$ which can be rewritten as

$$\text{SubOpt} := \{i \in [K] \mid \exists j \in F \setminus \{i\} : \mu_i \prec \mu_j\}. \quad (6.11)$$

To see this, note that if $\mu_i \prec \mu_j$ and $j \in F$, then, as $j \in F$, either $j \in \mathcal{S}_{\text{feas}}$ or there exists $j' \in \mathcal{S}_{\text{feas}}$ such that $\mu_j \prec \mu_{j'}$ and so $\mu_i \prec \mu_{j'}$. Thus

$$\exists j \in F \setminus \{i\} : \mu_i \prec \mu_j \implies \exists j \in \mathcal{S}_{\text{feas}} \setminus \{i\} : \mu_i \prec \mu_j \quad (6.12)$$

and the reverse inclusion is trivial, which justifies (6.11). Finally, it is simple to check that for the polyhedron P ,

$$\text{dist}(\mu, P^c) = \text{dist}(\mu, \partial P) \quad \forall \mu \in P. \quad (6.13)$$

To prove the correctness, we have to show that the returned partition (O_t, S_t, I_t) at $t = \tau$ satisfies $O_t = \mathcal{S}_{\text{feas}}$, $S_t \subset \text{SubOpt}$ and $I_t \subset F^c$. In this proof, we assume \mathcal{E} holds and we let $t = \tau$.

Step 1: Proving that $S_t \subset \text{SubOpt}$. Following (6.12), we have to show that

$$\forall i \in S_t, \exists j \in F \text{ such that } \mu_i \prec \mu_j. \quad (6.14)$$

By design, $S_t := (F_t \cup G_t) \setminus O_t$ and at stopping

$$Z_2(t) := \min_{i \in [K] \setminus O_t} \left[\mathbb{1}(i \in F_t) \left(\max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) + \mathbb{1}(i \notin F_t) \left(\gamma_i(t) \vee \max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) \right]$$

satisfies $Z_2(t) \geq 0$ so that

$$\forall i \in S_t, \mathbb{1}(i \in F_t) \left(\max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) + \mathbb{1}(i \notin F_t) \left(\gamma_i(t) \vee \max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) \geq 0. \quad (6.15)$$

Note that if $i \in S_t \cap G_t$, then $i \notin F_t$ and $\gamma_i(t) \leq 0$ so by (6.15),

$$\max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \geq 0,$$

which also holds if $i \in F_t$. Thus, it holds that

$$\forall i \in S_t, \max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \geq 0 \quad (6.16)$$

which on the event \mathcal{E} , using Lemma 6.3.1 and by the definition of $m(i, j)$ yields

$$\forall i \in S_t, \exists j \in (F_t \cup G_t) \setminus \{i\} \text{ such that } \mu_i \prec \mu_j. \quad (6.17)$$

However, this does not directly imply (6.14) as the stopping rule only guarantees that $O_t \subset F$ and at this point, some arms in $(F_t \cup G_t)$ could be infeasible (for the actual means). Thus, we shall prove that (6.17) holds for some $j \in O_t$. Note that by their definition, we have $F_t \cup G_t = S_t \cup O_t$. Let

$$H_t := \text{Par}(S_t, \mu) = \{i \in S_t \mid \forall j \in S_t \setminus \{i\}, \mu_i \not\prec \mu_j\},$$

the Pareto set of S_t based on the true means.

As $H_t \subset S_t$, (6.17) applies and

$$\forall i \in H_t, \exists j \neq i, j \in (F_t \cup G_t) = (S_t \cup O_t) \text{ such that } \mu_i \prec \mu_j,$$

then, as H_t is the Pareto set of S_t , arms in H_t cannot be dominated by another arm of S_t . Then, the equation above implies

$$\forall i \in H_t, \exists j \neq i, j \in O_t \text{ such that } \mu_i \prec \mu_j,$$

which as $O_t \subset F$ yields

$$\forall i \in H_t, \exists j \in F \setminus \{i\} \text{ such that } \mu_i \prec \mu_j. \quad (6.18)$$

Now, for $i \in S_t \setminus H_t$, there exists (by definition of H_t), $j \in H_t$ such that $\mu_i \prec \mu_j$ and by (6.18) there exists $j_2 \in F$ such that $\mu_j \prec \mu_{j_2}$ so $\mu_i \prec \mu_{j_2}$. Combining these results, we have proved that

$$\forall i \in S_t, \exists j \in F \setminus \{i\} \text{ such that } \mu_i \prec \mu_j,$$

hence $S_t \subset \text{SubOpt}$.

Step 2: Proving that $I_t \subset F^c$. The proof of this claim simply follows from the definition of I_t and Lemma 6.3.1. Indeed, as $I_t := G_t^c \cap F_t^c \subset O_t^c$ and by definition of G_t , it holds that for all arm $i \in I_t$, $\gamma_i(t) \geq 0$. Then by definition of γ_i , the previous yields $\eta_i(t) \geq \sqrt{\frac{2\sigma^2 g(t, \delta)}{N_{t,i}}} := U(t, \delta)$ and since $i \in F_t^c$ we invoke Lemma 6.3.1, and

$$\begin{aligned} \text{dist}(\mu_i, P) &> \text{dist}(\hat{\mu}_{t,i}, P) - U_i(t, \delta) \\ &= \eta_i(t) - U_i(t, \delta) \geq 0, \end{aligned}$$

so $\mu_i \notin P$ that is $i \in F^c$ which completes the proof that $I_t \subset F^c$.

Step 3: Proving that $O_t = \mathcal{S}_{\text{feas}}$. First, note that (O_t, S_t, I_t) is a partition of $[K]$. Since we have proved that $S_t \subset \text{SubOpt}$ and $I_t \subset F^c$, and by definition $\mathcal{S}_{\text{feas}} \cap (\text{SubOpt} \cup F^c) = \emptyset$, thus it holds that $\mathcal{S}_{\text{feas}} \subset O_t$.

Then, noting that as $Z_1(t) \geq 0$, proceeding similarly to Step 2, we have $O_t \subset F$, so at this step it holds that

$$\mathcal{S}_{\text{feas}} \subset O_t \text{ and } O_t \subset F. \quad (6.19)$$

Moreover, from $Z_1(t) \geq 0$, we derive

$$\forall i \in O_t, \forall j \in O_t \setminus \{i\}, M(i, j) \stackrel{\mathcal{E}}{\geq} M^-(i, j; t) \geq 0$$

which from the definition of $M(i, j)$ translates to

$$\forall i \in O_t, \forall j \in O_t \setminus \{i\}, \mu_i \not\prec \mu_j. \quad (6.20)$$

Therefore, as $\mathcal{S}_{\text{feas}} \subset O_t$ (6.19), we have

$$\forall i \in O_t, \forall j \in \mathcal{S}_{\text{feas}} \setminus \{i\}, \mu_i \not\prec \mu_j$$

which implies (using the transitivity of the Pareto dominance and the fact that any sub-optimal element in F is dominated by an element in $\mathcal{S}_{\text{feas}}$) that

$$\forall i \in O_t, \forall j \in F \setminus \{i\}, \mu_i \not\prec \mu_j.$$

Moreover $O_t \subset F$, so $O_t \subset \mathcal{S}_{\text{feas}}$ hence the conclusion follows as $O_t = \mathcal{S}_{\text{feas}}$.

Combining these results, we have shown that the recommendation is a valid partition of $[K]$, which concludes the proof of the correctness of Algorithm 6.1 for e-cPSI (which also implies correctness of cPSI). \square

Remark 6.1. The proof of the correctness can be adapted to generic time-uniform confidence bounds on M^\pm, m^\pm, η^\pm using an event $\mathcal{E} = \mathcal{E}_1 \cap \mathcal{E}_2$ where

$$\mathcal{E}_1 := \left\{ \forall t \geq 1, \forall (i, j) \in [K]^2, M(i, j) \in (M^-(i, j; t), M^+(i, j; t)), \right. \\ \left. \text{and } m(i, j) \in (m^-(i, j; t), m^+(i, j; t)) \right\},$$

and

$$\mathcal{E}_2 := \left(\forall t \geq 1, \forall i \in [K], \eta_i \in (\eta_i^-(t), \eta_i^+(t)) \right)$$

such that \mathcal{E} holds with probability at least $1 - \delta$.

6.5 Numerical study and discussion

Experimental setup. We now evaluate the performance of e-CAPE on instances inspired by real-world data, with a particular focus on clinical trial applications where feasibility and interpretability are essential. Throughout, we compare e-CAPE against three baselines: **(A-A)** the two-stage algorithm that combines the feasible-arm identification procedure of [Katz-Samuels & Scott 2018](#) with the Pareto set identification algorithm of [Kone, Kaufmann, et al. 2023](#); **(U)** a uniform sampling policy executed with the same stopping rule as e-CAPE; and **(R-CP)** an adaptation of the racing algorithm of [Auer et al. 2016](#) to the constrained setting (see Appendix F of [Kone, Kaufmann, et al. 2025b](#)). All algorithms use identical confidence bonuses $\beta_i(t, \delta) = \sqrt{2\sigma_i^2 \log(\log(t)/\delta)/N_{t,i}}$ and pairwise bounds $\beta_{i,j} = \sqrt{\beta_i^2 + \beta_j^2}$ with $\delta = 0.1$, yielding negligible empirical errors.

Clinical trial datasets. We begin with two experiments derived from published biomedical studies, each representing a typical multi-objective dose-finding or vaccine evaluation problem. The first dataset originates from the phase 2 clinical trial of [Mark C et al. 2013](#), which assessed the safety and efficacy of various doses (25–300mg) of Secukinumab for rheumatoid arthritis. Each arm corresponds to a dosage with measured efficacy and toxicity probabilities, and the feasible region is defined by a minimum efficacy of 40% and acceptable toxicity. This formulation mirrors early-stage dose-escalation studies in which the goal is to identify all doses that are both effective and safe, rather than a single optimum. Figure 6.2 shows the feasible region and mean responses of each dosage level.

The second dataset is based on the CovBoost study of [Munro et al. 2021](#), which compared 20 COVID-19 booster strategies across multiple immunological indicators. Following [Crepon et al. 2024](#), we retain three outcomes—neutralizing antibody titre (NT₅₀), immunoglobulin G (IgG), and cellular response—forming a three-dimensional response space. The feasible region corresponds to vaccine strategies achieving IgG levels above 8.25 titre. Figure 6.3 illustrates the feasible region and the mean response of each vaccine candidate.

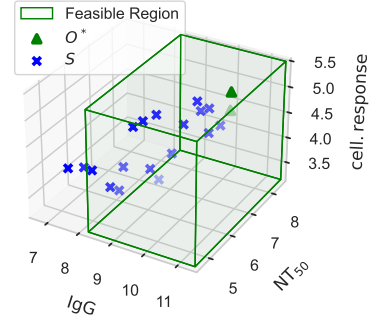
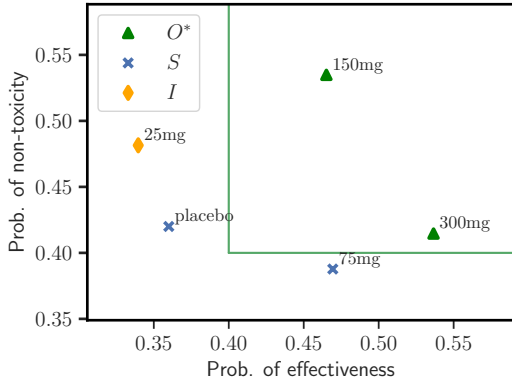


Figure 6.2: Average response in the Secukinumab dose-finding trial. The feasible region corresponds to efficacy $\geq 40\%$ and acceptable toxicity.

Figure 6.3: Average response in the CovBoost vaccine trial. The feasible region is defined by an IgG response above 8.25 titer.

Results. Figure 6.4 reports the empirical sample complexity over 500 runs for both datasets (the uniform sampling algorithm was omitted for CovBoost due to excessive runtime). In these experiments, e-cAPE consistently achieves lower sample complexity than its competitors. In the Secukinumab case, the 75mg dose lies near the feasibility boundary yet is dominated by the 150mg arm. While the two-stage algorithm (A-A) incurs additional cost verifying feasibility for such borderline arms, e-cAPE integrates feasibility and dominance testing, discarding them earlier and more efficiently. This behavior is fully consistent with the theoretical analysis of Section 6.2 which shows that decoupling feasibility from dominance can lead to sub-optimal sample allocation.

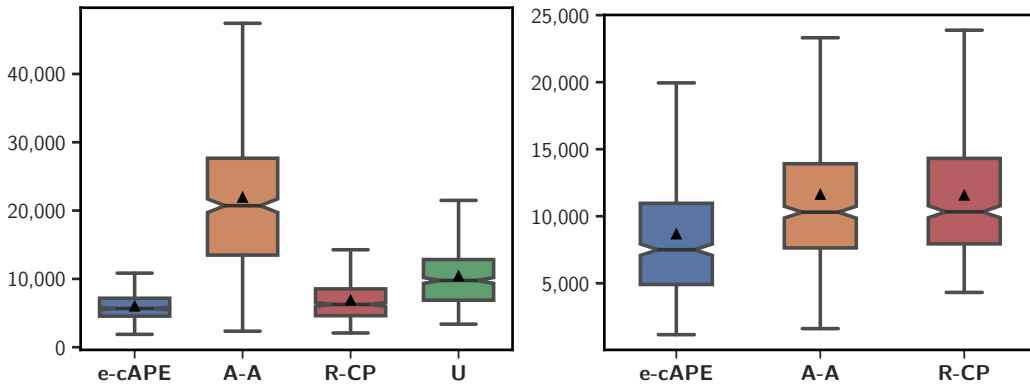


Figure 6.4: Empirical sample complexity averaged over 500 runs for the Secukinumab (left) and CovBoost (right) experiments.

Synthetic benchmarks. To further assess robustness, we test e-cAPE on synthetic three-dimensional instances defined over different feasible regions: an ordered polyhedron $P_1 = \{x \in \mathbb{R}^3 : x_i \leq x_{i+1}, i \in [2]\}$, a simplex $P_2 = \{x \in \mathbb{R}^3 : x_i \geq 0, \sum_i x_i \leq 1\}$, and a cube $P_3 = \{x \in \mathbb{R}^3 : 0.15 \leq x_i \leq 1, i \in [3]\}$. Figure 6.5 shows the geometry of these

instances, and Table 6.1 reports the sample complexity averaged over 500 runs. In all cases, e-cAPE achieves the best performance, often reducing the number of required samples by up to a factor 3 compared to the two-stage baseline. Interestingly, the uniform sampling strategy combined with our stopping rule sometimes approaches the two-stage baseline’s performance, confirming that some of the gain of e-cAPE stems from its improved stopping condition rather than aggressive sampling alone.

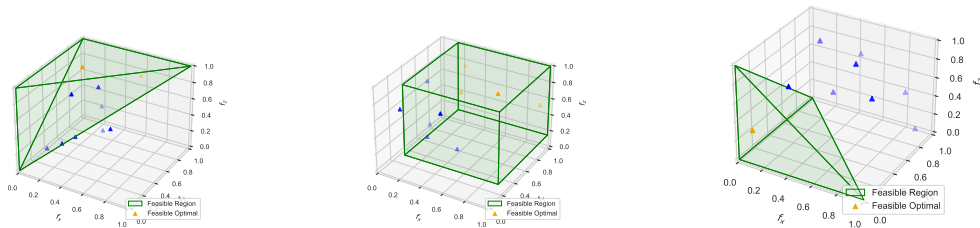


Figure 6.5: Synthetic constrained PSI instances with (left to right): ordered polyhedron, cube, and simplex. The green region denotes feasibility.

Table 6.1: Empirical sample complexity averaged over 500 runs.

Experiment	e-cAPE		MD-APT+APE (A-A)		Uniform (U)	
	Mean	Std	Mean	Std	Mean	Std
Simplex	2679	517	4886	873	12156	2428
Hypercube	19034	3158	59841	10186	56184	14392
Ordered polyhedron	4313	693	7183	792	11895	1932

Challenging instances for racing algorithms. To further highlight the limitations of elimination-based methods, we constructed a “hard” constrained PSI instance designed to challenge racing algorithms. In this two-dimensional Bernoulli setting, all arms are feasible but only one is Pareto-optimal, while the feasible region is defined by a half-space aligned with the y-axis. We considered $K \in \{5, 10\}$ arms, set $\delta = 0.01$, and averaged results over 250 runs with negligible empirical error. As anticipated from our theoretical discussion, racing algorithms must first certify that the dominating arm is feasible before discarding any dominated arm, which leads to a sample complexity that scales roughly linearly with K . In contrast, e-cAPE avoids this redundancy by jointly evaluating feasibility and dominance, resulting in substantially smaller sample requirements. This gap becomes even more pronounced as the number of arms increases, as illustrated in Figure 6.6 and Figure 6.7.

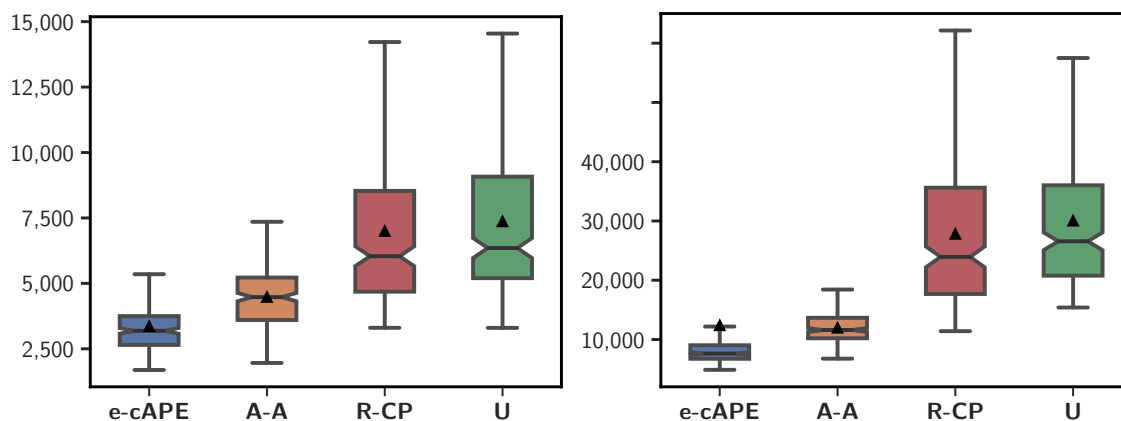


Figure 6.6: Empirical distribution of the sample complexity on the difficult constrained 5-armed instance. **Figure 6.7:** Empirical distribution of the sample complexity on the difficult constrained 10-armed instance.

Discussion. The advantage of e-cAPE is most pronounced in clinical scenarios where some candidate treatments lie close to safety or efficacy thresholds. By balancing feasibility testing and dominance exploration within a unified confidence-based framework, e-cAPE efficiently identifies all clinically acceptable strategies while explicitly explaining why others are rejected.

This feature of explainability is crucial in medical and regulatory settings, where investigators must justify exclusion decisions. In contrast to traditional racing or two-stage procedures, e-cAPE produces interpretable partitions of candidate arms into feasible, dominated, and infeasible subsets. Such interpretability supports transparent and data-driven decision-making during early-stage trials.

Conclusion. This chapter introduced the problem of constrained Pareto Set Identification (cPSI) and its explainable variant (e-cPSI), established information-theoretic lower bounds, and proposed e-cAPE, an algorithm achieving near-optimal performance in both theory and practice. Empirical results on both real and synthetic datasets confirm its efficiency and interpretability, especially in clinically motivated scenarios. Looking ahead, an exciting avenue for future research lies in designing adaptive exploration strategies that tolerate limited constraint violations during learning. Such extensions would be particularly relevant for ethical clinical trial design, where safety must be ensured throughout the study while still enabling rapid, data-efficient identification of promising treatment candidates.

6.6 Additional proofs

6.6.1 Stopping time

We upper bound the expected stopping time of Algorithm 6.1. We recall that P is fixed and known to the algorithm,

$$\mathcal{M}(P, \mu) := \{(S, I) : S \subset \text{SubOpt}, I \subset F^c \text{ and } S \cap I = \emptyset, S \cup I = (\mathcal{S}_{\text{feas}})^c\}$$

is the set of correct answers. The idea of the proof is to show that if the algorithm has not stopped after round t and \mathcal{E}_t holds, then at least one of b_t, c_t has not been sufficiently pulled *w.r.t.* the cost of identifying any correct response of \mathcal{M} . More precisely, we introduce the sets

$$W_t^1(S) := \{i \in \mathcal{S}_{\text{feas}} : \Delta_i(\mathcal{S}_{\text{feas}} \cup S) \leq 4\beta_i(t, \delta) \text{ or } \eta_i \leq 2U_i(t, \delta)\}, \quad (6.21)$$

$$W_t^2(I) := \{i \in I : \eta_i \leq 2U_i(t, \delta)\}, \quad (6.22)$$

$$W_t^3(S) := \{i \in S : \Delta_i(\mathcal{S}_{\text{feas}} \cup S) \leq 4\beta_i(t, \delta)\}, \quad (6.23)$$

and finally,

$$W_t(S, I) = W_t^1(S) \cup W_t^2(I) \cup W_t^3(S).$$

We may omit the dependency in t when it is clear from the context. At round t , any arm that belongs to $W_t(S, I)$ is called under-explored *w.r.t.* (S, I) . In the sequel, as $\mathcal{S}_{\text{feas}}$ is fixed and uniquely determined, we simplify notation: given $S \subset \text{SubOpt}$, we write $\Delta_i(S)$ to denote $\Delta_i(S \cup \mathcal{S}_{\text{feas}})$.

Proposition 6.6.1. *If Algorithm 6.1 has not stopped at time t and \mathcal{E}_t holds then for all $(S, I) \in \mathcal{M}$, $\{b_t, c_t\} \cap W_t(S, I)$ is non-empty.*

Proof of Proposition 6.6.1

In the remaining of this proof, we fix a correct answer $(S, I) \in \mathcal{M}$ and we show that if the algorithm has not stopped at time t and \mathcal{E}_t holds then one of b_t, c_t is under-explored *i.e.*, $\{b_t, c_t\} \cap W_t(S, I) \neq \emptyset$. First, we address the case $(F_t \cup G_t) \setminus \{b_t\} = \emptyset$. Note that $(F_t \cup G_t) \setminus \{b_t\} = \emptyset$, implies that arms in $[K] \setminus \{b_t\}$ can all be classified as confidently infeasible at round t . Then,

$$\mathcal{S}_{\text{feas}} = \begin{cases} \{b_t\} & \text{if } b_t \in \mathcal{S}_{\text{feas}} \\ \emptyset & \text{else.} \end{cases}$$

In both cases, the gap of b_t involves the quantity η_{b_t} . Assume $b_t \in F$, then $b_t \in \mathcal{S}_{\text{feas}}$.

Case 1.a If $b_t \in F_t$ then, as $G_t = \emptyset$ and $F_t \cup G_t = \{b_t\}$, we have

$$\begin{aligned}
 Z_2(t) &:= \min_{i \in [K] \setminus O_t} \left[\mathbb{1}(i \in F_t) \left(\max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) + \right. \\
 &\quad \left. \mathbb{1}(i \notin F_t) \left(\gamma_i(t) \vee \max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) \right] \\
 &= \min_{i \neq b_t} \left[\mathbb{1}(i \in F_t) \left(\max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) + \right. \\
 &\quad \left. \mathbb{1}(i \notin F_t) \left(\gamma_i(t) \vee \max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) \right] \\
 &= \min_{i \neq b_t} \left[\mathbb{1}(i \notin F_t) \left(\gamma_i(t) \vee \max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) \right] \quad (\text{since } F_t = \{b_t\}) \\
 &\geq 0
 \end{aligned}$$

which follows by definition of G_t , as $G_t = \emptyset$ and $F_t = \{b_t\}$. So if the algorithm has not stopped, $Z_1(t) < 0$. In this case, as O_t is a singleton, we simply have

$$Z_1(t) = \gamma_{b_t}(t) < 0$$

(the reader can check the formula of $Z_1(t)$ as $F_t = \{b_t\}$, we recall $\min_{\emptyset} = \infty$). $\gamma_{b_t}(t) < 0$ then implies $\eta_i(t) \leq U_i(t, \delta)$.

Next, assume $b_t \in F$, then $b_t \in \mathcal{S}_{\text{feas}}$ and

$$\begin{aligned}
 \eta_{b_t} &= \text{dist}(\mu_{b_t}, P^c) \stackrel{\varepsilon_t}{\leq} \text{dist}(\hat{\mu}_{t, b_t}, P^c) + U_{b_t}(t, \delta) \\
 &= \eta_i(t) + U_{b_t}(t, \delta) \\
 &\leq 2U_{b_t}(t, \delta)
 \end{aligned}$$

where the first inequality follows from Lemma 6.3.1. Thus, if $b_t \in F$ we have $b_t \in W_t^1(S)$. Now we assume $b_t \notin F$. We observe that

$$\begin{aligned}
 \eta_{b_t} &= \text{dist}(\mu_{b_t}, P) \stackrel{\varepsilon_t}{\leq} \text{dist}(\hat{\mu}_{t, b_t}, P) + U_{b_t}(t, \delta) \\
 &= 0 + U_{b_t}(t, \delta)
 \end{aligned}$$

which follows as by assumption $b_t \in F_t$. So we have $b_t \in W_t^2(I)$.

Case 1.b If $b_t \in G_t$, then by definition $\gamma_{b_t}(t) = \text{dist}(\hat{\mu}_{t, i}, P) - U_i(t, \delta) \leq 0$. By Lemma 6.3.1, either $b_t \notin F$ and

$$\begin{aligned}
 \eta_{b_t} &= \text{dist}(\mu_{b_t}, P) \stackrel{\varepsilon_t}{\leq} \text{dist}(\hat{\mu}_{t, b_t}, P) + U_{b_t}(t, \delta) \\
 &\leq 2U_{b_t}(t, \delta)
 \end{aligned}$$

or $b_t \in F$ (so $\mathcal{S}_{\text{feas}} = \{b_t\}$) and

$$\begin{aligned}
 \eta_{b_t} &= \text{dist}(\mu_{b_t}, P^c) \stackrel{\varepsilon_t}{\leq} \text{dist}(\hat{\mu}_{t, b_t}, P^c) + U_{b_t}(t, \delta) \\
 &= 0 + U_{b_t}(t, \delta).
 \end{aligned}$$

Therefore, either $b_t \notin F$ and $b_t \in W_t^2(I)$ or $b_t \in F = \{b_t\} = \mathcal{S}_{\text{feas}}$ and $b_t \in W_t^1(S)$. This concludes the analysis for the case $(F_t \cup G_t) = \{b_t\}$.

In the sequel, we assume at round t , $(F_t \cup G_t) \setminus \{b_t\} \neq \emptyset$. For this proof, we treat different cases for b_t, c_t in a series of lemmas summarized in the Table 6.2 which covers all the possible cases that could happen regarding b_t, c_t .

Case	Reference
$b_t \in S$, and $c_t \in \mathcal{S}_{\text{feas}}$	Lemma 6.6.5
$b_t \in \mathcal{S}_{\text{feas}}$ and $c_t \in \mathcal{S}_{\text{feas}}$	Lemma 6.6.4
$b_t \in I$ or $c_t \in I$	Lemma 6.6.2
$c_t \in S$	Lemma 6.6.3

Table 6.2: References to the exhaustive list of cases analyzed

The following lemmas are proved under the condition that cAPE has not stopped and \mathcal{E}_t holds.

Lemma 6.6.2. *If the cAPE has not stopped and $b_t \in I$ or $c_t \in I$ then one of them satisfies $\eta_i \leq 2U_i(t, \delta)$, i.e., $W_t^2(I) \cap \{b_t, c_t\} \neq \emptyset$.*

Proof. By design of b_t it holds that

$$\text{dist}(\hat{\mu}_{t,b_t}, P) \leq U_i(t, \delta). \quad (6.24)$$

To see this, observe that if $Z_2(t) < 0$, then, as in this case

$$b_t := \underset{i \in O_t}{\text{argmin}} \left[\gamma_i(t) \wedge \min_{j \in O_t} M^-(i, j; t) \right]$$

and $O_t \subset F_t$ we have $\hat{\mu}_{t,b_t} \in P$, so (6.24) trivially holds. If otherwise $Z_2(t) > 0$ and $Z_1(t) \leq 0$ then, recalling that in this case

$$b_t = \underset{i \in [K] \setminus O_t}{\text{argmin}} \left[\mathbb{1}(i \in F_t) \left(\max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) + \mathbb{1}(i \notin F_t) \left(\gamma_i(t) \vee \max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) \right]$$

it holds that either i) $b_t \in F_t$ so $\hat{\mu}_{t,b_t} \in P$ and (6.24) holds or ii) $b_t \notin F_t$ and $\gamma_{b_t}(t) \leq 0$ (otherwise $Z_2(t)$ would be non-negative). which by definition of γ_{b_t} yields Equation (6.24).

Next, if $b_t \in I$, then, as the event \mathcal{E}_t holds, Lemma 6.3.1 yields

$$\eta_{b_t} := \text{dist}(\mu_{b_t}, P) \leq \text{dist}(\hat{\mu}_{t,b_t}, P) + U_{b_t}(t, \delta)$$

and combining with (6.24) yields $\eta_{b_t} \leq 2U_{b_t}(t)$.

Assume instead $c_t \in I$. In this case, as $c_t \in F_t \cup G_t$ we either have $\text{dist}(\hat{\mu}_{t,c_t}, P) = 0$ (when $c_t \in F_t$) or $\gamma_{b_t}(t) \leq 0$ (when $c_t \in G_t \subset F_t^c$) i.e., $\text{dist}(\hat{\mu}_{t,b_t}, P) \leq U_{b_t}(t, \delta)$ so that (6.24) holds. Thus, the proof follows as in the previous case, leading to the conclusion that

$$\eta_{c_t} \leq 2U_{c_t}(t, \delta),$$

and so $W_t^2(I) \cap \{b_t, c_t\} \neq \emptyset$. □

Lemma 6.6.3. *If cAPE has not stopped and $c_t \in S$ then one of b_t, c_t is under-explored.*

Proof. If $c_t \in S \subset \text{SubOpt}$, then there exists $c_t^* \in \mathcal{O}^*$ such that

$$\Delta_{c_t}(S) = m(c_t, c_t^*).$$

As $c_t^* \in F$, it is easy to see that, on the event \mathcal{E}_t , c_t^* can not be declared as confidently infeasible at time t , so $c_t^* \in F_t \cup G_t$. If $b_t \neq c_t^*$ then by definition of

$$c_t := \underset{j \in (F_t \cup G_t) \setminus \{b_t\}}{\text{argmax}} m^+(b_t, j; t)$$

it holds that $m^+(b_t, c_t; t) \geq m^+(b_t, c_t^*; t)$, which by definition of $m^+(i, j; t)$ yields

$$\exists c \in [d] : \hat{\mu}_{t,c_t}^c - \hat{\mu}_{t,b_t}^c + \beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta) \geq \hat{\mu}_{t,c_t^*}^c - \hat{\mu}_{t,b_t}^c + \beta_{b_t}(t, \delta) + \beta_{c_t^*}(t, \delta)$$

thus when \mathcal{E}_t holds this implies that

$$\exists c \in [d] : \mu_{c_t^*}^c - \mu_{c_t}^c \leq 2\beta_{c_t}(t, \delta)$$

which in turn, implies

$$\Delta_{c_t}(S, I) \leq 2\beta_{c_t}(t, \delta). \tag{6.25}$$

If otherwise $b_t = c_t^* \in \mathcal{S}_{\text{feas}}$ then as the algorithm has not stopped, either $Z_2(t) < 0$ or $Z_1(t) < 0$ holds. We analyze both cases below.

Case 1: $Z_2(t) < 0$. In this case, as by design $b_t \in O_t^c$, either

- a) $b_t \in F_t$ and there exists $j \in O_t$ such that $\hat{\mu}_{t,b_t} \prec \hat{\mu}_{t,j}$, that is b_t is empirically feasible sub-optimal or
- b) $b_t \notin F_t$ holds.

In the case a), for the arm $j \in O_t$, as $\hat{\mu}_{t,b_t} \prec \hat{\mu}_{t,j}$, it holds that $M(b_t, j; t) \leq 0$ so $M^-(b_t, j; t) \leq 0$. So, as $j \in (F_t \cup G_t)$ and by definition of c_t it holds that $M^-(b_t, c_t; t) \leq 0$ i.e.,

$$M(b_t, c_t; t) \leq \beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta), \tag{6.26}$$

then, note that as $b_t = c_t^*$, $\Delta_{b_t}(S) \leq \Delta_{c_t}(S)$ (this follows from (6.10)) so

$$\begin{aligned} \max(\Delta_{b_t}(S), \Delta_{c_t}(S)) &\leq \Delta_{c_t}(S) := m(c_t, c_t^*) = m(c_t, b_t) \\ &\leq M(b_t, c_t) \\ &\stackrel{(i)}{\leq} M(b_t, c_t; t) + \beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta) \\ &\stackrel{(6.26)}{\leq} 2(\beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta)) \end{aligned}$$

where (i) follows from Lemma 6.3.1. Thus, for a_t , the least explored among b_t, c_t , it holds that

$$\Delta_{a_t}(S) \leq 4\beta_{a_t}(t, \delta),$$

which implies that $a_t \in W_t(S, I)$ in the case a). In the sub-case b), we have $b_t \notin F_t$, as $b_t = c_t^* \star \in \mathcal{S}_{feas}$, we have

$$\begin{aligned} \eta_{b_t} &= \text{dist}(\mu_{b_t}, P^c) \\ &\leq \text{dist}(\hat{\mu}_{t, b_t}, P^c) + U_{b_t}(t, \delta) \quad (\text{by Lemma 6.3.1 on } \mathcal{E}_t) \\ &\leq U_{b_t}(t, \delta), \end{aligned}$$

$b_t \in W_t(S, I)$ and this concludes the analysis of Case 1.

Case 2: $Z_1(t) < 0$ and $Z_2(t) \geq 0$. In this case, by design, $b_t \in O_t$. As

$$Z_1(t) := \min_{i \in O_t} \left[\gamma_i(t) \wedge \min_{j \in O_t \setminus \{i\}} [M^-(i, j; t)] \right]$$

is negative, and b_t is its minimizer, either a): there exists $j \in O_t$ such that $M^-(b_t, j; t) \leq 0$, which further yields

$$M^-(b_t, c_t, t) \leq 0$$

or b): $\gamma_{b_t}(t) < 0$. In the case a), proceeding similarly to Case 1.a) will lead to the conclusion that a_t , the least explored arm among b_t, c_t belongs to $W(S, I)$. The latter sub-case b) leads to $\text{dist}(\hat{\mu}_{t, b_t}, P^c) \leq U_{b_t}(t, \delta)$ so, as $b_t = c_t^* \in \mathcal{O}^*$, again we have

$$\begin{aligned} \eta_{b_t} &= \text{dist}(\mu_{b_t}, P^c) \\ &\leq \text{dist}(\hat{\mu}_{t, b_t}, P^c) + U_{b_t}(t, \delta) \quad (\text{by Lemma 6.3.1 on } \mathcal{E}_t) \\ &\leq 2U_{b_t}(t, \delta), \end{aligned}$$

that is $b_t \in W_t(S, I)$. Therefore, we conclude that one of b_t or c_t is under-explored. \square

Lemma 6.6.4. *If $b_t \in \mathcal{S}_{feas}$ and $c_t \in \mathcal{S}_{feas}$, then one of them is under-explored.*

Proof. Note that in this case, by definition of the gaps as $b_t, c_t \in \mathcal{S}_{feas}$,

$$\Delta_{b_t}(S) \leq M(b_t, c_t) \text{ and } \Delta_{c_t}(S) \leq M(b_t, c_t). \quad (6.27)$$

We analyze first the case $Z_1(t) < 0, Z_2(t) \geq 0$. By design,

$$b_t := \operatorname{argmin}_{i \in O_t} \left[\gamma_i(t) \wedge \min_{j \in O_t \setminus \{i\}} M^-(i, j; t) \right],$$

so, as $Z_1(t) < 0$, either b_t is not confidently feasible (i.e., $\gamma_{b_t}(t) \leq 0$ in which case the proof is similar to Lemma 6.6.2) or there exists $j \in O_t \setminus \{b_t\}$ such that $M^-(b_t, j; t) \leq 0$. Which as $j \in F_t$ (since $O_t \subset F_t$) and by design

$$c_t := \operatorname{argmin}_{j \in (F_t \cup G_t) \setminus \{b_t\}} [M^-(b_t, j; t)],$$

it follows that

$$M^-(b_t, c_t; t) \leq 0. \tag{6.28}$$

Then, using concentration properties of $M(i, j; t)$ (cf Lemma 6.3.1) and from (6.27), it follows

$$\begin{aligned} \max(\Delta_{b_t}, \Delta_{c_t}) &\leq M^+(b_t, c_t, t) \\ &\leq M(b_t, c_t; t) + \beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta) \\ &\leq 2(\beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta)) \end{aligned}$$

where the last inequality follows from $M^-(b_t, c_t; t) \leq 0$. For a_t the least-explored among b_t, c_t , the latter inequality implies

$$\Delta_{a_t}(S) \leq 4\beta_{a_t}(t, \delta),$$

that is $a_t \in W_t(S, I)$, which concludes our proof in the case $Z_1(t) < 0, Z_2(t) \geq 0$.

In the case $Z_2(t) < 0$, we recall that

$$b_t = \operatorname{argmin}_{i \in [K] \setminus O_t} \left[\mathbb{1}(i \in F_t) \left(\max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) + \mathbb{1}(i \notin F_t) \left(\gamma_i(t) \vee \max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) \right].$$

Assume $b_t \notin F_t$ (i.e., $\hat{\mu}_{t, b_t} \notin P$). Then by design, since $Z_2(t) < 0$ we have $\gamma_{b_t}(t) < 0$ and using concentration properties of Lemma 6.3.1, and as $b_t \in \mathcal{S}_{\text{feas}}$,

$$\begin{aligned} \eta_{b_t} &= \operatorname{dist}(\mu_{b_t}, P^c) \\ &\leq \operatorname{dist}(\hat{\mu}_{t, b_t}, P^c) + U_{b_t}(t, \delta), \quad \text{Lemma 6.3.1 on } \mathcal{E} \\ &\leq U_{b_t}(t, \delta), \end{aligned}$$

which follows as $\hat{\mu}_{t, b_t} \in P^c$. Thus $b_t \in W_t(S, I)$.

If $b_t \in F_t$ (i.e., $\hat{\mu}_{t, b_t} \in P$), as by design $b_t \in O_t^c$, then $b_t \in O_t^c \cap F_t$ so b_t is empirically feasible sub-optimal and

$$\exists j \in O_t \text{ such that } \hat{\mu}_{t, b_t} \prec \hat{\mu}_{t, j}$$

that is there exists $\exists j \in F_t \setminus \{b_t\}$ such that $M(b_t, j; t) \leq 0$, which as $M^-(b_t, j; t) \leq M(b_t, j; t)$ for all j yields

$$\exists j \in F_t \setminus \{b_t\} \text{ such that } M^-(b_t, j; t). \quad (6.29)$$

Therefore, by design of c_t ,

$$M^-(b_t, c_t; t) \leq 0 \quad (6.30)$$

so that we have with (6.27),

$$\begin{aligned} \max(\Delta_{b_t}(S), \Delta_{c_t}(S)) &\leq M^+(b_t, c_t, t) \\ &\stackrel{\mathcal{E}_t}{\leq} M(b_t, c_t; t) + \beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta) \\ &\stackrel{(6.30)}{\leq} 2(\beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta)) \end{aligned}$$

so, for a_t , the least explored arm among b_t, c_t

$$\Delta_{a_t}(S) \leq 4\beta_{a_t}(t, \delta),$$

that is $a_t \in W_t(S, I)$. □

Lemma 6.6.5. *If $b_t \in S$ and $c_t \in \mathcal{S}_{feas}$ then either b_t or c_t is under explored.*

Proof. If $Z_2(t) \leq 0$ then by design, as

$$\begin{aligned} b_t = \operatorname{argmin}_{i \in [K] \setminus \mathcal{O}_t} &\left[\mathbb{1}(i \in F_t) \left(\max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) + \right. \\ &\left. \mathbb{1}(i \notin F_t) \left(\gamma_i(t) \vee \max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) \right], \end{aligned}$$

and

$$\begin{aligned} Z_2(t) := \min_{i \in [K] \setminus \mathcal{O}_t} &\left[\mathbb{1}(i \in F_t) \left(\max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) + \right. \\ &\left. \mathbb{1}(i \notin F_t) \left(\gamma_i(t) \vee \max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) \right) \right] \end{aligned}$$

plus $c_t \in (F_t \cup G_t) \setminus \{b_t\}$, we have $m^-(b_t, c_t; t) \leq 0$. We invoke concentration properties on \mathcal{E}_t to have

$$\Delta_{b_t}(S) = m(b_t, b_t^*) \leq m^+(b_t, b_t^*; t)$$

where $b_t^* := \operatorname{argmax}_{j \in F} m(b_t, j)$ and note that on the event \mathcal{E}_t , it holds that $b_t^* \in (F_t \cup G_t)$ so by definition of c_t , $\Delta_{b_t} \leq m^+(b_t, c_t; t)$ and replacing by its expression, we further get

$$\Delta_{b_t}(S) \leq 2(\beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta)).$$

On the other side, observe that by definition of the gaps and as $b_t \in S$,

$$\Delta_{c_t}(S) \leq M(b_t, c_t)_+ + \Delta_{b_t}(S) \quad (6.31)$$

which then yields

$$\begin{aligned}\Delta_{c_t}(S) &\leq (-m(b_t, c_t; t) + \beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta))_+ + m(b_t, c_t; t) + \beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta) \\ &\leq \max(m(b_t, c_t; t) + \beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta), 2(\beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta)))\end{aligned}$$

thus

$$\Delta_{c_t}(S) \leq 2(\beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta)) \text{ and } \Delta_{b_t}(S) \leq 2(\beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta)),$$

so we conclude in the case $Z_2(t) \leq 0$.

Now, we assume $Z_2(t) > 0$; then $Z_1(t) < 0$ (otherwise, the algorithm would have stopped). We still have

$$\Delta_{b_t}(S) = m(b_t, b_t^*) \leq m^+(b_t, b_t^*; t) \leq m^+(b_t, c_t; t),$$

then, the idea is to show that $m^-(b_t, c_t, t) \leq 0$. We will prove that for any $i \in (F_t \cup G_t)$, $m^-(b_t, i, t) \leq 0$. To see this, first, note that, as $b_t \in O_t$, for any $i \in F_t$, $\hat{\mu}_{t, b_t} \not\prec \hat{\mu}_{t, i}$, thus for any $i \in F_t$, $m(b_t, i, t) \leq 0$.

Now let $i \in G_t$. By Lemma 6.6.6, there exists $j \in O_t$ such that $\hat{\mu}_{t, i} \prec \hat{\mu}_{t, j}$. Having $\hat{\mu}_{t, b_t} \prec \hat{\mu}_{t, i}$ would yield either $\hat{\mu}_{t, b_t} \prec \hat{\mu}_{t, i} \prec \hat{\mu}_{t, b_t}$ (if $j = b_t$) or by transitivity

$$\hat{\mu}_{t, b_t} \prec \hat{\mu}_{t, j},$$

with $j \neq b_t$, which is not possible for $b_t, j \in O_t$. Therefore, for any $i \in F_t$, $m(b_t, i; t) \leq 0$, so combined with the previous display,

$$\forall i \in (F_t \cup G_t), m(b_t, i; t) \leq 0.$$

In particular, $m(b_t, c_t; t) \leq 0$. Moreover, recalling that on \mathcal{E}_t ,

$$\Delta_{b_t}(S) \leq m^+(b_t, c_t; t),$$

we have $\Delta_{b_t}(S) \leq \beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta)$. On the other side, as in (6.31),

$$\Delta_{c_t}(S) \leq M(b_t, c_t)_+ + \Delta_{b_t}(S),$$

therefore,

$$\begin{aligned}\Delta_{c_t}(S) &\leq (-m(b_t, c_t; t) + \beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta))^+ + m(b_t, c_t; t) + \beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta) \\ &\leq 2(\beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta))\end{aligned}$$

Combined, we have proved that

$$\Delta_{c_t}(S) \leq \beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta) \text{ and } \Delta_{b_t}(S) \leq 2(\beta_{b_t}(t, \delta) + \beta_{c_t}(t, \delta))$$

which again makes it possible to conclude, as letting a_t the least explored among b_t and c_t ,

$$\Delta_{a_t}(S) \leq 4\beta_{a_t}(t, \delta)$$

that is $a_t \in W_t(S, I)$ □

Lemma 6.6.6. *If $Z_2(t) > 0$, then, for any $i \in G_t$, there exists $j \in O_t$ such that $\hat{\mu}_{t,i} \prec \hat{\mu}_{t,j}$.*

Proof. Recall that arms in G_t cannot yet be confidently classified as feasible or infeasible.

As $Z_2(t) > 0$, and $i \in G_t \cap O_t^c$, from the definition of $Z_2(t)$, it holds that

$$\max_{j \in (F_t \cup G_t) \setminus \{i\}} m^-(i, j; t) > 0,$$

so there exists $j \in (F_t \cup G_t) \setminus \{i\}$ such that $m^-(i, j; t) > 0$, in particular, $\hat{\mu}_{t,i} \prec \hat{\mu}_{t,j}$. Introducing

$$\mathfrak{J}_i := \{j \in F_t \setminus \{i\} : \hat{\mu}_{t,i} \prec \hat{\mu}_{t,j}\},$$

we will show that $\mathfrak{J}_i \neq \emptyset$. Let $\Omega_i := \{j \in (F_t \cup G_t) \setminus \{i\} : \hat{\mu}_{t,i} \prec \hat{\mu}_{t,j}\}$ and define $H_i := \text{Par } \Omega_i, \hat{\mu}_t$, the Pareto set of Ω_i based on the empirical means. By the result above, Ω_i is non-empty, so H_i is also non-empty.

We claim that $H_i \subset F_t$. Indeed, for $k \in H_t \cap F_t^c$ we have $k \in G_t \cap O_t^c$, then, as $Z_2(t) > 0$, as justified above for i , there exists $l \in (F_t \cup G_t) \setminus \{k\}$ such that $\hat{\mu}_{t,k} \prec \hat{\mu}_{t,l}$, however, as $k \in H_t := \text{Par } \Omega_i, \hat{\mu}_t$, the above is only possible if $l \notin \Omega_i$. Recalling that $l \in (F_t \cup G_t) \setminus \{k\}$ and $l \neq b_t$ we have $l \in (F_t \cup G_t) \setminus \{b_t\}$ so putting the above displays together, if $H_t \cap F_t^c \neq \emptyset$, there exists l and k such that

$$\hat{\mu}_{t,i} \prec \hat{\mu}_{t,k}; \hat{\mu}_{t,k} \prec \hat{\mu}_{t,l}$$

but $\hat{\mu}_{t,i} \not\prec \hat{\mu}_{t,l}$, which is not possible as Pareto dominance is transitive. Therefore, $H_t \subset F_t$ and since $H_t \neq \emptyset$, $\mathfrak{J}_i \neq \emptyset$. Thus, there exists $j \in F_t \setminus \{b_t\}$ such that $\hat{\mu}_{t,i} \prec \hat{\mu}_{t,j}$. Then, either $j \in O_t$ or there exists $j_2 \in O_t$ such that, $\hat{\mu}_{t,j} \prec \hat{\mu}_{t,j_2}$. In any case, there exists $j_3 \in O_t$ (either j or j_2) such that $\hat{\mu}_{t,i} \prec \hat{\mu}_{t,j_3}$. \square

6.6.2 Sample complexity: Proof of Theorem 6.4.1

Proof. Let $(S, I) \in \mathcal{M}$ be fixed. Let $T > 0$ be fixed and observe that

$$\min(\tau_\delta, T) \leq \left\lceil \frac{T}{2} \right\rceil + \sum_{t=\lceil T/2 \rceil}^T \mathbb{1}(t > \tau_\delta).$$

Then, assuming

$$\mathcal{E}^T := \bigcap_{\lceil T/2 \rceil \leq t \leq T} \mathcal{E}_t \tag{6.32}$$

holds, and using Proposition 6.6.1 we have

$$\begin{aligned}
 \min(\tau_\delta, T) &\leq \left\lceil \frac{T}{2} \right\rceil + \sum_{t=\lceil T/2 \rceil}^T \mathbb{1}(t > \tau_\delta) \\
 &\leq \left\lceil \frac{T}{2} \right\rceil + \sum_{t=\lceil T/2 \rceil}^T \mathbb{1}(b_t \in W_t(S, I) \vee c_t \in W_t(S, I)) \\
 &\leq \left\lceil \frac{T}{2} \right\rceil + \sum_{t=\lceil T/2 \rceil}^T \sum_{a=1}^K \mathbb{1}(b_t = a \vee c_t = a) \mathbb{1}(a \in W_t(S, I)) \\
 &= \left\lceil \frac{T}{2} \right\rceil + \sum_{a=1}^K \sum_{t=\lceil T/2 \rceil}^T \mathbb{1}(b_t = a \vee c_t = a) \mathbb{1}(a \in W_t(S, I)),
 \end{aligned}$$

then, we decompose the latter sum into three terms by defining the function R for $U \subset [K]$ as

$$R(U) := \sum_{a \in U} \sum_{t=\lceil T/2 \rceil}^T \mathbb{1}(b_t = a \vee c_t = a) \mathbb{1}(a \in W_t(S, I))$$

we have for the set $S \subset \text{SubOpt}$, using the definition of $W_t^3(S)$ for $a \in S$:

$$\begin{aligned}
 R(S) &= \sum_{a \in S} \sum_{t=\lceil T/2 \rceil}^T \mathbb{1}(b_t = a \vee c_t = a) \mathbb{1}(\Delta_a(S) \leq 4\beta_a(t, \delta)) \\
 &\leq \sum_{a \in S} \sum_{t=\lceil T/2 \rceil}^T \mathbb{1}(b_t = a \vee c_t = a) \mathbb{1}(N_{t,a} \leq 32\sigma^2 \Delta_a^{-2}(S) f(t, \delta)) \\
 &\leq \sum_{a \in S} \sum_{t=\lceil T/2 \rceil}^T \mathbb{1}(b_t = a \vee c_t = a) \mathbb{1}(N_{t,a} \leq 32\sigma^2 \Delta_a^{-2}(S) f(T, \delta)) \quad (\text{as } f \text{ is increasing})
 \end{aligned}$$

Therefore,

$$R(S) \leq \sum_{a \in S} \frac{32\sigma^2}{\Delta_a^2(S)} f(T, \delta). \quad (6.33)$$

Proceeding similarly for I , we obtain

$$R(I) \leq \sum_{a \in I} \frac{8\sigma^2}{\eta_a^2} g(T, \delta), \quad (6.34)$$

and for \mathcal{O}^* , we obtain

$$R(\mathcal{O}^*) \leq \sum_{a \in \mathcal{O}^*} \max \left(\frac{8\sigma^2 g(T, \delta)}{\eta_a^2}, \frac{32\sigma^2 f(T, \delta)}{\Delta_a^2(S)} \right). \quad (6.35)$$

Therefore, recalling the expressions of f and g

$$f(T, \delta) = \log \left(\frac{4k_1 K d T^\alpha}{\delta} \right) \text{ and } g(T, \delta) = 4 \log \left(\frac{4k_1 K 5^d T^\alpha}{\delta} \right),$$

we have

$$\begin{aligned} R(\mathcal{O}^*) &\leq \sum_{a \in \mathcal{S}_{\text{feas}}} 32\sigma^2 \max \left\{ \frac{\log \left(\frac{4k_1 K 5^d T^\alpha}{\delta} \right)}{\eta_a^2}, \frac{\log \left(\frac{4k_1 K d T^\alpha}{\delta} \right)}{\Delta_a^2(S)} \right\} \\ &\leq \sum_{a \in \mathcal{S}_{\text{feas}}} 32\sigma^2 \max \left\{ \frac{1}{\eta_a^2}, \frac{1}{\Delta_a^2(S)} \right\} \log \left(\frac{4k_1 K d T^\alpha}{\delta} \right) + \sum_{a \in \mathcal{S}_{\text{feas}}} \frac{32\sigma^2}{\eta_a^2} \log \left(\frac{5^d/d}{\delta} \right). \end{aligned}$$

Combining (6.35), (6.34) and (6.33) yields

$$\begin{aligned} R([K]) &\leq 32\sigma^2 \left(\sum_{a \in \mathcal{O}^*} \frac{1}{\min(\eta_a, \Delta_a(S))^2} + \sum_{a \in \mathcal{S}} \frac{1}{\Delta_a^2(S)} + \sum_{a \in \mathcal{I}} \frac{1}{\eta_a^2} \right) f(T, \delta) + \\ &= 32\sigma^2 C(S, I) f(T, \delta) + H(\mu, I), \end{aligned}$$

where $H(\mu, I) = \sum_{a \in \mathcal{I} \cup \mathcal{O}^*} 32\sigma^2 \frac{\log(5^d/d)}{\eta_a^2}$. Therefore, assuming (6.32), we have proved that

$$\min(\tau_\delta, T) < \left\lceil \frac{T}{2} \right\rceil + 32\sigma^2 C(S, I) f(T, \delta) + H(\mu, I)$$

Then, taking T such that the RHS is smaller than T would yield that the algorithm stops before T . Applying Lemma 6.6.7 with $a = 128\sigma^2 C(S, I)\alpha$ and $b = (4k_1 K d/\delta)^{1/\alpha}$ yields

$$T \geq 2a \log(ab) \implies a \log(bT) < T,$$

that is, letting $\tilde{C}(S, I) = 128\sigma^2 C(S, I)\alpha$,

$$\begin{aligned} T \geq 2\tilde{C}(S, I) \log \left(\tilde{C}(S, I) (4k_1 K d/\delta)^{1/\alpha} \right) &\implies 128\sigma^2 C(S, I) \log((4k_1 K d/\delta)T) < T \\ &\implies 32\sigma^2 C(S, I) f(T, \delta) < T/4. \end{aligned} \tag{6.36}$$

Moreover, we trivially have for $T > 4H(\mu, I)$, $H(\mu, I) < T/4$. Combining with (6.36), for $T > T^*(S, I) := 2\tilde{C}(S, I) \log \left(\tilde{C}(S, I) (4k_1 K d/\delta)^{1/\alpha} \right) + 4H(\mu, I)$, we have $\min(\tau_\delta, T) < T$ so the algorithm must have stopped before round T . Therefore,

$$\forall T > T^*(S, I), \tau_\delta > T \implies (\mathcal{E}^T)^c$$

that is $\mathbb{E}_\mu[\tau_\delta] \leq T^*(S, I) + \sum_{t > T^*(S, I)} \mathbb{P}_\nu((\mathcal{E}^t)^c)$. Letting $\Lambda_\alpha := \sum_{T \geq 1} \mathbb{P}_\nu((\mathcal{E}^T)^c)$ and as S, I is arbitrary, we have

$$\mathbb{E}_\mu[\tau_\delta] \leq \min_{(S, I) \in \mathcal{M}} T^*(S, I) + \Lambda_\alpha$$

and from Lemma D.3 of [Kone, Kaufmann, et al. 2025b](#),

$$\Lambda_\alpha \leq \frac{2^{\alpha-1}\delta}{4k_1} \sum_{T \geq 1} (\log(T) + 1) \left(\frac{f(T, \delta \cdot d/5^d) + f(T, \delta)}{T^{\alpha-1}} \right).$$

Finally, as $x \mapsto x \log(x)$ is increasing on $(1, \infty)$, we have

$$\mathbb{E}_\mu[\tau_\delta] \leq 256\sigma^2 C_{\mathcal{M}}^*(\mu) \log(128\sigma^2 C_{\mathcal{M}}^*(\mu)(4k_1 K d/\delta)^{1/\alpha}) + 4H(\mu, \mathbf{F}^c) + \Lambda_\alpha.$$

□

6.6.3 Technical results

We state some concentration inequalities that are used in our proofs.

Proof. Letting $\mathcal{X} \subset \mathbb{R}^d$ non-empty we have for all arm i

$$\begin{aligned} |\text{dist}(\hat{\mu}_{t,i}, \mathcal{X}) - \text{dist}(\mu_i, \mathcal{X})| &= \left| \sup_{x \in \mathcal{X}} [-\|\mu_i - x\|_2] - \sup_{x \in \mathcal{X}} [-\|\hat{\mu}_{t,i} - x\|_2] \right| \\ &\stackrel{(i)}{\leq} \sup_{x \in \mathcal{X}} \left| \|\hat{\mu}_{t,i} - x\|_2 - \|\mu_i - x\|_2 \right| \\ &\stackrel{(ii)}{\leq} \|\hat{\mu}_{t,i} - \mu_i\| \\ &\leq U_i(t, \delta) \end{aligned}$$

where the last inequality follows from \mathcal{E}_t , (i) follows from the inequality

$$|\sup f - \sup g| \leq \sup |f - g|$$

and (ii) follows from the $|\|x\| - \|y\|| \leq \|x - y\|$. □

Lemma 6.6.7. *Let $a, b \geq e$. The following statement holds:*

$$t \geq 2a \log(ab) \implies a \log(bt) \leq t.$$

Proof. First, note that $t \mapsto \log(bt)/t$ is decreasing on $[e/b, +\infty)$ and observe for $t_0 = 2a \log(ab)$, it holds that

$$\begin{aligned} a \log(bt_0) &= a \log(2ab \log(ab)) \\ &= a \log(ab) + a \log(2 \log(ab)) \\ &< 2a \log(ab) \quad (\text{as } \log(x) < x/2) \\ &= t_0 \end{aligned}$$

therefore, $a \frac{\log(t_0)}{t_0} < 1$ and as the function is decreasing

$$t \geq t_0 \implies a \frac{\log(bt)}{t} \leq a \frac{\log(bt_0)}{t_0} < 1.$$

□

6.6.4 Lower bound of e-cPSI

In this section, we prove a lower bound on the sample complexity of any δ -correct algorithm for the constrained PSI problem. As this is a pure exploration problem with multiple correct answers, it might be tempting to apply Theorem 1 of [Degenne & W. Koolen 2019](#). However, their result does not directly apply to our setting as the authors studied the regime $\delta \rightarrow 0$ and considered answers that are singletons of $[K]$ while we consider partitions of $[K]$.

We state below the result we will prove in this section. We denote by $\tilde{\mathcal{D}}$ a family of bandit instances that will become explicit by the end of the section.

Theorem 6.6.8. *Let $P := \{x \in \mathbb{R}^d \mid Ax \leq b\}$ and let $\nu \in \tilde{\mathcal{D}}^K$ with parameter μ . The following holds for any δ -correct algorithm*

$$\mathbb{E}_\mu[\tau] \min_{(S', I') \in \mathcal{M}} C(S', I') \frac{g(\delta)}{2}, \text{ where } g(\delta) \geq (1 - \delta) \log(1/\delta) - \log(4).$$

In particular,

$$\liminf_{\delta \rightarrow 0} \frac{\mathbb{E}_\mu[\tau]}{\log(1/\delta)} \geq \frac{C_{\mathcal{M}}^*(\mu)}{4}.$$

In what follows, we consider \mathcal{A}_δ , a δ -correct algorithm (cf Definition 6.2.3) and μ , a bandit instance, fixed and unknown to the algorithm. We assume that $\tau < \infty$ almost surely (otherwise the lower bound trivially holds). Let $\hat{O}, \hat{S}, \hat{I}$ be the recommendation of the algorithm after τ (random) steps and \mathcal{H}_τ the natural filtration of the stochastic process (A_t, X_t) stopped at round τ . Since \mathcal{A} is δ -correct, the following statement holds

$$\mathbb{P}_\mu \left(\hat{O} = \mathcal{S}^* \text{ and } (\hat{S}, \hat{I}) \in \mathcal{M} \right) \geq 1 - \delta. \quad (6.37)$$

We define

$$\tilde{S} := \{i \in [K] \mid \mu_i \in P \text{ and } \exists j \in \mathcal{S}^* \mid \mu_i \prec \mu_j\}$$

the set of arms that are feasible but sub-optimal. Note that alternatively, $\tilde{S} = \text{SubOpt} \cap F$. Similarly, we introduce

$$\tilde{I} := \{i \in [K] \mid \mu_i \notin P \text{ and } \nexists j \in \mathcal{S}^* \mid \mu_i \prec \mu_j\},$$

the set of arms that are unsafe but not uniformly worse than any other optimal arm, which rewrites as $\tilde{I} = F^c \cap \text{SubOpt}^c$. In particular, $\tilde{S} \cap \tilde{I} = \emptyset$ and it can be observed that for any valid answer $(S, I) \subset \mathcal{M}$, $\tilde{S} \subset S$ and $\tilde{I} \subset I$. Thus, for any δ -correct algorithm it holds that

$$\forall i \in \tilde{I}, \mathbb{P}_\mu(i \in \hat{I}) \geq 1 - \delta \quad \text{and} \quad \forall i \in \tilde{S}, \mathbb{P}_\mu(i \in \hat{S}) \geq 1 - \delta. \quad (6.38)$$

Then by definition, every arm of $\mathcal{B} := (\mathcal{S}^*)^c \cap (\tilde{S} \cup \tilde{I})^c$ can be (correctly) classified either in \hat{S} or \hat{I} , resulting in $2^{|\mathcal{B}|}$ correct answers, which can be up to 2^{K-1} correct answers.

A high probability correct answer for \mathcal{A}_δ . We describe a particular correct answer under bandit μ that is very likely to be recommended by \mathcal{A}_δ .

We build on (\tilde{S}, \tilde{I}) to construct a partition with (S, I) for which a property similar to Eq. (6.38) holds for every arm of (S, I) . To construct such instance, note that as \mathcal{A}_δ is δ -correct, for any $i \in \mathcal{B} = (\mathcal{S}^*)^c \cap (\tilde{S} \cup \tilde{I})^c$, it holds that $\mathbb{P}_\mu(i \in \hat{S} \cup \hat{I}) \geq 1 - \delta$, which yields

$$\max \left(\mathbb{P}_\mu(i \in \hat{S}), \mathbb{P}_\mu(i \in \hat{I}) \right) \geq \frac{1 - \delta}{2}. \quad (6.39)$$

Then, we define S, I as

$$S := \tilde{S} \cup \left\{ i \in C : \mathbb{P}_\mu(i \in \hat{S}) > \mathbb{P}_\mu(i \in \hat{I}) \right\}$$

and we similarly define the set

$$I := \tilde{I} \cup \left\{ i \in C \mid \mathbb{P}_\mu(i \in \hat{S}) \leq \mathbb{P}_\mu(i \in \hat{I}) \right\}.$$

By construction, (S, I) satisfies $S \cap I = \emptyset$ and $S \cup I = (\mathcal{S}^*)^c$, i.e., (S, I) is a valid answer. Moreover, one can verify that

- a) $\mathbb{P}_\mu(\hat{O} = \mathcal{S}^*) \geq 1 - \delta$
- b) for all $i \in S$, $\mathbb{P}_\mu(i \in \hat{S}) \geq \frac{1-\delta}{2}$ and
- c) for all $i \in I$, $\mathbb{P}_\mu(i \in \hat{I}) \geq \frac{1-\delta}{2}$;

that is this particular (\mathcal{S}^*, S, I) is a likely response for the algorithm \mathcal{A}_δ . Using a change of distribution lemma will allow us to derive a lower bound on the number of pulls of each arm for \mathcal{A}_δ to identify (S, I) as a correct answer. In the sequel, we restrict ourselves to bandits with multivariate normal arms and identity covariance. Therefore, each arm will be identified with its mean vector.

Lemma 6.6.9 (Kaufmann, Cappé, et al. 2016). *For all bandit models μ, λ and for any \mathcal{H}_τ -measurable event \mathcal{E} ,*

$$\sum_{i=1}^K \mathbb{E}_\mu[N_{\tau,i}] \text{KL}(\mu_i, \lambda_i) \geq \text{kl}(\mathbb{P}_\mu(\mathcal{E}), \mathbb{P}_\lambda(\mathcal{E})),$$

where $\text{KL}(\mu_i, \lambda_i)$ is the Kullback-Leibler divergence between distributions identified by parameters μ_i, λ_i and kl is the relative binary entropy.

Change of distribution. The idea now is to apply the classical change of distribution lemma. We define instances with means $\lambda^1, \dots, \lambda^K$ such that λ^i differs from μ only in the mean of arm i i.e., the bandit μ and λ^i are identical except $\lambda^i \neq \mu_i$. For such instances, applying Lemma 6.6.9 yields

$$\mathbb{E}_\mu[N_{\tau,i}] \text{KL}(\mu_i, \lambda^i) \geq \sup_{\mathcal{E} \in \mathcal{H}_\tau} \text{kl}(\mathbb{P}_\mu(\mathcal{E}), \mathbb{P}_{\lambda^i}(\mathcal{E})) \quad (6.40)$$

where we recall that $\text{KL}(\mu_i, \lambda_i^i)$ is the Kullback-Leibler divergence between multivariate normals of means μ_i, λ_i^i and identity covariance. $\text{kl}(x, y) := x \log(x/y) + (1-x) \log((1-x)/(1-y))$, and noting that $x \mapsto \text{kl}(x, y)$ is increasing on $\{x > y\}$ and decreasing on $\{x < y\}$, making an event likely under μ unlikely under λ^i will increase the RHS of Eq. (6.40). Going back to the instance (S, I) constructed earlier, we introduce the event

$$\mathcal{E}_i := \begin{cases} \{i \in \widehat{S}\} & \text{if } i \in S \\ \{i \in \widehat{I}\} & \text{if } i \in I, \\ \{\widehat{O} = \mathcal{S}^*\} & \text{if } i \in \mathcal{S}^*. \end{cases}$$

We recall that (\mathcal{O}^*, S, I) forms a partition of $[K]$. From the definition of (S, I) we have shown that for any δ -correct algorithm,

$$\mathbb{P}_\mu(\mathcal{E}_i) \geq \begin{cases} 1 - \delta & \text{if } i \in \mathcal{O} \\ \frac{1-\delta}{2} & \text{else,} \end{cases}$$

that is, such an event is very likely under bandit μ . We will choose each instance λ^i such that \mathcal{E}_i is unlikely under bandit λ^i , *i.e.*, we assume that λ^i is chosen such as for all i

$$\mathbb{P}_{\lambda^i}(\mathcal{E}_i) \leq \delta. \quad (6.41)$$

We now discuss the choice of λ^i depending on i to ensure such a property. Thus, assuming this property holds, we have for all arms i

$$\mathbb{E}_\mu[N_{\tau,i}] \text{KL}(\mu_i, \lambda_i^i) \geq \text{kl}\left(\frac{1-\delta}{2}, \delta\right)$$

which, using the KL formula for multivariate normal, yields

$$\mathbb{E}_\mu[N_{\tau,i}] \geq \frac{2}{\|\mu_i - \lambda_i^i\|_2^2} \underbrace{\text{kl}\left(\frac{1-\delta}{2}, \delta\right)}_{g(\delta)}.$$

Case 1: $i \in I$ an unsafe arm. In this case, $\mu_i \notin P$. Since P is a closed convex set, by Hilbert projection onto closed convex sets, there exists (unique) $z_i \in P$ such that $\|\mu_i - z_i\|_2 = \text{dist}(\mu_i, P)$. We now define $\lambda_i^i = z_i$. In this bandit λ_i , i is now a safe arm so, as \mathcal{A}_δ is a δ -correct algorithm it holds that

$$\mathbb{P}_{\lambda_i}(i \in \widehat{I}) \leq \delta \quad \text{i.e.,} \quad \mathbb{P}_{\lambda_i}(\mathcal{E}_i) \leq \delta,$$

which by further noting that $\|\mu_i - \lambda_i^i\|_2 = \text{dist}(\mu_i, P) = \eta_i$ results in

$$\mathbb{E}_\mu[N_{\tau,i}] \geq \frac{2}{\eta_i^2} g(\delta).$$

Case 2: $i \in S$ a sub-optimal arm. Note that in this case, there exists a subset $\Omega_i \subset \mathcal{S}^*$ such that $\mu_i \prec \mu_j$ for all $j \in \Omega_i$: it is the set of optimal arms that dominate μ_i . For any $j \in \Omega_i$ we define

$$c_j = \underset{c}{\text{argmin}} [\mu_j(c) - \mu_i(c)]$$

the coordinate for which the margin between i and j is the lowest. In particular, $m(i, j) = \mu_j(c_j) - \mu_i(c_j)$. Moreover, for any $\varepsilon > 0$, the vector defined by

$$\tilde{\mu}_i(c) = \begin{cases} \mu_i(c) + (1 + \varepsilon) m(i, j) & \text{if } c = c_j \\ \mu_i(c) & \text{else} \end{cases}$$

is not dominated by μ_j . Now let us define the vector s_i such that for any $c \in [d]$,

$$s_i(c) := \max_{j \in \Omega_i: c_j = c} m(i, j) \quad (6.42)$$

with the convention that $\max_{\emptyset} = 0$. We argue that defining

$$\lambda_i^i := \mu_i + (1 + \varepsilon)s_i,$$

we have

$$\mathbb{P}_{\lambda^i}(i \in \hat{S}) \leq \delta.$$

To see this, note that as only the mean of arm i has changed between bandits μ and λ^i , the set of safe and unsafe arms of $[K] \setminus \{i\}$ under μ and λ^i are the same. Moreover, the change of distribution ensures that under the bandit λ^i , i is not dominated by any safe arm. Therefore, $\mathbb{P}_{\lambda^i}(i \in \hat{S}) \leq \delta$. So applying what precedes and letting $\varepsilon \rightarrow 0$ yield

$$\mathbb{E}_{\mu}[N_{\tau, i}] \geq \frac{2}{\|s_i\|^2} g(\delta). \quad (6.43)$$

Case 3: $i \in \mathcal{S}^*$ the optimal set. We recall that arms in \mathcal{S}^* are both feasible and Pareto optimal among feasible arms. We then build alternative instances where arm $i \in \mathcal{S}^*$ is either made infeasible or sub-optimal among feasible arms under bandit λ^i .

Letting $\eta_i = \text{dist}(\mu_i, P^c) = \text{dist}(\mu_i, \partial P)$ (as $\mu_i \in P$, cf (6.13)) and since ∂P is a closed, using Lemma F2 of [Katz-Samuels & Scott 2019](#), there exists $z_i \in \partial P$ such that $\text{dist}(\mu_i, \partial P) = \|\mu_i - z_i\|_2 = \eta_i$. Then, note that, as $z_i \notin P^\circ$ (the interior of P), for all $\varepsilon > 0$, $\exists z_i^\varepsilon \in P^c$ such that $\|z_i^\varepsilon - z_i\| \leq \varepsilon$. Then, letting $\lambda_i^i = z_i^\varepsilon$ (for some ε fixed). It comes that i is infeasible under λ^i , so that $\mathbb{P}_{\lambda^i}(\mathcal{E}_i) \leq \delta$. Therefore, we have

$$\begin{aligned} \mathbb{E}_{\mu}[N_{\tau, i}] &\geq \frac{2}{\|\mu_i - z_i^\varepsilon\|^2} g(\delta) \\ &= \frac{2}{\|\mu_i - z_i + (z_i - z_i^\varepsilon)\|^2} g(\delta). \end{aligned}$$

Letting $\varepsilon \rightarrow 0$ yields

$$\mathbb{E}_{\mu}[N_{\tau, i}] \geq \frac{2}{\|\mu_i - z_i\|^2} g(\delta) = \frac{2}{\eta_i^2} g(\delta).$$

Alternative change of distribution will prove the dependency on the other quantities involved in the gaps. Assume the gap of i is attained for an arm $j \in \mathcal{S}^*$ with $\Delta_i = M(i, j)$. In this case, j is close to dominating i so that decreasing i will make it dominated by j , as result, i becomes either feasible but dominated or non-feasible, *i.e.*, $i \notin \mathcal{S}^*$ in both cases. We define in this case

$$\lambda_i^i := \mu_i - (1 + \varepsilon)(\mu_i - \mu_j)_+$$

where $(x)_+$ is defined component-wise for $x \in \mathbb{R}^d$. Then it can be easily checked that under bandit λ^i , arm i is now dominated by j and j is still feasible. Therefore, $\mathbb{P}_{\lambda^i}(\mathcal{E}_i) \leq \mathbb{P}_{\lambda^i}(i \in \widehat{O}) \leq \delta$. So that proceeding as before we prove that

$$\mathbb{E}_\mu[N_{\tau,i}] \geq \frac{2}{\|(\mu_i - \mu_j)_+\|^2} g(\delta). \quad (6.44)$$

Similarly, if we had $\Delta_i = M(j, i)$ for some $j \in \mathcal{S}^*$. Then, increasing i will make it dominate j . Thus, defining

$$\lambda_i^i := \mu_i + (1 + \varepsilon)(\mu_j - \mu_i)_+,$$

one can observe that arm j is now dominated by i under bandit λ^i . Note that j is still feasible in bandit λ^i (since its mean has not changed). So two things may occur. Either i is still feasible under λ^i in which case j is feasible and sub-optimal in bandit λ^i or i is now infeasible under bandit λ^i so that i is no longer an optimal arm. In both cases, the optimal set has changed, so that $\mathbb{P}_{\lambda^i}(\mathcal{E}_i) \leq \delta$. Reasoning as in the previous cases,

$$\mathbb{E}_\mu[N_{\tau,i}] \geq \frac{2}{\|(\mu_j - \mu_i)_+\|^2} g(\delta). \quad (6.45)$$

Therefore, by defining the quantity we have proved that

$$\mathbb{E}_\mu[N_{\tau,i}] \geq \frac{2}{\min(\tilde{\delta}_i^+, \eta_i)^2} g(\delta),$$

where

$$\tilde{\delta}_i^+ = \min_{j \in \mathcal{S}^* \setminus \{i\}} \min(\|(\mu_j - \mu_i)_+\|_2, \|(\mu_i - \mu_j)_+\|_2). \quad (6.46)$$

Note that this is a strictly positive quantity as $i \in \mathcal{S}^*$ and from the definition of non-dominance for all $j \in \mathcal{S}^*$.

Relation to the gaps in PSI. We recall the expressions of the PSI gaps (*w.r.t.* \mathcal{S}^*) as introduced in [Auer et al. 2016](#). For a sub-optimal arm $i \in S$,

$$\Delta_i := \Delta_{i^*} := \max_{j \in \mathcal{S}^*} m(i, j), \quad (6.47)$$

For an optimal arm $i \in \mathcal{S}^*$,

$$\Delta_i(S \cup \mathcal{S}^*) := \min(\delta_i^+, \delta_i^-(S))$$

where

$$\delta_i^+ := \min_{j \in \mathcal{S}^* \setminus \{i\}} \min(M(i, j), M(j, i)) \quad \text{and} \quad \delta_i^-(S) := \min_{j \in S} [(M(j, i))^+ + \Delta_j],$$

and $M(i, j) = \max_{c \in [d]} [\mu_i(c) - \mu_j(c)]$. Since \mathcal{S}^* is fixed, we ease notation and write $\Delta_i(S \cup \mathcal{S}^*) = \Delta_i(S)$ for all $i \in \mathcal{S}^*$.

As justified in Remark 18 of [Auer et al. 2016](#), $\delta_i^-(S)$, can be compensated by both $\delta_i^+(S)$ and Δ_{j^*} for some arms $j \in \text{SubOpt}$ (in our case). So we focus on matching the gaps in δ_i^+

for $i \in \mathcal{S}^*$. Moreover, from the definition of $M(i, j)$ for $i, j \in \mathcal{S}^*$ (and from the definition of Pareto dominance), it can be easily checked that for any $i \in \mathcal{S}^*$,

$$\delta_i^+ := \min_{j \in \mathcal{S}^* \setminus \{i\}} \min(\|(\mu_j - \mu_i)_+\|_\infty, \|(\mu_i - \mu_j)_+\|_\infty). \quad (6.48)$$

Intuitively, (in the case of independent objectives), (6.46) suggests to measure the distance between Pareto optimal arms in Euclidean whereas (6.48) measures in sup norm. As these quantities are computed over Pareto optimal arms, it can be checked that for $d = 2$, both terms are equal. In the worst case, the discrepancy between δ_i^+ and $\tilde{\delta}_i^+$ is at most \sqrt{d} . However, as we discuss later, there are classes of PSI/constrained PSI instances where these quantities are equal.

Similarly, one can check that for a sub-optimal arm, the discrepancy between Δ_i^* and $\|s_i\|^2$ (6.42) is at most $\min(\sqrt{|\Omega_i|}, \sqrt{d})$, where we recall that

$$\Omega_i := \{j \in \mathcal{S}^* : \mu_i \prec \mu_j\}.$$

Thus, in instances where Ω_i is singleton, that is each of SubOpt is dominated by a unique feasible optimal arm, $\|s_i\|_2$ and Δ_i^* are equal. To summarize, we have been justifying that under the following,

$$\forall i \in \mathcal{S}^*, \tilde{\delta}_i^+ = \delta_i^+ \quad \text{and} \quad \forall i \in \text{SubOpt}, \|s_i\|_2 = \Delta_i^* \quad (6.49)$$

combined with the results of sub-section 6.6.4 we would have

$$\mathbb{E}_\mu[N_{\tau,i}]/(2g(\delta)) \gtrsim \begin{cases} \frac{1}{\eta_i^2} & \text{if } i \in I \\ \frac{1}{(\Delta_i^*)^2} & \text{if } i \in S \\ \frac{1}{\min(\eta_i, \Delta_i(S))^2} & \text{if } i \in \mathcal{S}^*. \end{cases} \quad (6.50)$$

Then, as $\mathbb{E}_\mu[\tau] = \sum_{i=1}^K \mathbb{E}_\mu[N_{\tau,i}]$,

$$\mathbb{E}_\mu[\tau] \geq 2C(S, I)g(\delta).$$

where we recall that

$$C(S, I) = \sum_{i \in \mathcal{S}^*} \frac{1}{\min(\Delta_i(\mathcal{S}^* \cup S), \eta_i)^2} + \sum_{i \in S} \frac{1}{(\Delta_i(\mathcal{S}^* \cup S))^2} + \sum_{i \in I} \frac{1}{\eta_i^2}$$

Finally observing that (S, I) is an arbitrary correct response yields

$$\mathbb{E}_\mu[\tau] \gtrsim \min_{(S', I') \in \mathcal{M}} C(S', I')2g(\delta) \quad (6.51)$$

and the conclusion follows by observing that $2g(\delta) = 2 \text{kl}((1-\delta)/2, \delta) \geq (1-\delta) \log(1/\delta) - 2 \log(2)$ (as $\text{kl}(x, y) \geq x \log(1/y) - \log(2)$).

Conclusion. We now give examples of families of instances constructed to satisfy (6.49). This includes, e.g., instances where the means vectors are close in many objectives except a few (2 or 3). To give more intuition, we build some examples below. We let $\tilde{\mathcal{D}}$ be the family of multivariate normals with identity covariance and whose mean vector μ_1, \dots, μ_K is as below.

Let $\alpha > 0$ and let $\mu_0 \in P$, and $d \geq 1$ (arbitrary). Define $\mu_1 = \mu_0$ and for any $1 < i \leq K$, $\mu_k^1 = \mu_{k-1}^1 + \alpha$ and $\mu_k^2 = \mu_{k-1}^2 - \alpha$ and $\mu_k^c = \mu_0^c$ for all $c \notin \{1, 2\}$. Such sequences of vectors are non-dominated by each other and direct computation yields for any $i \in [K]$

$$\tilde{\delta}_i^+ = 2\alpha \quad \text{and} \quad \delta_i^+ = \alpha.$$

Thus, for $i \in \mathcal{S}^*$ the discrepancy between $\tilde{\delta}_i^+$ and δ_i^+ is 2, irrespective of the dimension. As for such instances, for any P , we either have $i \in \mathcal{S}^*$ or $\mu_i \notin P$, so, (6.51) is satisfied with an extra multiplicative constant 1/4.

Together, these imply that on this set of instances, we have

$$\mathbb{E}_\mu[\tau] \gtrsim \min_{(S', I') \in \mathcal{M}} C(S', I') \frac{g(\delta)}{2} \quad (6.52)$$

and $2g(\delta) = 2 \text{kl}((1 - \delta)/2, \delta) \geq (1 - \delta) \log(1/\delta) - 2 \log(2)$ which gives the claimed result,

$$\liminf_{\delta \rightarrow 0} \frac{\mathbb{E}_\mu[\tau]}{\log(1/\delta)} \geq \frac{C_{\mathcal{M}}^*(\mu)}{4}.$$

6.6.5 Complexity of Best response for cPSI

We focus on describing the procedure for computing the "best response"

$$\lambda_t := \operatorname{argmin}_{\lambda \in \text{Alt}(O_t)} \sum_i N_{t,i} \|\hat{\mu}_{t,i} - \lambda_i\|_2^2$$

, which is the most critical part in the implementation of Game-cPSI. The algorithm itself is described and analyzed in [Kone, Kaufmann, et al. 2025b](#).

In a recent work, [Crepon et al. 2024](#) proposed an algorithm to compute the best response λ_t in the case of unconstrained PSI (i.e., for $A = 0, b \in \mathbb{R}_+^d$). The computational cost of their algorithm is $O\left((K(p + d) + d^3 p) \binom{p + d - 1}{d - 1}\right)$ where p is the size of the Pareto set. We will show that their algorithm can be adapted to the case of constrained PSI. We prove the following result, which allows us to decompose Alt into two simpler sets $\text{Alt}^+(O)$ and $\text{Alt}^-(O)$.

Lemma 6.6.10. *It holds that*

$$\text{Alt}(O) = \text{Alt}^-(O) \cup \text{Alt}^+(O), \text{ where}$$

$$\begin{aligned} \text{Alt}^-(O) &:= \bigcup_{i,j \in O} (\{\lambda \in \mathcal{I}^K \mid A\lambda_i \not\leq b\} \cup \{\lambda \mid \lambda_i \leq \lambda_j\}) \\ \text{Alt}^+(O) &:= \bigcup_{i \notin O} \{\lambda \in \mathcal{I}^K \mid A\lambda_i \leq b \text{ and } \forall j \in O, \lambda_i \not\leq \lambda_j\}. \end{aligned}$$

Proof. We have

$$\begin{aligned}
 \lambda \in \text{Alt}(O) &\iff \mathcal{S}^*(\lambda) \neq O \\
 &\iff (a) : \exists i \in O \setminus \mathcal{S}^*(\lambda) \text{ or } (b) : \exists i \notin O \text{ such that } i \in \mathcal{S}^*(\lambda) \\
 &\iff (a) : \exists (i, j) \in O^2, i \neq j \text{ s.t. } \lambda_i \notin P \text{ or } \lambda_i \prec \lambda_j \text{ or} \\
 &\quad (b) : \exists i \notin O \text{ s.t. } \lambda_i \in P \text{ and } \forall j \in O, \lambda_i \not\prec \lambda_j
 \end{aligned}$$

To see the direct inclusion, let $\lambda \in \text{Alt}(O)$. If $O \not\subset F(\lambda)$ then (a) follows. Next, we assume $O \subset F(\lambda)$.

As $\lambda \in \text{Alt}(O)$, either $O \setminus \mathcal{S}^*(\lambda) \neq \emptyset$ or $\mathcal{S}^*(\lambda) \setminus O \neq \emptyset$. Assume there exists $i \in O$ such that $i \notin \mathcal{S}^*(\lambda)$.

If $\lambda_i \notin P$, then, as in the case above, the inclusion follows. Assume $\lambda_i \in P$, i.e., i is still a feasible arm in λ . In this case, there exists $j \in \mathcal{S}^*(\lambda)$ such that $\lambda_i \prec \lambda_j$, otherwise, we would have $i \in \mathcal{S}^*(\lambda)$.

If $j \in O$, then the inclusion (a) follows. If $j \notin O$, then, as $j \in \mathcal{S}^*(\lambda)$, we have: $\lambda_j \in P$ and $\forall k \in F(\lambda), \lambda_j \not\prec \lambda_k$. In particular, as $O \subset F(\lambda)$ it holds that for $k \in O$: $\lambda_j \not\prec \lambda_k$ and $j \in \mathcal{S}^*(\lambda) \subset F(\lambda)$, so (b) follows.

Now we assume there exists $i \in \mathcal{S}^*(\lambda)$ such that $i \notin O$, as we have $i \in \mathcal{S}^*(\lambda)$, it holds that $\lambda_i \in P$ and $\forall j \in F(\lambda), \lambda_i \not\prec \lambda_j$, in particular, as $O \subset F(\lambda)$, we have $\forall j \in O, \lambda_i \not\prec \lambda_j$, then (b) follows.

For the reverse inclusion, assume (a) holds. Then, it follows directly that we cannot have $O = \mathcal{S}^*(\lambda)$. Similarly, suppose (b) holds and $O = \mathcal{S}^*(\lambda)$. Then $\exists i \notin O = \mathcal{S}^*(\lambda)$ such that $i \in F(\lambda)$ and $\forall j \in O = \mathcal{S}^*(\lambda), \lambda_i \not\prec \lambda_j$, i.e., i is feasible in the instance λ , it does not belong to the optimal set $\mathcal{S}^*(\lambda)$ and it is not dominated by any arm of $\mathcal{S}^*(\lambda)$, which is impossible if $O = \mathcal{S}^*(\lambda)$. Therefore, when (a) or (b) holds, we have $\lambda \in \text{Alt}(O)$. \square

Lemma 6.6.11. *It holds that*

$$\inf_{\lambda \in \text{Alt}^-(O)} \sum_i \frac{1}{2} w_i \|\mu_i - \lambda_i\|_2^2 = \frac{1}{2} \min(\phi_1, \phi_2) \text{ where}$$

$$\phi_1 := \min_{i \in O} w_i \text{dist}(\mu_i, P^c)^2,$$

$$\phi_2 := \min_{i, j \in O^2, i \neq j} \frac{w_i w_j}{w_i + w_j} \sum_{c \leq d} (\mu_i^c - \mu_j^c)_+^2,$$

and a minimizer can be computed in time $O(p^2 d + pdm)$ where m is the number of constraints and $p = |O|$.

Proof. We have

$$\text{Alt}^-(O) = \left(\bigcup_{i \in O} \Gamma_i \right) \cup \left(\bigcup_{\substack{i, j \in O^2 \\ i \neq j}} \Lambda_{i, j} \right),$$

with

$$\begin{aligned}\Gamma_i &:= \{\lambda \in \mathcal{I}^K \mid \lambda_i \notin P\} \\ \Lambda_{i,j} &:= \{\lambda \in \mathcal{I}^K \mid \lambda_i \prec \lambda_j\}.\end{aligned}$$

Therefore, letting $D_w(\lambda) := \sum_i \frac{1}{2} w_i \|\mu_i - \lambda_i\|_2^2$ we have

$$\inf_{\lambda \in \text{Alt}^-(O)} D_w(\lambda) = \left(\min_{i \in O} \inf_{\lambda \in \Gamma_i} D_w(\lambda) \right) \wedge \min_{i,j \in O^2, i \neq j} \inf_{\lambda \in \Lambda_{i,j}} D_w(\lambda).$$

Next, we observe that

$$\inf_{\lambda \in \Gamma_i} D_w(\lambda) = \frac{1}{2} w_i \text{dist}(\mu_i, P^c)^2$$

and a minimizer of this quantity can be computed in $O(md)$. Indeed, since $\mu_i \in P$ (as $i \in O$),

$$\text{dist}(\mu_i, P^c) = \text{dist}(\mu_i, \partial P)$$

and the latter can be computed in $O(md)$ for a polyhedron with matrix of constraints $A \in \mathbb{R}^{m,d}$ (cf Lemma H.1 of [Katz-Samuels & Scott 2018](#)). Finally Lemma 2 of [Crepon et al. 2024](#) shows that the value of $\inf_{\lambda \in \Lambda_{i,j}} D_w(\lambda)$ and its minimizer can be computed in time $\mathcal{O}(p^2d)$. \square

Lemma 6.6.12. *There exists a polynomial-time algorithm (in $|O|$) that computes the value and a minimizer of*

$$\min_{\lambda \in \text{Alt}^+(O)} \sum_i \frac{1}{2} w_i \|\mu_i - \lambda_i\|_2^2.$$

Proof. In the instances of $\lambda \in \text{Alt}^+(O)$ a novel arm is added to the Pareto set of feasible arms. From its definition in Lemma 6.6.10, we have

$$\text{Alt}^+(O) = \bigcup_{i \in O^c} \text{Alt}_i^+(O)$$

where

$$\text{Alt}_i^+(O) := \{\lambda \in \mathcal{I}^K \mid \lambda_i \in P \text{ and } \forall j \in O, \lambda_i \not\prec \lambda_j\}.$$

Then, note that to guarantee that arm i is not dominated by j while minimizing the transportation cost, it is sufficient to move the arm (or both arms) along one objective. Therefore, to minimize the transportation cost $D(w, \lambda)$, it is sufficient to move arm i to a novel point $\lambda \in \mathbb{R}^d$ (satisfying the constraints) and move each arm $j \in O$ only along the coordinate with minimal cost to make λ not dominated by the $(\lambda_j)_{j \in O}$. This amounts to solving the following optimization problem:

$$\inf_{\lambda \in \text{Alt}_i^+(O)} D(w, \lambda) = \inf_{\substack{\lambda \in \mathbb{R}^d \\ A\lambda \leq b}} \frac{1}{2} w_i \|\mu_i - \lambda\|_2^2 + \sum_{j \in O} \frac{1}{2} w_j \min_{c \leq d} (\mu_j^c - \lambda^c)_+^2. \quad (6.53)$$

The rightmost problem is related to the optimization problem studied by [Crepon et al. 2024](#), except that now we have additional linear constraints $A\lambda \leq b$ due to the constrained PSI

setting. However, we will show that their algorithm can still be adapted to efficiently solve (6.53). To see this, we let ϕ be a mapping from O to $[d]$ and introduce $h_i^\phi : \mathbb{R}^d \rightarrow \mathbb{R}_+$, defined as

$$h_i^\phi(\lambda) := \frac{1}{2}w_i\|\mu_i - \lambda\|_2^2 + \sum_{j \in O} \frac{1}{2}w_j(\mu_j^{\phi(j)} - \lambda^{\phi(j)})_+^2.$$

Then note that the optimization problem in Equation 6.53 rewrites as

$$\min_{\phi \in [d]^p} \inf_{\substack{\lambda \in \mathbb{R}^d \\ A\lambda \leq b}} h_i^\phi(\lambda), \tag{6.54}$$

where the notation $[d]^p$ denotes the set of mapping from O to $[d]$ and $p = |O|$. However, as shown by Crepon et al. 2024, not all such maps are valid, and there is no need to optimize h_i^ϕ for an invalid map. The authors proposed an efficient graph-based algorithm to enumerate all the valid maps. They showed that the number of valid maps is bounded by $\binom{p+d-1}{d-1}$.

Now, it remains to minimize h_i^ϕ under the linear constraints for each valid map. Note that up to a reordering of the quantities, $(\mu_j^{\phi(j)})_{j \in O}$, h_i^ϕ is piecewise quadratic, and the problem can be decomposed into at most p^d convex quadratic problems with constraints $\tilde{A}\lambda \leq \tilde{b}$ where $\tilde{A} \in \mathbb{R}^{m+2d,d}$ and $\tilde{b} \in \mathbb{R}^{m+2d}$. It is known that solving a convex QP can be done in polynomial time (Ye & Tse 1989). Thus, when the number of arms is large and d is small, the overall complexity of problem 6.54 is $\mathcal{O}(p^{2d}s(m+d, d))$, where $s(m+2d, d)$ is the time complexity of a convex QP with $m+2d$ constraints in dimension d . \square

Chapter 7

Conclusion and Perspectives

Summary of Contributions

This dissertation investigated the design and analysis of algorithms for *pure exploration* in stochastic multi-objective decision problems, with a particular focus on identifying the *Pareto set* of optimal arms. Our goal was to develop theoretically grounded and computationally tractable strategies that achieve near-optimal performance in both the fixed-confidence and fixed-budget settings, while accommodating correlations and structural assumptions relevant to real-world applications such as early-stage clinical trials and multi-criteria recommender systems.

Unstructured Pareto Set Identification. The first part of the dissertation introduced the Pareto Set Identification (PSI) problem in its general, unstructured form. We analyzed PSI in the fixed-budget regime and proposed the Empirical Gap Elimination (EGE) algorithm, which estimates sub-optimality gaps associated with the difficulty of certifying each arm and progressively eliminates arms with larger gaps. We established an exponentially decreasing upper bound on its error probability. We proved that two specific variants, EGE-SH and EGE-SR, achieve near-optimal decay rates, matching the theoretical lower bounds up to logarithmic factors. Empirical evaluations confirmed these theoretical results, demonstrating that adaptive elimination markedly improves the efficiency of Pareto set identification.

In the fixed-confidence setting, we proposed Adaptive Pareto Exploration (APE), the first fully adaptive LUCB-type algorithm for PSI. We formalized several practical relaxations of the PSI objective, such as $(\varepsilon_1, \varepsilon_2)$ -PSI and ε_1 -PSI- k , and proved that they offer substantial reductions in sample complexity while fitting practical purposes.

Applications to clinical-trial-inspired datasets illustrated how adaptive exploration can efficiently identify promising treatment strategies across multiple efficacy criteria.

Linear PSI. Building on these foundations, we extended PSI to the linear setting, where arm means are linearly related through a shared but unknown parameter matrix θ and known feature vectors. We developed the G-optimal Empirical Gap Elimination (GEGE) algorithm, which integrates G-optimal design exploration with gap-based elimination. Analyzed under both the fixed-budget and fixed-confidence regimes, GEGE achieves tight

instance-dependent sample complexity depending only on the h smallest sub-optimality gaps, where h is the feature dimension, leading to substantial gains when the number of arms is large compared to h . Extensive simulations corroborated these theoretical insights, showing that exploiting linear structure dramatically improves sample efficiency compared to unstructured approaches.

Asymptotic Optimality and Correlated Objectives. The fifth chapter of the dissertation addressed *asymptotic optimality* in PSI, particularly under *correlated objectives*. We introduced PSIPS (Posterior Sampling for PSI), a posterior sampling-based algorithm that leverages posterior draws for both sampling and stopping. PSIPS removes the need to predefine the confidence level δ , avoids heavy computations required by Chernoff stopping rules, and can exploit correlated objectives to reduce the sample complexity. We proved that PSIPS attains *information-theoretic optimality*, with sample complexity matching the asymptotic lower bound. Empirically, PSIPS not only achieved comparable or superior accuracy to confidence-based methods but also exhibited large computational speedups compared to existing asymptotically optimal gradient-based approaches. In correlated settings, PSIPS may efficiently exploit the dependency structure to reduce sampling costs.

Constrained Pareto Set Identification. Finally, we extended PSI to account for feasibility constraints, where only non-dominated arms satisfying some linear constraints should be identified. We studied two variants of this problem: the non-explainable setting, where the goal is solely to identify the feasible Pareto set, and the explainable setting, where the learner must additionally provide reasons for excluding each non-selected arm (infeasible or Pareto-dominated). We proposed e-CAPE, a constrained extension of APE adapted to the explainable case, and proved its near-optimality for identifying the feasible Pareto-optimal set.

Perspectives and Open Problems

The results presented in this dissertation open several promising directions for future research in multi-objective pure exploration and beyond. While the algorithms developed here provide nearly optimal solutions in terms of sample complexity, several extensions remain to be explored to make these methods more adaptive, scalable, and applicable to richer decision-making environments.

Structured Correlation Models and Beyond Gaussianity. The current framework assumes a known covariance matrix. Extending PSI to the case of unknown or arm-dependent correlations, possibly learned online, would be of high practical value. Such extensions would enable the algorithms to adapt not only to arm means but also to the structure of dependencies among objectives, which could further reduce sample complexity in applications where correlation patterns differ across arms and are difficult to estimate upfront. Also, while convenient in theory, the Gaussian assumption is too restrictive, as many real-world outcomes (especially clinical or biological data) exhibit heavy tails, discrete distributions, or

log-normal distributions. Deriving tight guarantees for PSI under heavy-tailed distributions could be an interesting direction.

Scalability and Computational Efficiency. While PSIPS is more efficient than previous methods, sampling-based stopping still involves repeated posterior draws, which can become burdensome for large K or high-dimensional objectives. Developing low-variance approximations or variance reduction techniques will improve scalability and make these methods more practical for large-scale and high-dimensional applications.

Cost-Efficient Multi-Objective Pure Exploration. A natural continuation of this work is to explicitly account for the *cost of observations*. In many applications, observing all marginals of a multi-dimensional outcome is neither feasible nor desirable, as each marginal may incur a distinct computational, financial, or ethical cost. Developing cost-efficient PSI algorithms where the learner adaptively selects both arms and marginals to observe, aiming to minimize the total sampling cost while preserving δ -correctness, leads to a novel class of algorithms. This setting raises several open questions: How can one design *adaptive marginal selection strategies* that optimally balance exploration and cost under correlated objectives? Can one derive *instance-dependent cost lower bounds* analogous to the information-theoretic complexity introduced in this work? How does correlation reduce the overall cost of identifying the Pareto set?

A deeper understanding of such cost-aware algorithms could substantially impact applications such as adaptive clinical trials and personalized recommendation systems, where each measurement carries a distinct burden, whether financial or logistical.

Extension to Reinforcement Learning. Extending multi-objective pure exploration principles to reinforcement learning (RL) remains an exciting research frontier. Many practical tasks require reasoning over trajectories and state transitions, for example, long-term treatment administration to a patient. An important next step is to formulate *multi-objective best policy identification* in finite-horizon Markov decision processes, where one seeks to identify policies that are Pareto-optimal with respect to multiple criteria (*e.g.*, reward, risk, energy consumption). Recent advances in pure exploration for MDPs provide strong methodological foundations (A. J. Wagenmaker et al. 2022), yet no work has fully solved this problem. The challenge will lie in designing algorithms that combine efficient exploration with appropriate extensions of the sub-optimality gaps studied in this dissertation to some notion of “policy gaps.” The algorithms developed here could provide a starting point for such extensions to RL settings.

Towards Projection-Free Pure Exploration. While PSIPS was introduced with a focus on simplicity and computational tractability, it demonstrates that the posterior-resampling stopping rule is a viable alternative to the Chernoff stopping rule for PSI. Nevertheless, several challenges remain. Recent advances (Kaufmann & W.-M. Koolen 2021) have considerably simplified and improved the analysis of correctness for Chernoff stopping rules, making them increasingly popular, yet they remain hardly tractable except for problems with simple

or no structure. The posterior-resampling stopping rule offers a computationally tractable alternative to the GLR stopping rule. However, its correct calibration is problem-specific and relies on non-trivial anti-concentration results. Extending the GLR stopping rule to more general pure exploration problems and combining it with posterior-based sampling rules, in particular, top-two sampling strategies, will provide a general framework for computationally efficient and projection-free pure exploration applicable to a broader class of problems.

List of Datasets

CovBoost Dataset

Source and scope. We use immunogenicity outcomes from the Cov-Boost clinical trial of [Munro et al. 2021](#). For each vaccine strategy (“arm”) and for each of three immunogenicity indicators—anti-spike IgG, neutralizing antibody titre (NT₅₀), and wild-type cellular response—the paper reports: (i) the geometric mean response, (ii) cohort size, and (iii) a 95% confidence interval (CI) around the geometric mean, under a log-normal measurement model.

Modeling assumptions. Consistent with [Munro et al. 2021](#), we assume that each immunogenicity indicator is log-normally distributed. Consequently, the *log*-responses are approximately Gaussian. Empirically, the three indicators are weakly correlated; we therefore model them as *independent* and take a diagonal covariance matrix.

Recovering means on the log scale. Let x_1, \dots, x_n denote the raw (positive) biological measurements for a given arm and indicator, assumed log-normal. The reported geometric mean is

$$\bar{x}_{\text{geo}} = \left(\prod_{i=1}^n x_i \right)^{1/n}.$$

On the log scale, the empirical mean is

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n \log x_i = \log(\bar{x}_{\text{geo}}),$$

which we use as the (Gaussian) mean parameter for that arm/indicator.

Recovering variances on the log scale. The paper provides a 95% CI for the geometric mean. On the log scale, this induces a CI for \bar{x} of the form

$$\bar{x} \pm 1.96 \cdot \text{SE}_{\log}, \quad \text{with } \text{SE}_{\log} = \frac{\log U - \log L}{2 \cdot 1.96},$$

where $[L, U]$ is the reported CI in the original (positive) scale. From SE_{\log} and the cohort size n , we estimate the sample variance on the log scale via $s^2 \approx n \cdot \text{SE}_{\log}^2$. To reduce heteroscedasticity across arms, we pool per-indicator variances across all arms using the

usual unbiased pooled estimator

$$s_{\text{pooled}}^2 = \frac{\sum_{a=1}^K (n_a - 1) s_a^2}{\sum_{a=1}^K (n_a - 1)}.$$

The resulting pooled log-variances (one per indicator) are reported in Table 7.2.

Bandit abstraction. We thus obtain a $K = 20$ -armed, $d = 3$ -dimensional Gaussian bandit on the *log* scale: for arm i , $X_i \sim \mathcal{N}(\mu_i, \Sigma)$ with

$$\mu_i \in \mathbb{R}^3 \text{ given by the three log-means in Table 7.1, } \Sigma = \text{diag}(s_{\text{IgG}}^2, s_{\text{NT50}}^2, s_{\text{cell}}^2).$$

Sampling an arm simulates measuring the (log) response of a new patient on the three indicators.

Arm nomenclature and grouping. Arms correspond to *three-dose* strategies and are grouped by the first two doses (“Prime BNT/BNT” or “Prime ChAd/ChAd”). Within each group, the row label denotes the third (booster) dose. For example, in the “Prime BNT/BNT” group, the arm labeled “ChAd” means BNT as dose 1 and 2, and ChAd as the booster.

Table 7.1: Table of the empirical arithmetic mean of the log-transformed immune response for three immunogenicity indicators. Each acronym corresponds to a vaccine. There are two groups of arms corresponding to the first 2 doses: one with prime BNT/BNT (BNT as first and second dose) and the second with prime ChAd/ChAd (ChAd as first and second dose). Each row in the table gives the values of the 3 immune responses for an arm (*i.e.*, a combination of three doses).

Dose 1/Dose 2	Dose 3 (booster)	Immune response		
		Anti-spike IgG	NT ₅₀	cellular response
Prime BNT/BNT	ChAd	9.50	6.86	4.56
	NVX	9.29	6.64	4.04
	NVX Half	9.05	6.41	3.56
	BNT	10.21	7.49	4.43
	BNT Half	10.05	7.20	4.36
	VLA	8.34	5.67	3.51
	VLA Half	8.22	5.46	3.64
	Ad26	9.75	7.27	4.71
	m1273	10.43	7.61	4.72
	CVn	8.94	6.19	3.84
Prime ChAd/ChAd	ChAd	7.81	5.26	3.97
	NVX	8.85	6.59	4.73
	NVX Half	8.44	6.15	4.59
	BNT	9.93	7.39	4.75
	BNT Half	8.71	7.20	4.91
	VLA	7.51	5.31	3.96
	VLA Half	7.27	4.99	4.02
	Ad26	8.62	6.33	4.66
	m1273	10.35	7.77	5.00
	CVn	8.29	5.92	3.87

Table 7.2: Pooled variance of each group.

	Immune response		
	Anti-spike IgG	NT ₅₀	cellular response
Pooled sample variance	0.70	0.83	1.54

List of Notation

Global Notation

Acronyms and abbreviations

<i>i.i.d.</i>	independent and identically distributed
<i>i.e.</i>	<i>id est</i> (that is)
<i>s.t.</i>	such that
<i>w.r.t.</i>	with respect to
BAI	Best Arm Identification
PSI	Pareto Set Identification
FB	fixed-budget
FC	fixed-confidence
UCB	upper confidence bound
LCB	lower confidence bound
KL	Kullback–Leibler
GLR	generalized likelihood ratio
PAC	Probably Approximately Correct
APE	Adaptive Pareto Exploration
EGE	Empirical Gap Elimination
GEGE	G -optimal Empirical Gap Elimination
PSIPS	Pareto Set Identification with Posterior Sampling
e-cAPE	<i>explainable constrained</i> Adaptive Pareto Exploration

Multi-Armed Bandits

$K \in \mathbb{N}$	number of arms (candidate strategies/ treatments)
$d \in \mathbb{N}$	number of objectives (coordinates)
$[K] := \{1, \dots, K\}$	set of arms; indices $i, j, k, a \in [K]$
$[d] := \{1, \dots, d\}$	set of objectives; index $c \in [d]$
$\Sigma \in \mathbb{R}^{d \times d}$	covariance matrix of the d -dimensional outcomes (when applicable)
$\nu_i \in \mathcal{D}$	outcome distribution of arm i (supported on \mathbb{R}^d in the multi-objective case; scalar if $d = 1$)
$\nu := (\nu_1, \dots, \nu_K)$	bandit instance (collection of arm distributions)

$\mu_i := \mathbb{E}_{X \sim \nu_i}[X]$	mean outcome of arm i (vector in \mathbb{R}^d or scalar, μ_k^c its c -th coordinate)
$\mu := (\mu_1, \dots, \mu_K) \in \mathcal{I}^K$	mean parameter (with $\mathcal{I} \subseteq \mathbb{R}^d$ the set of feasible mean vectors)
$A_t \in [K]$	arm pulled at round t
$N_{t,i}$	number of pulls of arm i up to time $t - 1$
$Z_t \in \mathbb{R}^d$	observed outcome at round t ; $Z_t \mid A_t \sim \nu_{A_t}$
$\hat{\mu}_{t,i}$	empirical mean of arm i based on $N_{t,i}$ samples
$\hat{\mu}_t$	stacked vector $(\hat{\mu}_{t,i})_{i \in [K]}$
$\mathcal{H}_t := \sigma(A_1, Z_1, \dots, A_t, Z_t)$	history / filtration up to round t
σ^2	variance proxy for subgaussian arm distributions

Probability theory

$X \sim \nu$	random variable/vector X has distribution ν
$\mathbb{P}(E)$	probability of event E
$\mathbb{E}[X]$	expectation of X
\mathbb{P}_ν	probability under bandit instance ν
\mathbb{E}_ν	expectation under \mathbb{P}_ν
$\mathbb{1}(E)$	indicator of event E
$\text{KL}(P \parallel Q)$	Kullback–Leibler divergence between distributions P and Q
$\text{kl}(p, q)$	binary relative entropy (Bernoulli KL)
\triangle_K	$(K - 1)$ -dimensional probability simplex $\{w \in \mathbb{R}_+^K : \sum_{k=1}^K w_k = 1\}$

Pure exploration

$\mathcal{S}^* : \mathcal{I}^K \rightarrow 2^{[K]}$	query task; $\mathcal{S}^*(\mu)$ is the set of acceptable answers
$(A_t)_{t \geq 1}$	sampling rule (predictable w.r.t. $(\mathcal{H}_t)_{t \geq 1}$)
$(\hat{S}_t)_{t \geq 1}$	recommendation rule (answer at time t , may be an index, a subset of $[K]$, or another object)
τ	stopping time (with respect to (\mathcal{H}_t))
\hat{S}_τ	final output
$\delta \in (0, 1)$	fixed-confidence risk parameter
δ -PAC	$\mathbb{P}_\nu(\hat{S}_\tau \neq \mathcal{S}^*(\mu)) \leq \delta$ (single correct answer convention)

T

fixed-budget horizon (typically $\tau = T$)

Algebra and asymptotics

$(x)_+$

$\max\{x, 0\}$

$a \wedge b, a \vee b$

$\min\{a, b\}, \max\{a, b\}$

$\mathcal{O}(\cdot), \Omega(\cdot), o(\cdot)$

Landau notation

\log, \log_2

natural and base-2 logarithms

$\langle x, y \rangle, x^\top y$

inner product

$\|x\|_2, \|x\|$

Euclidean norm

$\|x\|_\infty$

ℓ_∞ norm: $\max_{c \in [d]} |x^c|$

I_d

$d \times d$ identity matrix

$\mathbf{1}/\mathbf{1}_d, \mathbf{0}/\mathbf{0}_d$

all-ones and all-zeros vectors in \mathbb{R}^d

$A, A^c, \partial A$

a set, its complement, and its boundary

Pareto dominance and Pareto set

$\mu_i \leq \mu_j$

componentwise order: $\mu_i^c \leq \mu_j^c$ for all $c \in [d]$

$\mu_i \preceq \mu_j$

Pareto dominance: $\mu_i \leq \mu_j$ and strict in at least one coordinate

$\mu_i \prec \mu_j$

strict Pareto dominance: $\mu_i^c < \mu_j^c$ for all $c \in [d]$

$i \preceq j, i \prec j$

shorthand for $\mu_i \preceq \mu_j$ and $\mu_i \prec \mu_j$

$\mathcal{S}^*(\mu)$

Pareto set: $\{i \in [K] : \nexists j \in [K] \text{ s.t. } \mu_i \prec \mu_j\}$

$\text{Alt}(\mathcal{S}^*(\mu)) := \{\lambda : \mathcal{S}^*(\lambda) \neq \mathcal{S}^*(\mu)\}$

alternative instances with a different Pareto set

$m(i, j) := \min_c [\mu_j^c - \mu_i^c]$

margin of non-domination of i by j

$M(i, j) := \max_c [\mu_i^c - \mu_j^c]$

margin by which i leads j

$\Delta_i^* := \max_{j \in \mathcal{S}^*} m(i, j)$

gap for sub-optimal i (distance to Pareto front)

$\delta_i^+ := \min_{j \notin \mathcal{S}^*} \min\{M(i, j), M(j, i)\}$

margin-to-being-dominated for optimal arm i

$\delta_i^- := \min_{j \in \mathcal{S}^* \setminus \{i\}} [M(j, i)_+ + \Delta_j^*]$

margin-to-sub-optimality for optimal arm i

$\Delta_i := \begin{cases} \Delta_i^*, & i \notin \mathcal{S}^*, \\ \min\{\delta_i^+, \delta_i^-\}, & i \in \mathcal{S}^*, \end{cases}$

combined suboptimality gap

$H(\nu) := \sum_{i=1}^K \Delta_i^{-2}$

instance complexity (governs fixed-confidence sample complexity)

$H_2(\nu) := \max_{k \in [K]} k \Delta_{(k)}^{-2}$

fixed-budget complexity ($\Delta_{(1)} \leq \dots \leq \Delta_{(K)}$ ordered gaps)

Fixed-Budget PSI

Fixed-budget objective

$$e_T(\nu) := \mathbb{P}_\nu(\widehat{S}_T \neq S^*(\nu))$$

$$e_{T,k}(\nu) := \mathbb{E}_\nu[L(\widehat{S}_T, k)]$$

$$\mathcal{L}(\widehat{S}, k) := \begin{cases} \mathbb{1}(\widehat{S} \not\subseteq S^*), & |\widehat{S}| = k, \\ \mathbb{1}(\widehat{S} \neq S^*), & \text{otherwise,} \end{cases}$$

misidentification probability (exact PSI)

k -relaxed expected loss (PSI- k objective)

loss function for PSI- k

EGE algorithm structure

R

total number of elimination rounds

λ_r

active set size threshold at the start of round r

t_r

number of samples per arm in round r

$\mathcal{A}_r \subseteq [K]$

active arm set at round r

$\mathcal{B}_r, \mathcal{D}_r$

arms classified optimal / sub-optimal through round r

$S_r := \text{EmpPareto}(\mathcal{A}_r)$

empirical Pareto set of active arms at round r

Empirical primitives and unified gap

$\widehat{\mu}_{r,i}$

empirical mean of arm i at the end of round r

$$m(i, j; r) := \min_{c \in [d]} [\widehat{\mu}_{r,j}^c - \widehat{\mu}_{r,i}^c]$$

empirical version of $m(i, j)$ at round r

$$M(i, j; r) := \max_{c \in [d]} [\widehat{\mu}_{r,i}^c - \widehat{\mu}_{r,j}^c]$$

empirical version of $M(i, j)$ at round r

$$\widehat{\Delta}_{i,r}^* := \max_{j \in \mathcal{A}_r} m(i, j; r)$$

gap for sub-optimal i

$$\widehat{\delta}_{i,r}^* := \min_{j \in \mathcal{A}_r \setminus \{i\}} [M(i, j; r) \wedge (M(j, i; r) + (\widehat{\Delta}_{j,r}^*)_+)]$$

gap for optimal arm k

$$\widehat{\Delta}_{i,r} := \max\{\widehat{\Delta}_{i,r}^*, \widehat{\delta}_{i,r}^*\}$$

unified gap (estimable without prior knowledge of S^*)

$\widetilde{T}^{R,t,\lambda}(\nu)$

generic EGE error exponent

PSI- k relaxation

$k \in [K]$

target number of Pareto-optimal arms to return

$$\omega(k) := \max_{i \in S^*} \Delta_i^{(k)}$$

k -th largest gap among optimal arms

$$\Delta_i^{(k)} := \begin{cases} \max\{\Delta_i, \omega(k)\}, & i \in S^*, \\ \Delta_i, & \text{otherwise,} \end{cases}$$

k -relaxed gap

$$H_2^{(k)}(\nu) := \max_{i \in [K]} i (\Delta_i^{(k)})^{-2}$$

k -relaxed fixed-budget complexity

Metrics

$\text{HV}(S)$

hypervolume indicator of arms in S (empirical evaluation metric)

Fixed-Confidence PSI

Relaxed Pareto identification

$$\varepsilon_1, \varepsilon_2 \geq 0$$

$$\mathcal{S}_{\varepsilon_1}^* := \{i \in [K] : \nexists j, \mu_i + \varepsilon_1 \mathbf{1} \prec \mu_j\}$$

$$k \in [K]$$

Empirical primitives

$$m(i, j; t), M(i, j; t)$$

$$M^+(i, j; t), M^-(i, j; t)$$

$$m^+(i, j; t), m^-(i, j; t)$$

$$\beta_{i,j}(t, \delta)$$

$$L_i(t) := \hat{\mu}_{t,i} - \beta_i(t, \delta), \quad U_i(t) := \hat{\mu}_{t,i} + \beta_i(t, \delta)$$

APE stopping statistics

$$S(t) := \{i : \nexists j, \hat{\mu}_{t,i} \prec \hat{\mu}_{t,j}\}$$

$$\text{OPT}_{\varepsilon_1}(t) := \{i : \forall j \neq i, M^-(i, j; t) + \varepsilon_1 > 0\}$$

$$g_i^{\varepsilon_2}(t) := \max_j [m^-(i, j; t) + \varepsilon_2 \mathbb{1}(j \in \text{OPT}_{\varepsilon_1}(t))]$$

$$h_i^{\varepsilon_1}(t) := \min_j [M^-(i, j; t) + \varepsilon_1]$$

$$Z_1^{\varepsilon_1}(t) := \min_i h_i^{\varepsilon_1}(t)$$

$$Z_2^{\varepsilon_1}(t) := \min_{i \in S(t)^c} \max\{g_i(t), h_i^{\varepsilon_1}(t)\}$$

Stopping times

$$\tau_{\varepsilon_1} := \inf\{t \geq K : Z_1^{\varepsilon_1}(t) > 0 \wedge Z_2^{\varepsilon_1}(t) > 0\}$$

$$\tau_{\varepsilon_1, \varepsilon_2}$$

$$\tau^k := \inf\{t \geq K : |\text{OPT}_{\varepsilon_1}(t)| \geq k\}$$

$$\tau_{\varepsilon_1}^k := \min(\tau_{\varepsilon_1}, \tau^k)$$

PSI- k gaps and complexity

$$\omega_i$$

$$\omega^k := \min_{i \in \mathcal{S}^{*,k}} \omega_i$$

$$\tilde{\Delta}_a^k := \max\{\omega^k, \varepsilon_1, \Delta_a\}$$

$$H(\nu) := \sum_a (\tilde{\Delta}_a^k)^{-2}$$

ε_1 : Pareto-margin relaxation; ε_2 : domination relaxation

ε_1 -Pareto set (additive relaxation of the Pareto set)

size parameter for k -PSI: return at most k Pareto-optimal arms

empirical m, M computed from $\hat{\mu}_t$

upper / lower confidence bounds on $M(i, j)$

upper / lower confidence bounds on $m(i, j)$

confidence radius for pair (i, j) (LIL-type bound)

lower / upper confidence bounds for arm i

empirical Pareto set at time t

ε_1 -empirically optimal set at time t

stopping score for sub-optimal arm i

stopping score for optimal arm i

aggregate optimality statistic

aggregate sub-optimality statistic

ε_1 -PSI stopping time

$(\varepsilon_1, \varepsilon_2)$ -PSI stopping time

k -PSI stopping time

ε_1 -PSI- k stopping time

gap of optimal arm i to the k -th best

k -th worst gap among optimal arms

modified gap for PSI- k

sample complexity for APE with relaxation

Linear PSI

Multi-output linear bandit model

h	feature dimension ($h \leq K$)
$x_a \in \mathbb{R}^h$	feature vector of arm a
$\theta \in \mathbb{R}^{h \times d}$	unknown parameter matrix (shared regression model)
$\mu_a = \theta^\top x_a \in \mathbb{R}^d$	linear mean model
$X := (x_1, \dots, x_K)^\top \in \mathbb{R}^{K \times h}$	feature matrix (all arms)
$y_t = \theta^\top x_{A_t} + \eta_t$	observation model (signal + noise) at round t
η_t	noise vector at round t (centered, subgaussian marginals)

Least-squares estimation

$X_n \in \mathbb{R}^{n \times h}, Y_n \in \mathbb{R}^{n \times d}$	arm-selection / observation matrices after n pulls
$V_n := X_n^\top X_n = \sum_i N_{n,i} x_i x_i^\top$	information / Gram matrix after n samples
V_n^\dagger	Moore–Penrose pseudoinverse of V_n
$\hat{\theta}_n := \arg \min_A \ X_n A - Y_n\ _F^2$	least-squares estimate of θ
$\hat{\mu}_{i,r} := \hat{\theta}_r^\top x_i$	estimated mean of arm i at round r (via linear model)
$\ x\ _A^2 := x^\top A x$	(Mahalanobis) squared norm with respect to matrix A

Optimal experimental design

$w \in \Delta_K : V(w) := \sum_{a=1}^K w_a x_a x_a^\top$	design matrix for allocation w
$\ x\ _{V(w)^{-1}}^2$	statistical prediction variance of an arm under allocation w
G -optimal design	$\arg \min_{w \in \Delta_K} \max_{i \in [K]} \ x_i\ _{V(w)^{-1}}^2$ (minimizes worst-case prediction variance)
$X_S := (x_i : i \in S)^\top$	feature sub-matrix restricted to arm set S
$h_S := \dim(\text{span}\{x_i : i \in S\})$	effective dimension of active arms S

Linear PSI complexity

$H_{1,\text{lin}}(\nu)$	fixed-confidence complexity: $\sum_{i \leq h} \Delta_{(i)}^{-2}$
$H_{2,\text{lin}}(\nu)$	fixed-budget complexity: $\max_{i \leq h} i \Delta_{(i)}^{-2}$

GEGE algorithm

$S_r := \{i \in \mathcal{A}_r : \nexists j \in \mathcal{A}_r, \hat{\mu}_{i,r} \prec \hat{\mu}_{j,r}\}$	empirical Pareto set of active arms at round r
--	--

Bayesian PSI

Gaussian learning model

$$\nu_i = \mathcal{N}(\mu_i, \Sigma)$$

$$\Sigma \in \mathbb{R}^{d \times d}$$

$$\|x\|_{\Sigma^{-1}} := \sqrt{x^\top \Sigma^{-1} x}$$

$$U_t \sim \text{Uniform}([0, 1])$$

$$\mathcal{H}_t := \sigma(U_1, A_1, Z_1, \dots, U_t)$$

Bayesian notation

$$\Pi_t := \bigotimes_k \mathcal{N}(\hat{\mu}_{t,k}, \Sigma / N_{t,k})$$

$$\hat{\Pi}_t := \bigotimes_k \mathcal{N}(\hat{\mu}_{t,k}, c(t-1, \delta) \Sigma / N_{t,k})$$

$$\theta_t^1, \theta_t^2, \dots \stackrel{\text{i.i.d.}}{\sim} \hat{\Pi}_t | \mathcal{H}_{t-1}$$

PS stopping rule

$$c(t, \delta)$$

$$M(t, \delta)$$

$$\tau_\delta^{\text{PS}} := \inf\{t : \forall m \in [M(t-1, \delta)], S^*(\theta_t^m) = \hat{S}_t\}$$

$$m_{t,\delta} := \inf\{m : \theta_t^m \in \text{Alt}(\hat{S}_t)\}$$

Game-based sampling rule

$$w_t := N_t / t$$

$$w_{\text{exp}}$$

$$\tilde{w}_t := (1 - \gamma_t)w_t + \gamma_t w_{\text{exp}}$$

Additional notation

$$d_\Sigma := d_\Sigma = \|1_d\|_{(\bar{\sigma}\Sigma)^{-1}}^2, \bar{\sigma} = \|\|\Sigma^{-1}\|\|$$

$$R_\Sigma(x)$$

$$W^{-1}(x) := -W_{-1}(-e^{-x})$$

Gaussian arm distribution (shared, known covariance)

known covariance matrix

Mahalanobis norm

exogenous randomness (for randomized sampling rules)

augmented filtration (includes exogenous noise)

product posterior (Gaussian, flat prior)

inflated posterior for posterior-sampling stopping

posterior samples at round t

inflation / coverage factor

maximum number of posterior draws required at round t

posterior-sampling stopping time

first posterior draw contradicting the current answer

empirical allocation at round t

forced-exploration allocation

mixed allocation (optimal + forced exploration)

covariance-dependent constant in tail bounds

multivariate Mills ratio under Σ

Lambert W (negative branch), used in stopping analysis

Constrained PSI

Constrained PSI / feasibility

$$P := \{x \in \mathbb{R}^d : Ax \leq b\}$$

m

$$\text{dist}(x, \mathcal{X}) := \inf_{y \in \mathcal{X}} \|x - y\|_2$$

$$\text{cl}(P), \text{int}(P), \partial P := \text{cl}(P) \setminus \text{int}(P)$$

$$\text{F}(\mu) := \{k : \mu_k \in P\}$$

$$\mathcal{S}_{\text{feas}}^*(\mu) := \text{Par}(\text{F}(\mu), \mu)$$

$$\text{SubOpt}(\mu) := \{i : \exists j \in \text{F}(\mu), \mu_i \prec \mu_j\}$$

$$\eta_i := \begin{cases} \text{dist}(\mu_i, P^c), & \mu_i \in P, \\ \text{dist}(\mu_i, \partial P), & \mu_i \notin P, \end{cases}$$

e-cPSI problem formulation

$$\text{Alt}(S) := \{\lambda : \mathcal{S}_{\text{feas}}^*(\lambda) \neq S\}$$

$$\text{Alt}(S, I) := \{\lambda : (S, I) \notin \mathcal{M}(\mathcal{P}, \lambda)\}$$

$$\mathcal{M}(\mathcal{P}, \mu) := \{(S, I) : S \subset \text{SubOpt}, \\ I \subset \mathbb{F}^c, S \cap I = \emptyset, S \cup I = \mathcal{S}_{\text{feas}}^c\}$$

$$T(\mu, S, I)$$

$$C^*(\nu) := \min_{(S, I)} C(\mu, S, I)$$

$$C(\mu, S, I)$$

Constrained gaps and complexity

$$\Delta_i^*(S) := \max_{j \in \text{Par}(S)} m(i, j)$$

$$\delta_i^*(S) := \min\{\delta_i^+(S), \delta_i^-(S)\}$$

$$\delta_i^+(S) := \min_{j \notin \text{Par}(S)} \min(M(i, j), M(j, i))$$

$$\delta_i^-(S) := \min_{j \in \text{Par}(S) \setminus \{i\}} (M(j, i)_+ + (\Delta_j^*(S))_+)$$

$$\Delta_i(S) := \max(\delta_i^*(S), \Delta_i^*(S))$$

feasibility polyhedron ($A \in \mathbb{R}^{m \times d}, b \in \mathbb{R}^m$)

number of linear constraints

distance from point x to set \mathcal{X}

closure, interior, and boundary of P

set of (truly) feasible arms

feasible Pareto set: $\{i : \mu_i \in P \text{ and no feasible } j \text{ dominates } i\}$

arms dominated by a feasible arm

feasibility margin of arm i

alternatives for constrained PSI (cPSI)

alternatives for explainable cPSI (e-cPSI)

set of correct partition answers for e-cPSI

e-cPSI characteristic time for answer (S, I)

finite-time e-cPSI sample complexity coefficient

sample complexity measure for partition (S, I)

gap of suboptimal arms in S

gap of arm i relative to feasible set S

margin-to-being-dominated (constrained to S)

margin-to-sub-optimality (constrained)

(Pareto) sub-optimality gap of arm i , restricted to S

List of Figures

1.1	Illustration of the multi-arm trial model	3
1.2	Empirical (arithmetic) average of the log-transformed immune response for three immunogenicity indicators reported by Munro et al. 2021. Each acronym corresponds to a vaccine. There are two groups of arms corresponding to the first 2 doses: one with prime BNT/BNT (BNT as first and second dose) and the second with prime ChAd/ChAd (ChAd as first and second dose). Combining the first two doses received and the third dose investigated in the trial, there were $K = 20$ arms.	4
1.3	Pareto set in a bi-objective setting. Green points are Pareto-optimal, while blue points are sub-optimal.	6
1.4	PSI gaps and distances	15
2.1	Application 1: COV-BOOST trial	45
2.2	Application 2: Sorting Networks dataset.	45
2.3	Arms on a convex Pareto set.	45
2.4	Each sub-optimal i is only dominated by i^*	45
2.5	$K = 200$ arms on the unit circle.	45
2.6	High dimension ($d = 10$) with 2 group of arms.	45
2.7	Synthetic Experiment 1: Group of arms on a convex Pareto set.	46
2.8	Hyper-volume fraction of the returned set on Exp.1 (convex Pareto set). . .	46
2.9	Estimated PSI- k loss for different values of k on Exp.1 (convex Pareto set). . .	46
3.1	COV-BOOST: APE- k vs. PSI-Unif-Elim at $\delta = 0.1$ (2000 runs).	77
3.2	Toy instance with $\mathcal{S}^* = \{0, 1, 2\}$. The difference on the x- and y-axes is 0.1 between arms 0 and 1, and 0.05 between arms 1 and 2. Right: average sample complexity over 2000 runs.	78
3.3	Average size of the returned cover (left) and corresponding sample complexity (right), both averaged over 2000 runs. The empirical error probability was negligible in all cases.	79

3.4	Frequency of each arm in the recommended set across 2000 runs for three representative values of ε_2	79
3.5	Illustrative instance with $\mathcal{S}^* = \{0, 2, 3\}$. Right: average sample complexity across repeated trials. APE avoids oversampling already-certified optimal arms.	80
4.1	Average misidentification rate vs. K (synthetic data).	101
4.2	Average sample complexity vs. K (synthetic data).	101
4.3	Average misidentification rate vs. T (NoC experiment).	102
4.4	Empirical sample complexity (NoC experiment).	102
5.1	Empirical stopping times on the covid19 experiment with $\delta = 0.01$ (left) and $\delta = 0.001$ (right).	135
5.2	Empirical impact of the correlation ρ on the sample complexity.	136
5.3	Theoretical impact of the correlation coefficient on the complexity T^*	136
5.4	Empirical stopping time on random Gaussian (left) and Bernoulli (right) instances.	136
5.5	Average runtime for the first T iterations in the COVID-19 experiment.	137
5.6	Average per-round number of rejection samples $m_t, m_{t,\delta}$	137
6.1	A constrained PSI instance. A two-stage approach scales as $1/\eta_1^2 + 2/\varepsilon^2 + \sum_{i=1}^2 1/\Delta_i(\{1, 2\})^2$, whereas e-CAPE scales as $1/\eta_1^2 + \sum_{i=1}^3 1/\Delta_i(\{1, 2, 3\})^2$ for $\varepsilon \ll 1$	166
6.2	Average response in the Secukinumab dose-finding trial. The feasible region corresponds to efficacy $\geq 40\%$ and acceptable toxicity.	171
6.3	Average response in the CovBoost vaccine trial. The feasible region is defined by an IgG response above 8.25 titer.	171
6.4	Empirical sample complexity averaged over 500 runs for the Secukinumab (left) and CovBoost (right) experiments.	171
6.5	Synthetic constrained PSI instances with (left to right): ordered polyhedron, cube, and simplex. The green region denotes feasibility.	172
6.6	Empirical distribution of the sample complexity on the difficult constrained 5-armed instance.	173
6.7	Empirical distribution of the sample complexity on the difficult constrained 10-armed instance.	173

List of Tables

2.1	Upper bounds on $e_T(\nu)$ for different algorithms (up to constants). APE-FB is an "oracle" algorithm introduced in Kone, Kaufmann, et al. 2024 by converting a fixed-confidence PSI algorithm. APE-FB requires the parameter $H(\nu)$ to run.	39
3.1	Stopping conditions and associated recommendation	69
3.2	Average sample complexity over 2000 random instances ($K = 5$). Average $ \mathcal{S}^* $ is (2.30, 4.06, 4.93) for $d = 2, 4, 8$	77
3.3	Average number of pulls under PSI-Unif-Elim divided by that under 0-APE- K for each arm. The largest gap arises from arm 2, which PSI-Unif-Elim oversamples to maintain dominance guarantees.	80
4.1	Sample complexity up to constant multiplicative terms of different algorithms for PSI in the fixed-confidence setting.	98
6.1	Empirical sample complexity averaged over 500 runs.	172
6.2	References to the exhaustive list of cases analyzed	176
7.1	Table of the empirical arithmetic mean of the log-transformed immune response for three immunogenicity indicators. Each acronym corresponds to a vaccine. There are two groups of arms corresponding to the first 2 doses: one with prime BNT/BNT (BNT as first and second dose) and the second with prime ChAd/ChAd (ChAd as first and second dose). Each row in the table gives the values of the 3 immune responses for an arm (<i>i.e.</i> , a combination of three doses).	203
7.2	Pooled variance of each group.	203

List of Algorithms

1.1	Generic Pure Exploration Bandit Algorithm	7
2.1	EGE: Empirical Gap Elimination	35
2.2	EGE-SR- k : Any k -sized subset of the Pareto Set	37
3.1	APE: Adaptive Pareto Exploration	71
4.1	OPTESTIMATOR: Least-squares estimation from G-optimal design	93
4.2	GEGE: G-optimal Empirical Gap Elimination (fixed-confidence)	95
4.3	GEGE: G-optimal Empirical Gap Elimination (fixed-budget)	96
5.1	PSIPS: Pareto Set Identification with Posterior Sampling	126
6.1	e-CAPE: Explainable Constrained Adaptive Pareto Exploration	164

Bibliography

- Mills, J. P. (1926). "Table of the Ratio: Area to Bounding Ordinate, for Any Portion of Normal Curve". *Biometrika*.
- Thompson, W. R. (1933). "On the likelihood that one unknown probability exceeds another in view of the evidence of two samples". *Biometrika*.
- Birnbaum, Z. W. (1942). "An Inequality for Mill's Ratio". *The Annals of Mathematical Statistics*.
- Robbins, H. E. (1952). "Some aspects of the sequential design of experiments". *Bulletin of the American Mathematical Society* 58.
- Kiefer, J. & J. Wolfowitz (1960). "The Equivalence of Two Extremum Problems". *Canadian Journal of Mathematics*.
- Kung, H. T., F. Luccio & F. P. Preparata (1975). "On Finding the Maxima of a Set of Vectors". *J. ACM*.
- Lai, T. L. (1976). "On Confidence Sequences". *The Annals of Statistics* 4.
- Lai, T. & H. Robbins (1985). "Asymptotically efficient adaptive allocation rules". *Advances in Applied Mathematics* 6.1.
- Ye, Y. & E. Tse (July 1989). "An extension of Karmarkar projective algorithm for convex quadratic programming". *Math. Program.* 44, pp. 157–179.
- Jennison, C. & B. W. Turnbull (1993). "Group sequential tests for bivariate response: interim analyses of clinical trials with both efficacy and safety endpoints." *Biometrics*.
- Miettinen, K. (1998). "Nonlinear multiobjective optimization". *International Series in Operations Research and Management Science*.
- Even-Dar, E., S. Mannor & Y. Mansour (2002). "PAC Bounds for Multi-armed Bandit and Markov Decision Processes". *Annual Conference Computational Learning Theory*.
- Auer, P. (2003). "Using Confidence Bounds for Exploitation-Exploration Trade-Offs".
- Pukelsheim, F. (2006). *Optimal Design of Experiments (Classics in Applied Mathematics)* (Classics in Applied Mathematics, 50). Society for Industrial and Applied Mathematics.
- Baricz, Á. (2008). "Mills' ratio: Monotonicity patterns and functional inequalities". *Journal of Mathematical Analysis and Applications*.
- Bubeck, S., R. Munos & G. Stoltz (2009). "Pure exploration in multi-armed bandits problems". *Proceedings of the 20th International Conference on Algorithmic Learning Theory*. ALT'09. Springer-Verlag.
- Lu, D. & W. V. Li (2009). "A note on multivariate Gaussian estimates". *Journal of Mathematical Analysis and Applications*.
- Audibert, J.-Y. & S. Bubeck (2010). "Best Arm Identification in Multi-Armed Bandits". *COLT - 23th Conference on Learning Theory - 2010*.

- Almer, O., N. Topham & B. Franke (2011a). “A learning-based approach to the automated design of MPSoC networks”. *Proceedings of the 24th International Conference on Architecture of Computing Systems*. Springer-Verlag.
- (2011b). “A learning-based approach to the automated design of MPSoC networks”. *Proceedings of the 24th International Conference on Architecture of Computing Systems*. ARCS’11. Springer-Verlag.
- Garivier, A. & O. Cappé (2011). “The KL-UCB Algorithm for Bounded Stochastic Bandits and Beyond”. *Proceedings of the 24th Annual Conference on Learning Theory*. Vol. 19. Proceedings of Machine Learning Research. PMLR.
- Agrawal, S. & N. Goyal (2012). “Analysis of Thompson Sampling for the Multi-armed Bandit Problem”. *Proceedings of the 25th Annual Conference on Learning Theory*. Vol. 23. Proceedings of Machine Learning Research. PMLR.
- Gabillon, V., M. Ghavamzadeh & A. Lazaric (2012). “Best Arm Identification: A Unified Approach to Fixed Budget and Fixed Confidence”. *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- Kalyanakrishnan, S., A. Tewari, P. Auer & P. Stone (2012). “PAC Subset Selection in Stochastic Multi-Armed Bandits”. *Proceedings of the 29th International Conference on Machine Learning*. Omnipress.
- Kaufmann, E., N. Korda & R. Munos (2012). “Thompson sampling: An asymptotically optimal finite-time analysis”. *International conference on algorithmic learning theory*.
- Plotkin, S. A. & P. B. Gilbert (2012). “Nomenclature for Immune Correlates of Protection After Vaccination”. *Clinical Infectious Diseases* 54.11.
- Zuluaga, M., P. Milder & M. Püschel (2012). “Computer generation of streaming sorting networks”. *DAC Design Automation Conference 2012*.
- Bubeck, S., T. Wang & N. Viswanathan (2013). “Multiple Identifications in Multi-Armed Bandits”. *Proceedings of the 30th International Conference on Machine Learning*. PMLR.
- Drugan, M.-M. & A. Nowe (2013). “Designing multi-objective multi-armed bandits algorithms: A study”. *The 2013 International Joint Conference on Neural Networks (IJCNN)*.
- Karnin, Z., T. Koren & O. Somekh (2013). “Almost Optimal Exploration in Multi-Armed Bandits”. *Proceedings of the 30th International Conference on International Conference on Machine Learning*. JMLR.
- Kaufmann, E. & S. Kalyanakrishnan (2013). “Information Complexity in Bandit Subset Selection”. *Conference On Learning Theory*. JMLR: Workshop and Conference Proceedings.
- Mark C, G., D. Patrick, R. Hanno B, S. Jerzy, D. Eva, M. Vadim, A. Jacob A, L. Sang-Heon, C. Christine E, K. Herbert, I. Takashi, H. Sophie & M. Shephard (2013). “Efficacy and safety of secukinumab in patients with rheumatoid arthritis: a phase II, dose-finding, double-blind, randomised, placebo controlled study”. *Annals of the rheumatic diseases*.
- Zuluaga, M., G. Sergent, A. Krause & M. Püschel (2013). “Active Learning for Multi-Objective Optimization”. *Proceedings of the 30th International Conference on Machine Learning*. PMLR.

- Chen, S., T. Lin, I. King, M. R. Lyu & W. Chen (2014). “Combinatorial Pure Exploration of Multi-Armed Bandits”. *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- De Rooij, S., T. Van Erven, P. D. Grünwald & W. M. Koolen (2014). “Follow the leader if you can, hedge if you must”. *J. Mach. Learn. Res.*
- Jamieson, K., M. Malloy, R. Nowak & S. Bubeck (2014). “lil’ UCB: An Optimal Exploration Algorithm for Multi-Armed Bandits”. *Proceedings of The 27th Conference on Learning Theory*. PMLR.
- Jamieson, K. & R. Nowak (2014). “Best-arm identification algorithms for multi-armed bandits in the fixed confidence setting”. *2014 48th Annual Conference on Information Sciences and Systems (CISS)*.
- Soare, M., A. Lazaric & R. Munos (2014). “Best-arm identification in linear bandits”. *Advances in Neural Information Processing Systems*.
- Auer, P., C.-K. Chiang, R. Ortner & M.-M. Drugan (2016). “Pareto Front Identification from Stochastic Bandit Feedback”. *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*. PMLR.
- Carpentier, A. & A. Locatelli (2016). “Tight (Lower) Bounds for the Fixed Budget Best Arm Identification Bandit Problem”. *29th Annual Conference on Learning Theory*. PMLR.
- Degenne, R. & V. Perchet (2016). “Combinatorial semi-bandit with known covariance”. *Advances in Neural Information Processing Systems*.
- Garivier, A. & E. Kaufmann (2016). “Optimal Best Arm Identification with Fixed Confidence”. *29th Annual Conference on Learning Theory*. PMLR.
- Kaufmann, E., O. Cappé & A. Garivier (2016). “On the Complexity of Best-Arm Identification in Multi-Armed Bandit Models”. *Journal of Machine Learning Research*.
- Locatelli, A., M. Gutzeit & A. Carpentier (2016). “An optimal algorithm for the Thresholding Bandit Problem”. *Proceedings of The 33rd International Conference on Machine Learning*. PMLR.
- Russo, D. (2016). “Simple Bayesian Algorithms for Best Arm Identification”. *29th Annual Conference on Learning Theory*. Vol. 49. Proceedings of Machine Learning Research. PMLR.
- Zuluaga, M., A. Krause & M. Püschel (2016). “e-PAL: An Active Learning Approach to the Multi-Objective Optimization Problem”. *Journal of Machine Learning Research*.
- Allen-Zhu, Z., Y. Li, A. Singh & Y. Wang (2017). “Near-optimal design of experiments via regret minimization”. *Proceedings of the 34th International Conference on Machine Learning*. JMLR.org.
- Simchowitz, M., K. Jamieson & B. Recht (2017). “The Simulator: Understanding Adaptive Sampling in the Moderate-Confidence Regime”. *Proceedings of the 2017 Conference on Learning Theory*. PMLR.
- Howard, S. R., A. Ramdas, J. D. McAuliffe & J. S. Sekhon (2018). “Time-uniform, nonparametric, nonasymptotic confidence sequences”. *The Annals of Statistics*.
- Katz-Samuels, J. & C. Scott (2018). “Feasible Arm Identification”. *Proceedings of the 35th International Conference on Machine Learning*. PMLR.

- Tao, C., S. Blanco & Y. Zhou (2018). “Best Arm Identification in Linear Bandits with Linear Dimension Dependency”. *Proceedings of the 35th International Conference on Machine Learning*. PMLR.
- Xu, L., J. Honda & M. Sugiyama (2018). “A fully adaptive algorithm for pure exploration in linear bandits”. *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*. PMLR.
- Afshari, H., W. Hare & S. Tesfamariam (2019). “Constrained multi-objective optimization algorithms: Review and comparison with application in reinforced concrete structures”. *Applied Soft Computing*.
- Amani, S., M. Alizadeh & C. Thrampoulidis (2019). “Linear Stochastic Bandits Under Safety Constraints”. *Advances in Neural Information Processing Systems*.
- Degenne, R. & W. Koolen (2019). “Pure Exploration with Multiple Correct Answers”. *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- Degenne, R. (2019). “Impact of structure on the design and analysis of bandit algorithms”. PhD thesis.
- Degenne, R., W. M. Koolen & P. Ménard (2019). “Non-asymptotic pure exploration by solving games”. *Advances in Neural Information Processing Systems*.
- Fiez, T., L. Jain, K. G. Jamieson & L. Ratliff (2019). “Sequential Experimental Design for Transductive Linear Bandits”. *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- Garivier, A., P. Ménard & G. Stoltz (2019). “Explore first, exploit next: The true shape of regret in bandit problems”. *Mathematics of Operations Research*.
- Katz-Samuels, J. & C. Scott (2019). “Top Feasible Arm Identification”. *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*. PMLR.
- Lu, S., G. Wang, Y. Hu & L. Zhang (2019). “Multi-Objective Generalized Linear Bandits”. *Proceedings of the 28th International Joint Conference on Artificial Intelligence*. AAAI Press.
- Ménard, P. (2019). *Gradient Ascent for Active Exploration in Bandit Problems*.
- Roy Chaudhuri, A. & S. Kalyan Krishnan (2019). “PAC Identification of Many Good Arms in Stochastic Multi-Armed Bandits”. *Proceedings of the 36th International Conference on Machine Learning*. PMLR.
- Daulton, S., M. Balandat & E. Bakshy (2020). “Differentiable Expected Hypervolume Improvement for Parallel Multi-Objective Bayesian Optimization”. *Proceedings of the 34th International Conference on Neural Information Processing Systems*. Curran Associates Inc.
- Degenne, R., P. Ménard, X. Shang & M. Valko (2020). “Gamification of pure exploration for linear bandits”. *International Conference on Machine Learning*. PMLR.
- Jedra, Y. & A. Proutiere (2020a). “Optimal Best-arm Identification in Linear Bandits”. *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- (2020b). “Optimal Best-arm Identification in Linear Bandits”. *Advances in Neural Information Processing Systems*. Curran Associates, Inc.

- Katz-Samuels, J., L. Jain, z. karnin zohar & K. G. Jamieson (2020). “An Empirical Process Approach to the Union Bound: Practical Algorithms for Combinatorial and Linear Bandits”. *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- Lattimore, T. & C. Szepesvari (2020). *Bandit Algorithms*. Cambridge University Press.
- Mehrotra, R., N. Xue & M. Lalmas (2020). “Bandit Based Optimization of Multiple Objectives on a Music Streaming Platform”. *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. Association for Computing Machinery.
- Russo, D., B. V. Roy, A. Kazerouni, I. Osband & Z. Wen (2020). *A Tutorial on Thompson Sampling*.
- Shang, X., R. de Heide, P. Menard, E. Kaufmann & M. Valko (2020). “Fixed-confidence guarantees for Bayesian best-arm identification”. *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*. PMLR.
- Alieva, A., A. Cutkosky & A. Das (2021). “Robust Pure Exploration in Linear Bandits with Limited Budget”. *Proceedings of the 38th International Conference on Machine Learning*. PMLR.
- Camilleri, R., K. Jamieson & J. Katz-Samuels (2021). “High-dimensional Experimental Design and Kernel Bandits”. *Proceedings of the 38th International Conference on Machine Learning*. PMLR.
- Garivier, A. & E. Kaufmann (2021). *Non-Asymptotic Sequential Tests for Overlapping Hypotheses and application to near optimal arm identification in bandit models*.
- Kaufmann, E. & W.-M. Koolen (2021). “Mixture Martingales Revisited with Applications to Sequential Tests and Confidence Intervals”. *Journal of Machine Learning Research*.
- Kumar, A., G. Wu, M. Z. Ali, Q. Luo, R. Mallipeddi, P. N. Suganthan & S. Das (2021). “A Benchmark-Suite of real-World constrained multi-objective optimization problems and some baseline results”. *Swarm and Evolutionary Computation*.
- Munro, A.-P.-S., L. Janani, V. Cornelius & et al. (2021). “Safety and immunogenicity of seven COVID-19 vaccines as a third dose (booster) following two doses of ChAdOx1 nCov-19 or BNT162b2 in the UK (COV-BOOST): a blinded, multicentre, randomised, controlled, phase 2 trial”. *The Lancet*.
- Wang, P.-A., R.-C. Tzeng & A. Proutiere (2021). “Fast Pure Exploration via Frank-Wolfe”. *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- Azizi, M., B. Kveton & M. Ghavamzadeh (2022). “Fixed-Budget Best-Arm Identification in Structured Bandits”. *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence*. International Joint Conferences on Artificial Intelligence Organization.
- Camilleri, R., A. Wagenmaker, J. Morgenstern, L. Jain & K. Jamieson (2022). “Active learning with safety constraints”. *Proceedings of the 36th International Conference on Neural Information Processing Systems*. NIPS '22. Curran Associates Inc.
- Faizal, F. Z. & J. Nair (2022). *Constrained Pure Exploration Multi-Armed Bandits with a Fixed Budget*.
- Jourdan, M., R. Degenne, D. Baudry, R. de Heide & E. Kaufmann (2022). “Top Two Algorithms Revisited”. *Advances in Neural Information Processing Systems*. Curran Associates, Inc.

- Karpov, N. & Q. Zhang (2022). “Collaborative Best Arm Identification with Limited Communication on Non-IID Data”. *arXiv preprint, eprint:2207.08015*.
- Ramdas, A., J. Ruf, M. Larsson & W. Koolen (2022). *Admissible anytime-valid sequential inference must rely on nonnegative martingales*.
- Wagenmaker, A. J., M. Simchowitz & K. Jamieson (2022). “Beyond No Regret: Instance-Dependent PAC Reinforcement Learning”. *Proceedings of Thirty Fifth Conference on Learning Theory*. Vol. 178. Proceedings of Machine Learning Research. PMLR.
- Wang, S. & J. Zhu (2022). “Thompson sampling for (combinatorial) pure exploration”. *International Conference on Machine Learning*. PMLR.
- Yang, J. & V. Tan (2022). “Minimax Optimal Fixed-Budget Best Arm Identification in Linear Bandits”. *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- Ararat, C. & C. Tekin (2023). “Vector Optimization with Stochastic Bandit Feedback”. *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*. PMLR.
- Degenne, R. (2023). “On the Existence of a Complexity in Fixed Budget Bandit Identification”. *Proceedings of Thirty Sixth Conference on Learning Theory*. PMLR.
- Emmenegger, N., M. Mutny & A. Krause (2023). “Likelihood Ratio Confidence Sets for Sequential Decision Making”. *Thirty-seventh Conference on Neural Information Processing Systems*.
- Jourdan, M., R. Degenne & E. Kaufmann (2023). “An ε -Best-Arm Identification Algorithm for Fixed-Confidence and Beyond”. *Thirty-Seventh Conference on Neural Information Processing Systems*.
- Kim, W., G. Iyengar & A. Zeevi (2023). *Pareto Front Identification with Regret Minimization*.
- Kone, C., E. Kaufmann & L. Richert (2023). “Adaptive Algorithms for Relaxed Pareto Set Identification”. *Thirty-seventh Conference on Neural Information Processing Systems*.
- You, W., C. Qin, Z. Wang & S. Yang (2023). “Information-Directed Selection for Top-Two Algorithms”. *Proceedings of Thirty Sixth Conference on Learning Theory*. Vol. 195. Proceedings of Machine Learning Research. PMLR.
- Crepon, É., A. Garivier & W. M Koolen (2024). “Sequential learning of the Pareto front for multi-objective bandits”. *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*. PMLR.
- Kanarios, K., Q. Zhang & L. Ying (2024). “Cost Aware Best Arm Identification”. *Reinforcement Learning Conference*.
- Karagözlü, E. M., Y. C. Yıldırım, C. Ararat & C. Tekin (2024). “Learning the Pareto Set Under Incomplete Preferences: Pure Exploration in Vector Bandits”. *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*.
- Kone, C., E. Kaufmann & L. Richert (2024). “Bandit Pareto Set Identification: the Fixed Budget Setting”. *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*. PMLR.
- Li, D., F. Zhang, C. Liu & Y. Chen (2024). *Constrained Multi-objective Bayesian Optimization through Optimistic Constraints Estimation*.

- Li, Z., K. Jamieson & L. Jain (2024). “Optimal Exploration is no harder than Thompson Sampling”. *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*. Proceedings of Machine Learning Research. PMLR.
- Poiani, R., R. Degenne, E. Kaufmann, A. M. Metelli & M. Restelli (2024). “Optimal Multi-Fidelity Best-Arm Identification”. *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- Shukla, A. & D. Basu (2024). “Preference-based Pure Exploration”. *Advances in Neural Information Processing Systems*. Vol. 37. Curran Associates, Inc.
- Das, U., A. Shukla & D. Basu (2025). *FraPPE: Fast and Efficient Preference-based Pure Exploration*.
- Kejriwal, K., N. Karamchandani & J. Nair (2025). “On the Asymptotic Optimality of Confidence Interval Based Algorithms for Fixed Confidence MABs”. *Proceedings of the AAAI Conference on Artificial Intelligence* 39.17.
- Kirschner, J., A. Krause, M. Meziu & M. Mutny (2025). *Confidence Estimation via Sequential Likelihood Mixing*.
- Kone, C., M. Jourdan & E. Kaufmann (2025). “Pareto Set Identification With Posterior Sampling”. *Proceedings of The 28th International Conference on Artificial Intelligence and Statistics*. Vol. 258. Proceedings of Machine Learning Research. PMLR.
- Kone, C., E. Kaufmann & L. Richert (2025a). “Bandit Pareto Set Identification in a Multi-Output Linear Model”. *The 28th International Conference on Artificial Intelligence and Statistics*.
- (2025b). “Constrained Pareto Set Identification with Bandit Feedback”. *Forty-second International Conference on Machine Learning*.