



HAL
open science

The price of unfairness in linear bandits with biased feedback

Solenne Gaucher, Alexandra Carpentier, Christophe Giraud

► **To cite this version:**

Solenne Gaucher, Alexandra Carpentier, Christophe Giraud. The price of unfairness in linear bandits with biased feedback. 2022. hal-03611628v2

HAL Id: hal-03611628

<https://hal.science/hal-03611628v2>

Preprint submitted on 1 Jun 2022 (v2), last revised 29 Nov 2022 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

The price of unfairness in linear bandits with biased feedback

Solenne Gaucher ^{*1}, Alexandra Carpentier², and Christophe Giraud¹

¹Laboratoire de Mathématiques d’Orsay, Université Paris-Saclay

²University of Potsdam

June 1, 2022

Abstract

In this paper, we study the problem of fair sequential decision making with biased linear bandit feedback. At each round, a player selects an action described by a covariate and by a sensitive attribute. The perceived reward is a linear combination of the covariates of the chosen action, but the player only observes a biased evaluation of this reward, depending on the sensitive attribute. To characterize the difficulty of this problem, we design a phased elimination algorithm that corrects the unfair evaluations, and establish upper bounds on its regret. We show that the worst-case regret is smaller than $\mathcal{O}(\kappa_*^{1/3} \log(T)^{1/3} T^{2/3})$, where κ_* is an explicit geometrical constant characterizing the difficulty of bias estimation. We prove lower bounds on the worst-case regret for some sets of actions showing that this rate is tight up to a possible sub-logarithmic factor. We also derive gap-dependent upper bounds on the regret, and matching lower bounds for some problem instance. Interestingly, these results reveal a transition between a regime where the problem is as difficult as its unbiased counterpart, and a regime where it can be much harder.

1 Introduction

Artificial intelligence is increasingly used in a wide range of decision making scenarii with higher and higher stakes, with application in online advertisement [27], credit [3], health care [10], education [24] and job interviews [30], in the hope of improving accuracy and efficiency. Recent works have shown that the decisions made by algorithms can be dangerously biased against certain categories of people, and have endeavored to mitigate this behavior [19, 12, 6, 23]. Studies have underlined that the main cause of algorithmic unfairness is the presence of bias in the training set [23], which led to the development of methods aiming to guarantee the fairness of the algorithms. This paper, in lines with these works, addresses the problem of online decision making under biased feedback.

Linear bandits have become a very popular tool in online decision making problems, when side information on the actions is available in the form of covariates. In the present paper, we consider a variant of this problem, where the agent only has access to an unfair assessment of the action taken, that is systematically biased against a group of actions. For example, examiners may be prejudiced against people from a minority group, and give them lower grades; similarly, algorithms trained on biased data may produce unfair assessments of the credit risk of individuals belonging to a minority group. Note that not correcting biased evaluation can have adverse effects for all parties: on the one hand, actions disadvantaged by the evaluation mechanism will be unfairly discriminated against; on the other hand, the agent may spend his budget on an unfairly advantaged action that is actually sub-optimal. The problem of sequential decision making under biased feedback can be formalized as follows.

Biased linear bandit problem A player is presented with a set of k distinct actions characterized by covariates $x \in \mathcal{X} \subset \mathbb{R}^d$, and by known sensitive attributes $z_x \in \{-1, 1\}$ indicating the group of the action. At

*solenne.gaucher@math.u-psud.fr

each round $t \leq T$, the player chooses the action x_t and receives an unobserved reward $x_t^\top \gamma^*$, where $\gamma^* \in \mathbb{R}^d$ is the regression parameter specifying the true value of the action. The regret of the player is given by

$$R_T = \mathbb{E} \left[\sum_{t \leq T} (x^* - x_t)^\top \gamma^* \right], \quad \text{where } x^* \in \operatorname{argmax}_{x \in \mathcal{X}} x^\top \gamma^*. \quad (1)$$

By contrast to the classical linear bandit, the player does not observe a noisy version of the unbiased reward $x_t^\top \gamma^*$. Instead, she observes an unfair evaluation y_t of the value of the action $x_t^\top \gamma^*$, given by the following biased linear model:

$$y_t = x_t^\top \gamma^* + z_{x_t} \omega^* + \xi_t$$

where $\xi_t \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$ is a noise term. The evaluation are systematically biased against a certain group: this unequal treatment of the groups is captured by the bias parameter $\omega^* \in \mathbb{R}$.

Preliminary discussion The biased linear bandit is a variant of the linear bandit. By contrast, in the classical linear bandit model, the agent observes a noisy version of the reward. Obviously, applying directly an algorithm designed for linear bandit to biased linear bandits without correcting the evaluations would lead to a linear regret if the evaluation mechanism is prejudiced against the group of the best action in terms of reward, and if the best action in terms of feedback belongs to the advantaged group. To avoid this pitfall, one must estimate the bias in order to correct the evaluations. This implies a change in the exploration-exploitation trade-off, as exploration becomes more expensive. Indeed, in classical bandit problems, one can compare the rewards of two actions by repeatedly sampling them - or, to put it differently, one can find the best action by sampling only those actions that seem optimal. This does not hold in the biased linear bandit: if, at some point, the set of potentially optimal actions contains representatives from both groups, and does not span \mathbb{R}^d , one is forced to sample sub-optimal actions to estimate the bias and improve the estimation of the unbiased rewards. For this reason, classical algorithm for linear bandit that only sample actions considered as potentially optimal, such as OFUL [1] or Phase Elimination [21], can suffer linear regret. This underlines the necessity to ensure sufficient estimation of the bias parameter, even when it implies sampling sub-optimal actions.

1.1 Related work

Fairness in bandit problems has mostly been studied from the perspective of fair budget allocation between actions. This problem is motivated by the fact that classical bandit algorithms select sub-optimal actions only a vanishing fraction of the time, which may be undesirable in many situations. To mitigate this problem and guarantee diversity in the actions selected, some papers [4, 25, 8, 13, 34] have proposed new algorithms ensuring fairness of the selection frequency of each action. The framework studied in this paper is different: we consider here that the mechanism for observing the rewards is unfair, and we aim at correcting it in order to maximize a (fair) true cumulative reward.

The biased linear model has been studied in the batch setting in [7], where the authors investigate the optimal trade-off between minimax risk and Demographic Parity. Detection of systematic bias, interpreted as a treatment effect, has been investigated in a batch setting in [15]. In [2], the authors consider a similar model, with unobserved sensitive attribute z and known bias parameter ω^* , under additional assumption that the sensitive attribute z is independent from the covariate x . By contrast, we show that bias estimation is one of the main difficulties of the biased bandit problem.

The linear bandit with biased feedback can be viewed as a stochastic partial monitoring game. With the terminology of partial monitoring, the biased problem considered in the present paper is globally observable but not locally observable: in this case, the optimal worst-case regret rate typically increases as $\tilde{O}(T^{2/3})$. This regret rate is for example achieved in the related problem of partial linear monitoring with linear feedback and linear reward using an Information Directed Sampling algorithm [17]. However, the dependence of the regret on the geometry of the action set and on the dimension d remains in most cases an open question [22, 5, 17]. In this paper, we characterize the geometry of the biased linear bandit problem, and we investigate dependence of the regret on the gaps.

1.2 Contribution and outline

In this paper, we introduce the linear bandit problem with biased feedback. We design a new algorithm based on optimal design for this problem. We derive an upper bound on the worst case regret of this algorithm of order $\kappa_*^{1/3} \log(T)^{1/3} T^{2/3}$ for large T , where κ_* is an explicit constant depending on the geometry of the action set. We provide matching lower bounds on some problem instances, showing that the constant κ_* characterizes the difficulty of the action set. Note that this regret is higher than the classical rates of order $\tilde{O}(dT^{1/2})$ obtained for d -dimensional linear bandits: this increase corresponds to the price to pay for debiasing the unfair evaluations.

We also characterize the gap-depend regret, showing that it is of order $(d/\Delta_{\min} \vee \kappa(\Delta)/\Delta_{\neq}^2) \log(T)$, where Δ_{\min} is the minimum gap, Δ_{\neq} is the gap between the best actions of the two groups, and $\kappa(\Delta)$ corresponds to the minimum regret to pay for estimating the bias with a given variance. This bound underlines the relative difficulties of the d -dimensional linear bandit and of the bias estimation. When $d/\Delta_{\min} \geq \kappa(\Delta)/\Delta_{\neq}^2$, i.e. when one group contains all near-optimal actions, the difficulty is dominated by that of the corresponding linear bandit problem. When both groups contain near-optimal actions, and $d/\Delta_{\min} \leq \kappa(\Delta)/\Delta_{\neq}^2$, the regret corresponds to the price of debiasing the rewards.

The rest of the paper is organized as follows. In Section 2, we present the FAIR PHASED ELIMINATION algorithm: we first discuss parameter estimation in Section 2.1, before presenting a sketch of the algorithm in Section 2.2 (a detailed version of this algorithm is provided in Appendix B). Then, in Section 3, we establish an upper bound on its worst-case regret. In Section 4, we derive a gap-dependent upper bound on the regret of our algorithm. In Section 5, we establish lower bounds on some action sets for both the worst-case and the gap-dependent regret, showing that these rates are sharp respectively up to a sub-logarithmic factor and an absolute multiplicative constant. Additional discussions on the geometry of bias estimation are postponed to Appendix A.

1.3 Notations and additional assumptions

We assume that all covariates $x \in \mathcal{X}$ are distinct, which implies that the group z_x of action x is well defined. We also assume that no group is empty, that the set $\left\{ \begin{pmatrix} x \\ z_x \end{pmatrix} : x \in \mathcal{X} \right\}$ spans \mathbb{R}^{d+1} (which guarantees identifiability of the parameters), and that the rewards are bounded: $\max_{x \in \mathcal{X}} |x^\top \gamma^*| \leq 1$.

When necessary, we underline the dependence of the regret on the parameter θ by denoting it R_T^θ . We denote by $a_x = \begin{pmatrix} x \\ z_x \end{pmatrix}$ the vector describing an action and its group, by $\theta^* = \begin{pmatrix} \gamma^* \\ \omega^* \end{pmatrix} \in \mathbb{R}^{d+1}$ the unknown parameter, and by $\mathcal{A} = \{a_x : x \in \mathcal{X}\}$ the set of actions and of corresponding sensitive attributes. We denote by $\Delta = (\Delta_x)_{x \in \mathcal{X}}$ the vector of gaps $\Delta_x = \max_{x' \in \mathcal{X}} (x' - x)^\top \gamma^*$, and by $\mathcal{C}(\mathcal{X}) = \{\gamma \in \mathbb{R}^d : \forall x \in \mathcal{X}, |x^\top \gamma| \leq 1\}$ the set of admissible parameters. Note that for all $x \in \mathcal{C}(\mathcal{X})$, $\Delta_x \leq 2$. For $i \leq d+1$, let e_i be the i -th vector of the canonical basis of \mathbb{R}^{d+1} , and for any matrix M , let M^+ be a generalized inverse of M . We denote by $\mathcal{P}^{\mathcal{X}}$ the set of probability measures on \mathcal{X} , and $\mathcal{M}^{\mathcal{X}} = \{\mu : \mathcal{X} \mapsto \mathbb{R}_+\}$. For any $\mu \in \mathcal{P}^{\mathcal{X}}$ or $\mu \in \mathcal{M}^{\mathcal{X}}$, we denote $V(\mu) = \sum_{x \in \mathcal{X}} \mu(x) a_x a_x^\top$ the covariance matrix corresponding to this allocation. Moreover, for $u \in \mathbb{R}^{d+1}$ (resp. $U \in \mathbb{R}^{d+1}$), we denote by $\mathcal{P}_u^{\mathcal{X}}$ (resp. $\mathcal{M}_u^{\mathcal{X}}$) the measures μ in $\mathcal{P}^{\mathcal{X}}$ (resp. in $\mathcal{M}^{\mathcal{X}}$) such that $u \in \text{Range}(V(\mu))$. For $U \subset \mathbb{R}^{d+1}$, we denote by $\mathcal{P}_U^{\mathcal{X}}$ (resp. $\mathcal{M}_U^{\mathcal{X}}$) the measures μ such that $\mu \in \mathcal{P}_u^{\mathcal{X}}$ (resp. $\mathcal{M}_u^{\mathcal{X}}$) for all $u \in U$.

2 Fair Phased Elimination algorithm

The Fair Phased Elimination algorithm belongs to the category of sequential elimination algorithms. Classical sequential elimination algorithms typically proceed by phases, indexed by $l = 1, 2, \dots$. At phase l , these algorithms consider a set of potentially optimal actions \mathcal{X}_l . The rewards of all actions $x \in \mathcal{X}_l$ are then estimated with a given precision $O(\epsilon_l)$, typically chosen as $\epsilon_l = 2^{2-l}$, by sampling actions in \mathcal{X}_l . Actions sub-optimal by a gap larger than the precision level are then removed from the set \mathcal{X}_{l+1} of potentially optimal actions for the phase $l+1$.

As underlined previously, sequential elimination algorithms may suffer linear regret in the biased linear bandit problem if actions allowing to estimate the bias are discarded by the algorithm before the best group is identified. To mitigate this problem, we first estimate the biased evaluations of the potentially optimal

actions, using ordinary least squares estimation. We then debias the estimations using an estimator for the bias relying on independent observations, which may be obtained by sampling sub-optimal actions. Before presenting the algorithm, let us discuss the estimation of the evaluations and of the bias parameter.

2.1 Optimal design for parameter estimation in the biased linear bandit

G-optimal design for biased evaluation estimation As in the Phased Elimination algorithm [21], we rely on G-optimal design to estimate the biased evaluations $a_x^\top \theta^*$ with small error uniformly over a set of actions \mathcal{X}_l . More precisely, for a given set of potentially optimal actions \mathcal{X}_l , we compute the G-optimal design solution to the problem

$$\underset{\pi \in \mathcal{P}_{\mathcal{X}_l}^{\mathcal{X}_l}}{\text{minimize}} \max_{x \in \mathcal{X}_l} a_x^\top (V(\pi))^+ a_x. \quad (\text{G-optimal design}) \quad (2)$$

This can be done using polynomial-time algorithms, relying for example on interior points method [32], or on mixed integer second-order cone programming [31]. The celebrated General Equivalence theorem of Kiefer [16] and Pukelsheim [29] states that the value of Equation (2) is bounded by $d + 1$. Let π^* denote any design solution to the G-optimal design problem (2), and let $\hat{\theta}$ denote the ordinary least square estimator obtained by sampling each action $x \in \mathcal{X}_l$ exactly $\lceil n\pi^*(x) \rceil$ times for a given $n > 0$. Then, for all $x \in \mathcal{X}_l$, the General Equivalence theorem implies that the variance of the estimate $a_x^\top \hat{\theta}$ is smaller than $(d+1)/n$. Moreover, the G-optimal design π^* can be chosen so that it is supported by at most $(d+1)(d+2)/2$ points, so the total number of samples is at most $n + (d+1)(d+2)/2$.

Δ -optimal design for bias evaluation In this paragraph, we introduce the Δ -optimal design, which is discussed in greater depth in Appendix A.5. To estimate the bias parameter ω^* , we use the estimator $\hat{\omega} = e_{d+1}^\top \hat{\theta}$, where $\hat{\theta}$ is the ordinary least square estimator for the full parameter θ^* . Now, if we sample each action $x \in \mathcal{X}$ exactly $\mu(x)$ time, the variance of $\hat{\omega}$ is equal to $e_{d+1}^\top V(\mu)^+ e_{d+1}$. Given the vector of gaps Δ , the design μ minimizing the regret of this exploration phase, while ensuring that the variance of $\hat{\omega}$ is smaller than 1, is solution of the problem

$$\underset{\mu \in \mathcal{M}_{\mathcal{X}}^{e_{d+1}}}{\text{minimize}} \sum_x \mu(x) \Delta_x \quad \text{such that} \quad e_{d+1}^\top V(\mu)^+ e_{d+1} \leq 1. \quad (\Delta\text{-optimal design}) \quad (3)$$

Let us denote μ^Δ a minimizer of (3), and $\kappa(\Delta) = \sum_{x \in \mathcal{X}} \mu^\Delta(x) \Delta_x$. Lemma 9 in Appendix A explains how to compute the design μ^Δ in polynomial time by adapting tools from c -optimal design. This lemma also shows that the support of μ^Δ can be chosen to be of cardinality at most $d + 1$. Then, choosing each action exactly $\lceil n\mu^\Delta(x) \rceil$ times for a given $n > 0$ allows us to estimate the bias with variance lower than n^{-1} and a regret no larger than $n\kappa(\Delta) + 2(d + 1)$. Obviously, we do not know the gap vector Δ beforehand, so we must estimate it as we go.

2.2 Outline of the Fair Phased Elimination algorithm

The Fair Phased Elimination algorithm, sketched in Algorithm 3, relies on the following key ideas. First, note that within a group, the order of the true rewards and of the biased evaluations are the same. Hence, within a group, we can use classical algorithms for linear bandits to choose the actions and estimate the biased evaluations with a controlled within-group regret: this is done using **G-exploration and elimination**. Second, to compare actions belonging to different groups, we independently estimate the bias parameter ω^* , using **Δ -exploration and elimination**. Finally, we underline that bias estimation may require to sample very sub-optimal actions. Therefore, it can be overly costly to estimate the bias up to the precision level required to identify the best group. To prevent this, we use a **stopping criteria**.

G-exploration and elimination At each phase $l = 1, 2, \dots$, we keep two sets of potentially optimal actions belonging to the groups $+1$ and -1 , denoted respectively $\mathcal{X}_l^{(+1)}$ and $\mathcal{X}_l^{(-1)}$. If we have not identified the group containing the best action, we run a G-EXP-ELIM routine 1 on each set $\mathcal{X}_l^{(z)}$ for $z = 1$ and $z = -1$.

This routine samples actions according to a rounded G-optimal design on $\mathcal{X}_l^{(z)}$, with a total number of observations chosen so that the biased evaluations of all actions in $\mathcal{X}_l^{(z)}$ are known with an error at most ϵ_l . The set $\mathcal{X}_{l+1}^{(z)}$ is obtained by removing from $\mathcal{X}_l^{(z)}$ actions whose estimated evaluations are sub-optimal by a gap larger than $3\epsilon_l$, compared to the empirical best action in the group. This allows to ensure that only actions sub-optimal by a gap $\mathcal{O}(\epsilon_l)$ remain in $\mathcal{X}_{l+1}^{(z)}$, and to estimate the gap vector Δ with a precision sufficient for Δ -optimal estimation.

If the group containing the best action has been identified, we discard the other group, and run a G-EXP-ELIM routine 1 on the set of potentially optimal actions in this group.

Routine 1 G-EXP-ELIM (\mathcal{X}, n, ϵ)

- 1: Compute G-optimal design π solution of (2) on \mathcal{X} , with $|\text{supp}(\pi)| \leq (d+1)(d+2)/2$
 - 2: Sample $\lceil n\pi(x) \rceil$ times each action a_x for $x \in \mathcal{X}$ ▷ G-optimal parameter estimation
 - 3: Compute the ordinary least square estimator $\hat{\theta}$
 - 4: $\mathcal{X}' \leftarrow \left\{ x \in \mathcal{X} : \max_{x' \in \mathcal{X}} (x' - x)^\top \hat{\theta} \leq 3\epsilon \right\}$ ▷ Suboptimal actions elimination
 - 5: **return** $\hat{\theta}$ and \mathcal{X}'
-

Δ -exploration and elimination If the group of the best action has not been found before phase l , we run the Δ -EXP-ELIM routine 2. More precisely, relying on a previous estimate $\hat{\Delta}^l$ of the gap vector Δ , we compute the $\hat{\Delta}^l$ -optimal design $\hat{\mu}$. We then estimate the bias using actions sampled according to a rounded version of this design, with a total number of observations chosen so that the error of bias estimation is smaller than ϵ_l , and use it to debias the reward estimation. If the debiased evaluation of the best action of each group are separated by a gap larger than $4\epsilon_l$, we consider that the best group is the one containing the empirical best action in terms of biased evaluation, and we discard the other group.

If we cannot find the best group, we rely on estimates of the bias and of the biased evaluations obtained during the previous round to update the estimate of the gap vector $\hat{\Delta}^{l+1}$.

Routine 2 Δ -EXP-ELIM ($\mathcal{X}, (\mathcal{X}^{(z)}, \hat{\theta}^{(z)})_{z \in \{-1, 1\}}, \hat{\Delta}, n, \epsilon$)

- 1: Compute $\hat{\Delta}$ -optimal design $(\hat{\mu}, \kappa(\hat{\Delta}))$ solution of (3) on \mathcal{X} , with $|\text{supp}(\hat{\mu})| \leq d + 1$
 - 2: Sample $\lceil n\hat{\mu}(x) \rceil$ times each action a_x for $x \in \mathcal{X}$ ▷ $\hat{\Delta}$ -optimal bias estimation
 - 3: Compute $\hat{\omega} = e_{d+1}^\top \hat{\theta}$, where $\hat{\theta}$ is the ordinary least square estimator
 - 4: **for** $z \in \{-1, 1\}$ and $x \in \mathcal{X}^{(z)}$ **do** $\hat{m}_x \leftarrow a_x^\top \hat{\theta}^{(z)} - z\hat{\omega}$ ▷ Debiased rewards estimation
 - 5: **if** $\exists z \in \{-1, 1\}$ such that $\max_{x \in \mathcal{X}^{(z)}} \hat{m}_x \geq \max_{x \in \mathcal{X}^{(-z)}} \hat{m}_x + 4\epsilon$ **then** $\mathcal{Z} \leftarrow \{z\}$ ▷ Group elimination
 - 6: **else** $\hat{\Delta}_x \leftarrow 2 \wedge (\max_{x' \in \mathcal{X}^{(-1)} \cup \mathcal{X}^{(1)}} \hat{m}_{x'} - \hat{m}_x + 4\epsilon)$ for all $x \in \mathcal{X}^{(-1)} \cup \mathcal{X}^{(1)}$
 - 7: **return** \mathcal{Z} and $\hat{\Delta}$
-

Stopping criteria As underlined previously, the Δ -EXP-ELIM routine samples actions that can be very sub-optimal. As a consequence, when the gap between the best two actions of each group is small, finding the best group can be overly costly in terms of regret. To prevent this, if the best group has not been found at stage l fulfilling $\epsilon_l \leq (\kappa(\hat{\Delta}^l) \log(T)/T)^{1/3}$, the bias estimation is stopped and the empirical best action in $\mathcal{X}_{l+1}^{(1)} \cup \mathcal{X}_{l+1}^{(-1)}$ is sampled for the remaining time (see Algorithm 3)

3 Upper bound on the worst-case regret of FAIR PHASED ELIMINATION

The regret of the FAIR PHASED ELIMINATION depends on the difficulty of estimating the bias parameter, captured by $\kappa(\Delta)$. Lemma 7 in Appendix A.5 shows that for all parameter $\gamma^* \in \mathcal{X}$, $\kappa(\Delta)$ is upper bounded

Algorithm 3 FAIR PHASED ELIMINATION (sketched)

```

1: input:  $\delta, T, \mathcal{X}, k = |\mathcal{X}|, \epsilon_l = 2^{2-l}$  for  $l \geq 1$ 
2: initialize:  $\mathcal{X}_1^{(+1)} \leftarrow \{x : z_x = 1\}, \mathcal{X}_1^{(-1)} \leftarrow \{x : z_x = -1\},$ 
3:    $\mathcal{Z}_1 \leftarrow \{-1, +1\}, \widehat{\Delta}^1 \leftarrow (2, \dots, 2), l \leftarrow 0$ 
4: while the budget is not spent do  $l \leftarrow l + 1$ 
5:   for  $z \in \mathcal{Z}_l$  do
6:      $(\widehat{\theta}^{(z)}, \mathcal{X}_{l+1}^{(z)}) \leftarrow \text{G-EXP-ELIM} \left( \mathcal{X}_l^{(z)}, \frac{2(d+1)}{\epsilon_l^2} \log \left( \frac{kl(l+1)}{\delta} \right), \epsilon_l \right)$ 
7:   if  $\mathcal{Z}_l = \{-1, +1\}$  then
8:     if  $\epsilon_l \leq \left( \kappa(\widehat{\Delta}^l) \log(T)/T \right)^{1/3}$  then ▷ Stop bias estimation
9:       Sample best action in  $\mathcal{X}_{l+1}^{(-1)} \cup \mathcal{X}_{l+1}^{(+1)}$  for the remaining time
10:    else
11:       $(\mathcal{Z}_{l+1}, \widehat{\Delta}^{l+1}) \leftarrow \Delta\text{-EXP-ELIM} \left( \mathcal{X}, \left( \mathcal{X}_{l+1}^{(z)}, \widehat{\theta}_l^{(z)} \right)_{z \in \{-1, 1\}}, \widehat{\Delta}^l, \frac{2}{\epsilon_l^2} \log \left( \frac{l(l+1)}{\delta} \right), \epsilon_l \right)$ 

```

by $2\kappa_*$, where κ_* is the *minimal variance of the bias estimator* given by

$$\kappa_* = \min_{\pi \in \mathcal{P}_{e_{d+1}}^{\mathcal{X}_{d+1}}} e_{d+1}^\top (V(\pi))^+ e_{d+1}.$$

The following theorem provides a bound on the worst case regret depending on κ_* . Proofs are postponed to Appendix C.1.

Theorem 1. *For the choice $\delta = T^{-1}$, there exists two numerical constants $C, C' > 0$ such that the following bound on the regret of the FAIR PHASED ELIMINATION algorithm 4 holds*

$$\begin{aligned}
 R_T &\leq C \left(\kappa_*^{1/3} T^{2/3} \log(T)^{1/3} + (d \vee \kappa_*) \log(T) + d^2 + d\kappa_*^{-1/3} T^{1/3} \log(kT) \log(T)^{-1/3} \right) \\
 &\leq C' \kappa_*^{1/3} T^{2/3} \log(T)^{1/3} \quad \text{for } T \geq \frac{((d \vee \kappa_*)^{3/2} \log(T)) \vee d^3}{\sqrt{\kappa_*}} \vee \frac{(d \log(kT))^3}{(\kappa_* \log(T))^2}.
 \end{aligned}$$

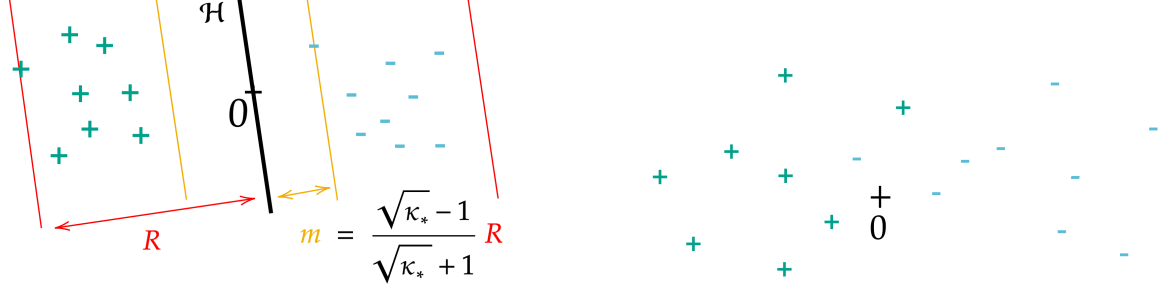
In Section 5.1, we show that the upper bound obtained in Theorem 1 is sharp in some settings, up to the sub-logarithmic factor $\log(T)^{1/3}$.

Theorem 1 shows that the worst-case regret of the Fair Phased Elimination algorithm asymptotically grows as $C\kappa_*^{1/3} T^{2/3} \log(T)^{1/3}$. This worst-case regret rate is higher than the typical rate $Cd \log(T) T^{1/2}$ obtained under unbiased feedback on the rewards (see, e.g., [1]). This increase in the regret corresponds to the cost of learning from unfair evaluations. It is due to the fact that the algorithm may need to sample actions that are sub-optimal in order to estimate the bias parameter. Note that this rate $\tilde{\mathcal{O}}(T^{2/3})$ is typical for globally observable bandit problems with partial linear monitoring, and can be obtained by applying results established in [17] for in the partial linear monitoring setting to the biased linear bandit problem.

By contrast to previous results, Theorem 1 characterizes precisely the dependence of the worst-case regret on the geometry of the action set. The relevant constant κ_* is the minimal variance for estimating the bias, which appears when considering the related c -optimal design problem. While the connection between G-optimal design and the linear bandit problem has already been exploited, it is to the best of our knowledge the first time that c -optimal design is related to a partial monitoring problem.

The constant κ_* corresponds to the minimum number of samples required for estimating the bias with a variance equal to 1 (up to rounding issues). Intuitively, if the actions are very correlated with their sensitive attributes, more samples will be needed to estimate the bias with the same precision. This situation corresponds to cases where κ_* is large, and leads to a higher regret. Lemma 1, illustrated in Figure 1, relates κ_* to the margin between the two groups of actions.

Lemma 1. *κ_* is the largest constant $\kappa \geq 0$ such that, there exists an hyperplane \mathcal{H} containing zero and separating the two groups, and such that, the margin to \mathcal{H} is at least $\sqrt{\kappa-1}/\sqrt{\kappa+1}$ times the maximum distance of all points to the hyperplane (see Figure 1). When no such hyperplane exists, then $\kappa_* = 1$.*



(a) The margin m is equal to $\sqrt{\kappa_*}-1/\sqrt{\kappa_*}+1$ times the maximum distance R of any action to the hyperplane. (b) $\kappa_* = 1$: the groups cannot be separated by a hyperplane containing 0.

Figure 1: Interpretation of κ_* in terms of separation of the groups.

Interestingly, Lemma 1 underlines that under reasonable assumptions, the constant κ_* may not depend on the ambient dimension d , and it can even be equal to 1. By contrast, the previous bounds obtained for an Information Directed Sampling algorithm are of order $\alpha^{1/3}d^{1/2}T^{2/3}\log(kT)^{1/2}$, where α is a measure of the complexity of the action set called the worst-case alignment constant. Lemma 6 in Appendix A shows that α is equivalent to the minimal variance of the bias estimator κ_* . Hence, our bound improves over previous results by a factor $d^{1/2}\log(T)^{1/6}(\log(kT)/\log(T))^{1/2}$.

The gaps are not involved in the definition of the minimal variance of bias estimation κ_* . The reader may have expected to get, instead of κ_* , the minimax regret for estimating the bias

$$\tilde{\kappa} = \max_{\gamma \in \mathcal{C}(\mathcal{X}), x' \in \mathcal{X}} \sum_{x \in \mathcal{X}} \tilde{\mu}(x)(x' - x)^\top \gamma, \quad \text{where}$$

$$\tilde{\mu} = \operatorname{argmin}_{\mu} \max_{x' \in \mathcal{X}, \gamma \in \mathcal{C}(\mathcal{X})} \sum_{x \in \mathcal{X}} \mu(x)(x' - x)^\top \gamma, \quad \text{such that } \mu \in \mathcal{M}_{e_{d+1}}^{\mathcal{X}} \text{ and } e_{d+1}^\top V(\mu)^+ e_{d+1} \leq 1.$$

Next lemma shows that κ_* and $\tilde{\kappa}$ are in equivalent up to a factor 2. We refer the interested reader to Appendix A, where further discussions on the geometry of bias estimation are postponed, due to space constraints.

Lemma 2. $\tilde{\kappa}/2 \leq \kappa_* \leq 2\tilde{\kappa}$.

4 Upper bound on the gap-depend regret of FAIR PHASED ELIMINATION

In this section, we provide an upper bound on the worst-case regret that depends on the gap between the two best actions, and on the gap between the best actions of the two groups. Compared to instance-dependent bounds, established in the linear bandit problem in [20, 18], gap-dependent bounds characterize the dependence of the regret on a small number of parameters. They are typically less sharp than instance-dependent bounds, but allow to better highlight the influence of the parameters on the difficulty of the problem. The bound established in the following theorem relates the difficulty of the biased linear bandit to that of bias estimation, and to that of the corresponding d -dimensional linear bandit. Proofs are postponed to Appendix C.1.

Theorem 2. *Assume that $x^* \in \operatorname{argmax}_{x \in \mathcal{X}} x^\top \gamma^*$ is unique. Then, there exists two numerical constants $C, C' > 0$ such that, for the choice $\delta = T^{-1}$, the following bound on the regret of the FAIR PHASED ELIMINATION algorithm 4 holds*

$$\begin{aligned} R_T &\leq C \left(\left(\frac{d}{\Delta_{\min}} \vee \frac{\kappa(\Delta \vee \Delta_{\neq} \vee \varepsilon_T)}{\Delta_{\neq}^2} \right) \log(T) + d^2 + \frac{d}{\Delta_{\min}} \log(k) \right) \\ &\leq C' \left(\frac{d}{\Delta_{\min}} \vee \frac{\kappa(\Delta \vee \Delta_{\neq} \vee \varepsilon_T)}{\Delta_{\neq}^2} \right) \log(T) \quad \text{for } T \geq k \vee e^{d\Delta_{\min}} \end{aligned}$$

where $\Delta_{\min} = \min_{x \in \mathcal{X} \setminus x^*} \Delta_x$, $\Delta_{\neq} = \min_{x \in \mathcal{X}: z_x = -z_{x^*}} \Delta_x$, and $\varepsilon_T = (\kappa_* \log(T)/T)^{1/3}$.

The term $d/\Delta_{\min} \vee \kappa(\Delta \vee \Delta_{\neq} \vee \varepsilon_T)/\Delta_{\neq}^2$ highlights the two sources of difficulty of the problem. On the one hand, the term d/Δ_{\min} is unavoidable: even if the algorithm knew beforehand the group containing the best action, it would still need to play a game of d -dimensional linear bandits in this group, and suffer, in the worst-case, the corresponding gap-dependent regret [1]. Note that lower bounds on gap-dependent regret of classical linear bandits follow from considering a setting with one near-optimal action with gap Δ_{\min} in each of the d dimensions. Then, any algorithm needs to explore each dimension up to $\Delta_{\min}^{-2} \log(T)$ times in order to find the best action, but can do so by choosing the near-optimal actions, thus having a regret $\Delta_{\min}^{-1} \log(T)$ in each direction. By contrast, the term $\kappa(\Delta \vee \Delta_{\neq} \vee \varepsilon_T)/\Delta_{\neq}^2$ is characteristic of the biased linear bandit problem: it is due to the fact that the algorithm may need to sample very sub-optimal actions in order to find the group containing the best action. Indeed, to identify this group, one must estimate the bias with a precision Δ_{\neq} , i.e. sample sub-optimal actions with average regret $\kappa(\Delta)$ approximately $\Delta_{\neq}^{-2} \log(T)$ times.

When $d/\Delta_{\min} \leq \kappa(\Delta \vee \Delta_{\neq} \vee \varepsilon_T)/\Delta_{\neq}^2$, the regret corresponds to the regret of this bias estimation phase. In other words, when both groups contain near-optimal actions, the difficulty of the problem is dominated by the price to pay for debiasing the unfair evaluations. Interestingly, when $d/\Delta_{\min} > \kappa(\Delta \vee \Delta_{\neq} \vee \varepsilon_T)/\Delta_{\neq}^2$, the difficulty of the linear bandit with systematic bias is dominated by that of the classical d -linear bandit. In this case, the algorithm is able to find the group containing the best action, and the problem reduces to a linear bandit in dimension d . Thus, the linear bandit with systematic bias is a non trivial example of a globally observable game that can be locally observable around the best action.

Finally, we underline that the magnitude of the bias does not appear in the regret: intuitively, no matter its magnitude, the algorithm always need to estimate it up to the same precision (of order Δ_{\neq}) in order to find the best group and to be optimal in terms of gap-dependent regret. This indicates that our algorithm is robust against important discriminations in the evaluation mechanism.

5 Lower bounds on the regret

In this section, we derive lower bounds on the worst-case regret and the gap-dependent regret that respectively match the upper bounds established in Theorems 1 and 2 up to sub-logarithmic factors or numerical constants.

5.1 Lower bound on the worst-case regret

Theorems 1 and 2 underline the dependence of the regret on the geometry of the action set. Before stating our result, we begin by introducing the notion of κ_* -correlated action set.

Definition 1 (κ_* -correlated action set). *For $\kappa_* \geq 1$, a set of actions \mathcal{A} is κ_* -correlated if $\mathcal{A} \in \mathbf{A}_{\kappa_*, d}$, where*

$$\mathbf{A}_{\kappa_*, d} = \left\{ \begin{array}{l} \mathcal{A} = \{a_1, \dots, a_k\} \subset \left(\mathbb{R}^d \times \{-1, +1\}\right)^k : \\ k \in \mathbb{N}^*, \min_{\pi \in \mathcal{P}_{e_{d+1}}^{\mathcal{A}}} \left\{ e_{d+1}^{\top} \left(\sum_{a \in \mathcal{A}} \pi(a) a a^{\top} \right)^+ e_{d+1} \right\} \geq \kappa_* \end{array} \right\}$$

is the set of actions sets such that the minimal variance of the bias estimator is larger than κ_ .*

In the following theorem, we establish a lower bound on the regret valid for all $\kappa_* \geq 1$ by designing κ_* -correlated sets of actions $\mathcal{A} \in \mathbf{A}_{\kappa_*, d}$, and obtaining lower bounds on the regret of any algorithm on these sets of actions.

Theorem 3. *Let $\kappa_* \geq 1$, $d \geq 2$ and $T \geq 4^3 \kappa_*$. There exists an action set $\mathcal{A} \in \mathbf{A}_{\kappa_*, d}$ such that for any algorithm, there exists a bandit problem with parameter $\theta_T \in \mathbb{R}^{d+1}$ such that the regret of this algorithm on the problem characterized by θ_T satisfies $R_T^{\theta_T} \geq \kappa_*^{1/3} T^{2/3} / 8e$.*

Previous lower bounds on the regret of linear bandits with partial monitoring, established in [17], state that the regret must be at least $c_{\mathcal{A}} T^{2/3}$ for some parameter $\theta_T \in \mathbb{R}^{d+1}$, where $c_{\mathcal{A}} > 0$ is a constant depending (not explicitly) on \mathcal{A} . By contrast, Theorem 3 provides an explicit characterization of the dependence of the regret rate on the geometry of the problem, which matches the upper bound of Theorem 1 up to a sub-logarithmic factor. Note that the assumption $d \geq 2$ is necessary here: if $d = 1$, there are at most two potentially optimal actions (namely, $\max\{x : x \in \mathcal{X}\}$ and $\min\{x : x \in \mathcal{X}\}$). Then, the problem becomes locally observable, and regret of order $\tilde{O}(T^{1/2})$ can be achieved [17].

5.2 Lower bound on the gap-dependent regret

We now present a lower bound on the gap-dependent regret. More precisely, for given values of Δ_{\min} and Δ_{\neq} , we establish a lower bound on the worst case regret among parameters θ verifying $\Delta_{\min} \leq \min_{x \in \mathcal{X} \setminus x^*} \Delta_x$, and $\Delta_{\neq} \leq \min_{x \in \mathcal{X} : z_x = -z_{x^*}} \Delta_x$. Before stating formally the result, let us define the corresponding parameter set. For an action set $\mathcal{A} \in \mathbf{A}_{\kappa_*, d}$, and for $(\Delta_{\min}, \Delta_{\neq}) \in (0, 1)^2$ such that $\Delta_{\min} \leq \Delta_{\neq}$, we denote

$$\Theta_{\Delta_{\min}, \Delta_{\neq}}^{\mathcal{A}} = \left\{ \begin{array}{l} \theta = \begin{pmatrix} \gamma \\ \omega \end{pmatrix} : \gamma \in \mathcal{C}(\mathcal{X}), \exists ! \begin{pmatrix} x^* \\ z_{x^*} \end{pmatrix} \in \operatorname{argmax}_{\begin{pmatrix} x \\ z_x \end{pmatrix} \in \mathcal{A}} \{x^\top \gamma\}, \\ \forall \begin{pmatrix} x' \\ z_{x'} \end{pmatrix} \in \mathcal{A} \text{ such that } x' \neq x^*, (x^* - x')^\top \gamma \geq \Delta_{\min}, \\ \forall \begin{pmatrix} x' \\ z_{x'} \end{pmatrix} \in \mathcal{A} \text{ such that } z_{x'} \neq z_{x^*}, (x^* - x')^\top \gamma \geq \Delta_{\neq} \end{array} \right\}$$

the set of parameters with minimum gap Δ_{\min} , and minimum between-group-gap Δ_{\neq} .

The upper bounds established in Theorem 2 underline the dependence of the gap-dependent regret on the minimal regret $\kappa(\Delta)$ for estimating the bias. Before stating our results, we define a class of problems $\Theta_{\Delta_{\min}, \Delta_{\neq}, \kappa}^{\mathcal{A}}$ such that $\kappa(\Delta) \leq \kappa$. For a parameter $\gamma \in \mathcal{C}(\mathcal{X})$, let us denote $\Delta(\gamma)_x = \max_{x' \in \mathcal{X}} (x' - x)^\top \gamma$, and $\Delta(\gamma) = (\Delta(\gamma)_x)_{x \in \mathcal{X}}$. Moreover, for a given set \mathcal{A} , let us denote

$$\Theta_{\Delta_{\min}, \Delta_{\neq}, \kappa}^{\mathcal{A}} = \Theta_{\Delta_{\min}, \Delta_{\neq}}^{\mathcal{A}} \cap \left\{ \theta = \begin{pmatrix} \gamma \\ \omega \end{pmatrix} : \gamma \in \mathcal{C}(\mathcal{X}), \kappa(\Delta(\gamma)) \leq \kappa \right\}.$$

Theorem 4. *For all $\kappa \geq 2$ and all $d \geq 4$, there exists a set of actions $\mathcal{A} \in \mathbb{R}^{d+1}$ such that for all $(\Delta_{\min}, \Delta_{\neq}) \in (0, 1/8)^2$ with $\Delta_{\min} \leq \Delta_{\neq}$,*

$$\liminf_{T \rightarrow \infty} \sup_{\theta \in \Theta_{\Delta_{\min}, \Delta_{\neq}, \kappa}^{\mathcal{A}}} \frac{R_T^\theta}{\log(T)} \geq \left\lfloor \frac{d}{10\Delta_{\min}} \right\rfloor \vee \left\lfloor \frac{\kappa + 2}{8\Delta_{\neq}^2} \right\rfloor. \quad (4)$$

Theorem 4 shows that for some action sets \mathcal{A} , the gap-depend regret of the FAIR PHASED ELIMINATION algorithm is asymptotically optimal up to a numerical constant. Note that the assumption $d \geq 4$ is necessary in our proof to design an action set \mathcal{A} such that Equation (4) holds for all $\Delta_{\min}, \Delta_{\neq} \in (0, 1/8)$. On the other hand, as discussed in Appendix C.5, for $d \geq 2$, for all $\Delta_{\min}, \Delta_{\neq} \in (0, 1/8)$, we can show that there exists action sets \mathcal{A} and $\theta \in \Theta_{\Delta_{\min}, \Delta_{\neq}}^{\mathcal{A}}$ such that the lower bound in Equation (4) still holds, by considering separately the cases $d/\Delta_{\min} > \kappa/\Delta_{\neq}^2$ and $d/\Delta_{\min} \leq \kappa/\Delta_{\neq}^2$.

6 Conclusion

In this paper, we addressed the problem of online decision making under biased bandit feedback. We designed a new algorithm based on Δ - and G-optimal design, and obtained worst-case and gap-dependent upper bounds on its regret. We obtained lower bounds on the regret for some problem instances showing that these rates are tight up to sub-logarithmic factors in some settings. These rates highlight two behaviors: on the one hand, the worst case rate $\mathcal{O}(\kappa_*^{1/3} \log(T)^{1/3} T^{2/3})$ highlights the cost induced by the biased feedback, and the need to select sub-optimal actions in order to debias it. On the other hand, the gap-dependent bound shows that for some instance, the problem can be locally observable around the best action: then, the difficulty of the problem is dominated by the difficulty of the corresponding linear bandit problem, and is no more difficult than this problem. When this is not the case, the regret scales as $\kappa(\Delta) \Delta_{\neq}^{-2} \log(T)$, where Δ_{\neq} is the gap between the best actions of the two groups, and $\kappa(\Delta)$ is the minimum regret for estimating the bias with a given precision. This work paves the way for studying other bandit models with unfair feedback, considering for example continuous, multi-dimensional sensitive attributes.

References

- [1] Y. Abbasi-Yadkori, D. Pál, and C. Szepesvári. Improved algorithms for linear stochastic bandits. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011.

- [2] A. Barik and J. Honorio. Fair sparse regression with clustering: An invex relaxation for a combinatorial problem. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.
- [3] A. Byanjankar, M. Heikkilä, and J. Mezei. Predicting credit risk in peer-to-peer lending: A neural network approach. In *2015 IEEE Symposium Series on Computational Intelligence*, pages 719–725, 2015.
- [4] L. E. Celis, S. Kapoor, F. Salehi, and N. K. Vishnoi. An algorithmic framework to control bias in bandit-based personalization, 2018.
- [5] S. Chaudhuri and A. Tewari. Phased exploration with greedy exploitation in stochastic combinatorial partial monitoring games. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [6] S. Chawla and M. Jagadeesan. Individual Fairness in Advertising Auctions Through Inverse Proportionality. In M. Braverman, editor, *13th Innovations in Theoretical Computer Science Conference (ITCS 2022)*, volume 215 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 42:1–42:21, Dagstuhl, Germany, 2022. Schloss Dagstuhl – Leibniz-Zentrum für Informatik.
- [7] E. Chzhen and N. Schreuder. A minimax framework for quantifying risk-fairness trade-off in regression, 2020.
- [8] H. Claire, Y. Chen, J. Modi, M. F. Jung, and S. Nikolaidis. Multi-armed bandits with fairness constraints for distributing resources to human teammates. *2020 15th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pages 299–308, 2020.
- [9] G. Elfving. Optimum Allocation in Linear Regression Theory. *The Annals of Mathematical Statistics*, 23(2):255 – 262, 1952.
- [10] J. Fauw, J. Ledsam, B. Romera-Paredes, S. Nikolov, N. Tomasev, S. Blackwell, H. Askham, X. Glorot, B. O’Donoghue, D. Visentin, G. Driessche, B. Lakshminarayanan, C. Meyer, F. Mackinder, S. Bouton, K. Ayoub, R. Chopra, D. King, A. Karthikesalingam, and O. Ronneberger. Clinically applicable deep learning for diagnosis and referral in retinal disease. *Nature Medicine*, 24, 09 2018.
- [11] J. Fellman. *On the Allocation of Linear Observations*. Commentationes physico-mathematicae. Societas Scientiarum Fennica, 1974.
- [12] A. Fuster, P. Goldsmith-Pinkham, T. Ramadorai, and A. Walther. Predictably unequal? the effects of machine learning on credit markets. *The Journal of Finance*, 77(1):5–47, 2022.
- [13] H. Hadiji, S. Gerchinovitz, J.-M. Loubes, and G. Stoltz. Diversity-Preserving K-Armed Bandits, Revisited. working paper or preprint, Oct. 2020.
- [14] R. Harman and T. Jurik. Computing c-optimal experimental designs using the simplex method of linear programming. *Computational Statistics & Data Analysis*, 53(2):247–254, dec 2008.
- [15] A. Khademi, S. Lee, D. Foley, and V. Honavar. Fairness in algorithmic decision making: An excursion through the lens of causality. In *The World Wide Web Conference, WWW ’19*, page 2907–2914, New York, NY, USA, 2019. Association for Computing Machinery.
- [16] J. Kiefer. General Equivalence Theory for Optimum Designs (Approximate Theory). *The Annals of Statistics*, 2(5):849 – 879, 1974.
- [17] J. Kirschner, T. Lattimore, and A. Krause. Information directed sampling for linear partial monitoring. In J. D. Abernethy and S. Agarwal, editors, *Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pages 2328–2369. PMLR, 2020.
- [18] J. Kirschner, T. Lattimore, C. Vernade, and C. Szepesvari. Asymptotically optimal information-directed sampling. In M. Belkin and S. Kpotufe, editors, *Proceedings of Thirty Fourth Conference on Learning Theory*, volume 134 of *Proceedings of Machine Learning Research*, pages 2777–2821. PMLR, 15–19 Aug 2021.

- [19] A. Köchling and M. C. Wehner. Discriminated by an algorithm: a systematic review of discrimination and fairness by algorithmic decision-making in the context of hr recruitment and hr development. *Business Research*, pages 1–54, 2020.
- [20] T. Lattimore and C. Szepesvári. The End of Optimism? An Asymptotic Analysis of Finite-Armed Linear Bandits. In A. Singh and J. Zhu, editors, *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pages 728–737. PMLR, 20–22 Apr 2017.
- [21] T. Lattimore and C. Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- [22] T. Lin, B. Abrahao, R. Kleinberg, J. Lui, and W. Chen. Combinatorial partial monitoring game with linear feedback and its applications. In E. P. Xing and T. Jebara, editors, *Proceedings of the 31st International Conference on Machine Learning*, number 2, pages 901–909, Beijing, China, 22–24 Jun 2014. PMLR.
- [23] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan. A survey on bias and fairness in machine learning. *ACM Computing Surveys*, 54(6), jul 2021.
- [24] Z. Papamitsiou and A. A. Economides. Learning analytics and educational data mining in practice: A systematic literature review of empirical evidence. *Journal of Educational Technology & Society*, 17(4):49–64, 2014.
- [25] V. Patil, G. Ghalme, V. Nair, and Y. Narahari. Achieving fairness in the stochastic multi-armed bandit problem. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(04):5379–5386, Apr. 2020.
- [26] A. Pázman. *Foundations of Optimum Experimental Design*. Mathematics and its Applications. Springer Netherlands, 1986.
- [27] C. Perlich, B. Dalessandro, T. Raeder, O. Stitelman, and F. Provost. Machine learning for targeted display advertising: transfer learning in action. *Machine Learning*, 95(1):103–127, 2014.
- [28] L. Pronzato and G. Sagnol. Removing inessential points in c- and A-optimal design. *Journal of Statistical Planning and Inference*, 213:233–252, 2021.
- [29] F. Pukelsheim. On linear regression designs which maximize information. *Journal of statistical planning and inference*, 4:339–364, 1980.
- [30] M. Raghavan, S. Barocas, J. M. Kleinberg, and K. E. C. Levy. Mitigating bias in algorithmic hiring: evaluating claims and practices. *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 2020.
- [31] G. Sagnol and R. Harman. Computing exact d-optimal designs by mixed integer second order cone programming. *The Annals of Statistics*, 43, 07 2013.
- [32] L. Vandenberghe, S. Boyd, and S.-P. Wu. Determinant maximization with linear matrix inequality constraints. *SIAM Journal on Matrix Analysis and Applications*, 19(2):499–533, 1998.
- [33] M. Černý and M. Hladík. Two complexity results on c-optimality in experimental design. *Computational Optimization and Applications*, 51(3):1397–1408, apr 2012.
- [34] L. Wang, Y. Bai, W. Sun, and T. Joachims. Fairness of exposure in stochastic bandits. In M. Meila and T. Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 10686–10696. PMLR, 18–24 Jul 2021.

Appendix

The Appendix is organized as follows. We begin in Section A by further discussing the interpretation and computation of κ_* and $\kappa(\Delta)$, and their relation to the worst-case alignment constant of [17] and to the problem of optimal estimation of the bias against the worst parameter. Then, we provide in Section B a detailed version of the FAIR PHASED ELIMINATION algorithm 3. Then, in Section C, we prove the main results of this paper.

A On the geometry of bias estimation

The constants κ_* and $\kappa(\Delta)$ respectively characterize the difficulty of the worst-case problem, and of the gap-depend problem, and highlight the dependence of the regret on the geometry of the action set. In this section, we begin by discussing in Section A.1 the interpretation of the constant κ_* as the variance of the e_{d+1} -optimal design. Using Elfving's characterization of the e_{d+1} -optimal design, we then derive an alternative characterization of κ_* in terms of separation of the actions of the two groups in Section A.2

A.1 Bias estimation as a e_{d+1} -optimal design problem

Recall that κ_* is the minimal variance of the bias estimator related to the problem of e_{d+1} -optimal design.

e_{d+1} -optimal design Optimal design theory addresses the following problem: a scientist must design a set of n experiments $\{x_1, \dots, x_n\} \in \mathcal{X}^n$ so as to estimate at best a parameter of interest, where each experiment $x \in \mathcal{X}$ corresponds to a point $a_x \in \mathbb{R}^{d+1}$. The aim of the scientist is to choose a design, i.e. a function $\mu : \mathcal{X} \mapsto \mathbb{N}$ indicating the budget $\mu(x)$ to be allocated to each experiment $x \in \mathcal{X}$. Each experiment x is then repeated exactly $\mu(x)$ times, and the corresponding observations $y_{x,1}, \dots, y_{x,\mu(x)}$ are collected for each $x \in \mathcal{X}$. The law of the observations corresponding to experiment x at point a_x is given by

$$y_{x,i} = a_x^\top \theta^* + \xi_{x,i},$$

where $\xi_{x,i} \sim \mathcal{N}(0, 1)$ are independent noise terms, and $\theta^* \in \mathbb{R}^{d+1}$ is an unknown parameter. The aim of the scientist is to choose the design μ so as to best estimate (some features of) the parameter θ^* , under a constraint on the total number of experiments $\sum_{x \in \mathcal{X}} \mu(x) \leq n$ for some $n \in \mathbb{N}$.

Different criteria can be used to characterize the optimality of a design μ . For example, one may need to estimate the full parameter θ^* , in order to predict the outcomes of the experiments $x \in \mathcal{X}$ with a small uniform error: this leads to the G-optimal design problem (2). Alternatively, for c a vector in \mathbb{R}^{d+1} , one may aim at finding the best design $\mu \in \mathcal{N}^{\mathcal{X}}$ for estimating the scalar product $c^\top \theta^*$ under a budget constraint $\sum_{x \in \mathcal{X}} \mu(x) \leq n$, where $\mathcal{N}^{\mathcal{X}} = \{\mu : \mathcal{X} \rightarrow \mathbb{N}\}$. This problem is known as *c-optimal design*. Unbiased linear estimation of $c^\top \theta^*$ is possible only when c belongs to the image of $V(\mu)$, and in this case the best linear unbiased estimator of the scalar product $c^\top \theta^*$ is given by $c^\top \hat{\theta}$, where $\hat{\theta}$ is the least-square estimator defined as

$$\hat{\theta} = V(\mu)^+ \sum_{x \in \mathcal{X}} a_x \left(\sum_{i \leq \mu(x)} y_{x,i} \right) \quad \text{for} \quad V(\mu) = \sum_{x \in \mathcal{X}} \mu(x) a_x a_x^\top. \quad (5)$$

The variance of the estimator $c^\top \hat{\theta}$ is then equal to $c^\top V(\mu)^+ c$.

Exact *c-optimal design* aims at choosing the allocation $\mu \in \mathcal{N}^{\mathcal{X}}$ minimizing the variance of $c^\top \hat{\theta}$ for a given budget $\sum_x \mu(x) \leq n$, under the constraint that $c \in \text{Range}(V(\mu))$. Let us define the normalized design $\pi : x \in \mathcal{X} \mapsto \mu(x)/n$, and let us underline that π defines a probability on \mathcal{X} . The variance of $c^\top \hat{\theta}$ is then equal to $n^{-1} c^\top V(\pi)^+ c$. In the limit $n \rightarrow +\infty$, the problem is equivalent to the problem of approximate *c-optimal design* (sometimes simply referred to as *c-optimal design*), that aims at finding a probability measure $\pi \in \mathcal{P}_c^{\mathcal{X}} := \{\pi \in \mathcal{P}^{\mathcal{X}} : c \in \text{Range}(V(\pi))\}$ solution to the following problem

$$\min_{\pi \in \mathcal{P}_c^{\mathcal{X}}} c^\top V(\pi)^+ c. \quad (\text{c-optimal design})$$

Note that when $\{a_x : x \in \mathcal{X}\}$ spans \mathbb{R}^{d+1} , for any $c \in \mathbb{R}^{d+1}$, there exists a design π such that $c \in \text{Range}(V(\pi))$, and hence the c -optimal design problem admits a solution.

Computation of the e_{d+1} -optimal design Finding an exact optimal allocation $\mu \in \mathcal{N}^{\mathcal{X}}$ under the constraint that $\sum_{x \in \mathcal{X}} \mu(x) \leq n$ is unfortunately NP-complete. However, finding an approximate optimal design $\pi \in \mathcal{P}_c^{\mathcal{X}}$ can be done in polynomial time [33]. Several algorithms, including multiplicative algorithms [11] and a simplex method of linear programming [14], have been proposed to iteratively approximate the optimal design. More recently, [28] suggested using screening tests to remove inessential points to accelerate optimization algorithms.

Classical results from e_{d+1} -optimal design show that there exists a c -optimal design supported by at most $d + 1$ points (see, e.g., [26, 14] for a proof of this result). The following Lemma indicates how to obtain an exact design by rounding an approximate design supported by at most $d + 1$ points.

Lemma 3. *For any $\pi \in \mathcal{M}_{e_{d+1}}^{\mathcal{X}}$ and any $m > 0$, the estimator $e_{d+1}^\top \hat{\theta}_\mu$ computed from the design $\mu : x \mapsto \lceil m\pi(x) \rceil$ is an unbiased estimator of $e_{d+1}^\top \theta$ and it has a variance at most $m^{-1} e_{d+1}^\top V(\pi)^+ e_{d+1}$.*

Obviously, similar results also hold for G-optimal design.

Lemma 4. *Let π be a solution of the G-optimal design problem (2). Then, for any $m > 0$ and any $x \in \mathcal{X}$, the estimator $a_x^\top \hat{\theta}_\mu$ computed from the design $\mu : x \mapsto \lceil m\pi(x) \rceil$ is an unbiased estimator of the evaluation $a_x^\top \theta$, and it has a variance*

$$a_x^\top V(\mu)^+ a_x \leq m^{-1}(d + 1).$$

A.2 Interpretation of κ_* in terms of separation of the groups

Next theorem, due to Elfving, characterizes solutions to the c -optimal design problem.

Theorem 5 ([9]). *Let $\mathcal{S} = \text{convex hull}\{+a_x, -a_x : x \in \mathcal{X}\}$ be the Elfving's set of $\{a_x : x \in \mathcal{X}\} \subset \mathbb{R}^{d+1}$, and let $\partial\mathcal{S}$ denote the boundary of \mathcal{S} . A design $\pi \in \mathcal{P}_c^{\mathcal{X}}$ is c -optimal for $c \in \mathbb{R}^{d+1}$ if and only if there exists $\zeta \in \{-1, +1\}^{\mathcal{X}}$ and $t > 0$ such that*

$$tc = \sum_{x \in \mathcal{X}} \pi(x) \zeta_x a_x \in \partial\mathcal{S}.$$

Moreover, $t^{-2} = c^\top (V(\pi))^+ c$ is value of the c -optimal design problem.

Elfving's characterization of the e_{d+1} -optimal design allows us to derive the following equivalent characterization of κ_* .

Lemma 5. $\kappa_* = \max_{u \in \mathbb{R}^d} \frac{1}{\max_{x \in \mathcal{X}} (x^\top u + z_x)^2}$.

Lemma 1 follows from the characterization in Lemma 5. When $\kappa_* > 1$, the vector \tilde{u} defined as $\tilde{u} = \text{argmax}_{u \in \mathbb{R}^d} \frac{1}{\max_{x \in \mathcal{X}} (x^\top u + z_x)^2}$ is a normal vector of the separating hyperplane \mathcal{H} in Figure 1. Moreover, as shown in the proof of Lemma 1, the margin is in this case equal to $1 - \kappa_*^{-1/2}$, while the maximum distance of all points to the hyperplane is $1 + \kappa_*^{-1/2}$.

Lemma 5 also allows us to compare the bound in Theorem 1 with previous results on linear bandit with partial monitoring, expressed in terms of the worst-case alignment constant.

A.3 Comparison to the worst-case alignment constant

Previous work on linear bandit with partial linear monitoring measures the difficulty of the bandit game using the *worst-case alignment constant* α , defined as

$$\alpha = \max_{u \in \mathbb{R}^d} \frac{\max_{x, x' \in \mathcal{X}} ((x - x')^\top u)^2}{\max_{x \in \mathcal{X}} (z_x x^\top u + 1)^2}.$$

The following Lemma shows that this constant is essentially equivalent to the minimal variance of the bias estimator κ_* .

Lemma 6. $\frac{\kappa_*}{3} \leq \alpha \leq 16\kappa_*$.

On the one hand, Lemma 6 shows that κ_* and α are essentially equivalent. In particular, Theorem 3 implies that the large T regret is of order $\alpha^{1/3} \log(T)^{1/3} T^{2/3}$. This improves over previous known rates, obtained in [17], by a factor $d^{1/2} \log(T)^{1/6} (\log(kT)/\log(T))^{1/2}$.

On the other hand, as underlined, the constant κ_* appears when considering the well-studied problem of c -optimal design. Therefore, classical results and algorithms for optimal design can be used to characterize and compute this constant.

A.4 Optimal bias estimation against the worst parameter

The constant κ_* also appears naturally when considering the related problem of optimal bias estimation against the worst parameter.

Regret of e_{d+1} -optimal design Recall that κ_* denotes the *minimal variance of the bias estimator*, i.e. the value of the solution of the e_{d+1} -optimal design problem

$$\kappa_* = \min_{\pi \in \mathcal{P}_{e_{d+1}}^{\mathcal{X}}} e_{d+1}^{\top} (V(\pi))^+ e_{d+1},$$

The e_{d+1} -optimal design can be equivalently defined as the solution of the problem

$$\text{minimize } \sum_{x \in \mathcal{X}} \mu(x) \quad \text{such that } \mu \in \mathcal{M}_{e_{d+1}}^{\mathcal{X}} \text{ and } e_{d+1}^{\top} V(\mu)^+ e_{d+1} \leq \kappa_*. \quad (6)$$

The characterization given in Equation (6) underlines that the e_{d+1} -optimal design provides (up to discretization issues) the minimal number of samples required for estimating ω^* with a variance κ_* . Let us denote by μ^* the optimal design for estimating ω^* with a variance 1, defined as

$$\mu^* = \underset{\mu}{\operatorname{argmin}} \sum_{x \in \mathcal{X}} \mu(x) \quad \text{such that } \mu \in \mathcal{M}_{e_{d+1}}^{\mathcal{X}} \text{ and } e_{d+1}^{\top} V(\mu)^+ e_{d+1} \leq 1.$$

Note that from the definition of κ_* , we have $\sum_x \mu^*(x) = \kappa_*$.

A first (naive) approach to obtain an estimate of the bias parameter ω^* with precision level $\epsilon > 0$ would consist in sampling actions according to $\epsilon^{-2} \mu^*$, rounded according to the procedure defined in Lemma 3. Let us denote by Δ_x the gap $\Delta_x = \max_{x' \in \mathcal{X}} (x' - x)^{\top} \gamma^*$ between the (non-observed) reward of the best action and the reward of the action x . The regret corresponding to this estimation phase would then be

$$\epsilon^{-2} \sum_{x \in \mathcal{X}} \mu^*(x) \Delta_x,$$

which can be as large as $\kappa_* \epsilon^{-2} \max_x \Delta_x$. Interestingly, we show that the regret corresponding to the e_{d+1} -optimal design is equivalent (up to a small multiplicative constant) to the minimax regret.

Optimal worst-case estimation The minimax regret corresponds to the regret of the best sampling scheme against the worst admissible parameter γ . Note that, for a given design μ , this worst-case regret is given by

$$\max_{x' \in \mathcal{X}, \gamma \in \mathcal{C}(\mathcal{X})} \sum_x \mu(x) (x' - x)^{\top} \gamma,$$

where we recall that $\mathcal{C}(\mathcal{X}) = \{\gamma \in \mathbb{R}^d : \forall x \in \mathcal{X}, |x^{\top} \gamma| \leq 1\}$ is the set of admissible parameters. To achieve the lowest regret against the worst parameter, we must use the minimax optimal design $\tilde{\mu}$ solution to the problem

$$\tilde{\mu} = \underset{\mu}{\operatorname{argmin}} \max_{x' \in \mathcal{X}, \gamma \in \mathcal{C}(\mathcal{X})} \sum_{x \in \mathcal{X}} \mu(x) (x' - x)^{\top} \gamma \quad \text{such that } \mu \in \mathcal{M}_{e_{d+1}}^{\mathcal{X}} \text{ and } e_{d+1}^{\top} V(\mu)^+ e_{d+1} \leq 1.$$

Lemma 2 underlines that the regret corresponding to the e_{d+1} -optimal design is no larger than twice the minimax regret.

A.5 On the Δ -optimal design

Recall that for a vector of gaps $\Delta = (\Delta_x)_{x \in \mathcal{X}}$, μ^Δ denotes the Δ -optimal design, defined as the solution of the following problem

$$\mu^\Delta = \underset{\mu}{\operatorname{argmin}} \sum_{x \in \mathcal{X}} \mu(x) \Delta_x \quad \text{such that } \mu \in \mathcal{M}_{e_{d+1}}^{\mathcal{X}} \text{ and } e_{d+1}^\top V(\mu)^+ e_{d+1} \leq 1. \quad (\Delta\text{-optimal design})$$

If we knew the gaps Δ_x , we could sample the actions according to the Δ -optimal design μ^Δ , and pay the regret $\epsilon^{-2} \kappa(\Delta)$ (up to rounding error) for estimating ω^* with an error smaller than ϵ , where

$$\kappa(\Delta) = \sum_{x \in \mathcal{X}} \mu^\Delta(x) \Delta_x.$$

Lemma 7. *If $\gamma^* \in \mathcal{C}(\mathcal{X})$, then $\kappa(\Delta) \leq 2\kappa_*$*

Proof. Be definition of $\mathcal{C}(\mathcal{X})$, for all $\gamma^* \in \mathcal{C}(\mathcal{X})$, all $x, x' \in \mathcal{X}$, we have

$$(x - x')^\top \gamma^* \leq |x^\top \gamma^*| + |x'^\top \gamma^*| \leq 2.$$

Then,

$$\kappa(\Delta) \leq 2 \min_{\mu} \sum_{x \in \mathcal{X}} \mu(x) \quad \text{such that } \mu \in \mathcal{M}_{e_{d+1}}^{\mathcal{X}} \text{ and } e_{d+1}^\top V(\mu)^+ e_{d+1} \leq 1.$$

Let μ_* be the solution of the e_{d+1} -optimal design problem

$$\underset{\mu}{\operatorname{minimize}} e_{d+1}^\top V(\mu)^+ e_{d+1} \quad \text{such that } \mu \in \mathcal{P}_{e_{d+1}}^{\mathcal{X}}.$$

By definition of κ_* , we see that $e_{d+1}^\top V(\mu_*)^+ e_{d+1} = \kappa_*$. This implies that the measure $\kappa_* \times \mu_*$ verifies the constraints $e_{d+1}^\top V(\kappa_* \times \mu_*)^+ e_{d+1} \leq 1$ and $\kappa_* \mu_* \in \mathcal{M}_{e_{d+1}}^{\mathcal{X}}$. Thus,

$$\kappa(\Delta) \leq 2 \sum_{x \in \mathcal{X}} \kappa_* \mu_*(x) = 2\kappa_*.$$

□

On the regret $\kappa(\Delta)$ The function κ verifies the following properties.

Lemma 8. *For two vectors of gaps Δ, Δ' , denote by $\Delta \wedge \Delta'$ (respectively $\Delta \vee \Delta'$) the vector of gaps given by $(\Delta \wedge \Delta')_x = \Delta_x \wedge \Delta'_x$ (respectively $(\Delta \vee \Delta')_x = \Delta_x \vee \Delta'_x$) for all $x \in \mathcal{X}$. Moreover, denote $\Delta \leq \Delta'$ if $\Delta_x \leq \Delta'_x$ for all $x \in \mathcal{X}$. Then, the following properties hold :*

- i) for all $c > 0$, $\kappa(c\Delta) = c\kappa(\Delta)$;
- ii) if $\Delta \leq \Delta'$, then $\kappa(\Delta) \leq \kappa(\Delta')$;
- iii) $\kappa(\Delta \vee \Delta') \geq \kappa(\Delta) \vee \kappa(\Delta')$;
- iv) the function $\epsilon \mapsto \kappa(\Delta \vee \epsilon)$ is continuous at 0.

Computation of the Δ -optimal design In practice, the Δ -optimal design can be computed by adapting algorithms designed for finding the e_{d+1} -optimal design. Indeed, the next lemma shows that the computation of the Δ -optimal design amounts to computing an e_{d+1} -optimal design for some rescaled features.

Lemma 9. *For any vector $\Delta \in (0, +\infty)^{\mathcal{X}}$, let π^Δ be the e_{d+1} -optimal design relative to the set $\mathcal{A}^\Delta = \left\{ \Delta_x^{-1/2} \begin{pmatrix} x \\ z_x \end{pmatrix} : x \in \mathcal{X} \right\}$ and let $\kappa^\Delta = e_{d+1}^\top V(\pi^\Delta)^+ e_{d+1}$ be the e_{d+1} -optimal variance relative to \mathcal{A}^Δ . Then, the Δ -optimal design μ^Δ is given by $\mu^\Delta(x) = \kappa^\Delta \pi^\Delta(x) \Delta_x^{-1}$ for all $x \in \mathcal{X}$. In addition, the support of μ^Δ can be chosen to be of cardinality at most $d + 1$.*

B Detailed Fair Phased Elimination algorithm

We present the notations used in Algorithm 4. The phases are indexed by $l \in \mathbb{N}^*$. The sets $\mathcal{X}_l^{(z)}$ for $z \in \{-1, +1\}$ corresponds to actions in group z that are considered as potentially optimal in phase l . The variable \widehat{z}_l^* encodes the group determined as optimal: it is 0 as long as this group has not been determined. The subscript (z) refer to the group z when $z \in \{-1, +1\}$, and otherwise to the estimation of the bias ω^* : for example, the probability $\pi_l^{(z)}$ for $z \in \{-1, +1\}$ and $l > 1$ corresponds to the approximate G-optimal design on $\mathcal{X}_l^{(z)}$. Then, for $z \in \{-1, +1\}$, allocations $\mu^{(z)}$ (resp. $\mu^{(0)}$) correspond to allocation of samples in the exploration phase $\text{Exp}_l^{(z)}$ (resp. $\text{Exp}_l^{(0)}$). Similarly, $V_l^{(z)}$ (resp $V_l^{(0)}$) denotes the variance matrix of the estimator $\begin{pmatrix} \widehat{\gamma}_l^{(z)} \\ \widehat{\omega}_l^{(z)} \end{pmatrix}$ (resp. $\widehat{\omega}_l^{(0)}$) obtained from observations made during phase $\text{Exp}_l^{(z)}$ (resp. $\text{Exp}_l^{(0)}$). Finally, $\text{Explore}_l^{(z)}$ (resp. $\text{Explore}_l^{(0)}$) is a Boolean variable indicating whether the exploration at phase l for group z (resp. for the bias parameter) has been performed. It is used in the proofs to ensure that the corresponding estimators are well defined.

Algorithm 4 Fair Phased Elimination (detailed version)

1: **Input:** $\delta, T, k = |\mathcal{X}|$
2: **Initialize:** Recovery $\leftarrow \emptyset, t \leftarrow 0, l \leftarrow 1, \hat{z}_1^* \leftarrow 0,$
3: $\mathcal{X}_1^{(+1)} \leftarrow \{x : z_x = 1\}, \mathcal{X}_1^{(-1)} \leftarrow \{x : z_x = -1\}, \hat{\Delta}_x^1 \leftarrow 2$ for $x \in \mathcal{X}$
4: **while** $t < T$ **do**
5: **Initialize:** $\epsilon_l \leftarrow 2^{2-l}, \hat{z}_{l+1}^* \leftarrow \hat{z}_l^*, \hat{\Delta}^{l+1} \leftarrow \hat{\Delta}^l, \text{Explore}_l^{(z)} \leftarrow \text{False}$ for $z \in \{-1, 0, +1\}$
6: **for** $z \in \{-1, +1\}$ such that $z \neq -\hat{z}_l^*$ **do** ▷ G-optimal Exploration and Elimination
7: $\pi_l^{(z)} \leftarrow \underset{\pi}{\operatorname{argmin}} \left\{ \max_{x \in \mathcal{X}_l^{(z)}} a_x^\top V(\pi) + a_x : \pi \in \mathcal{P}_{\mathcal{X}_l^{(z)}}, |\operatorname{supp}(\pi)| \leq \frac{(d+1)(d+2)}{2} \right\}$
8: $\mu_l^{(z)}(x) \leftarrow \left\lceil \frac{2(d+1)\pi_l^{(z)}(x)}{\epsilon_l^2} \log \left(\frac{kl(l+1)}{\delta} \right) \right\rceil$ for all $x \in \mathcal{X}_l^{(z)}$
9: $n_l^{(z)} \leftarrow \sum_{x \in \mathcal{X}_l^{(z)}} \mu_l^{(z)}(x), \text{Exp}_l^{(z)} \leftarrow \{t+1, \dots, T \wedge (t+n_l^{(z)})\}$
10: **if** $t+n_l^{(z)} \leq T$ **then**
11: $\text{Explore}_l^{(z)} \leftarrow \text{True}$, choose each action $x \in \mathcal{X}_l^{(z)}$ exactly $\mu_l^{(z)}(x)$ times
12: $V_l^{(z)} \leftarrow \sum_{t \in \text{Exp}_l^{(z)}} a_{x_t} a_{x_t}^\top, \hat{\theta}_l^{(z)} \leftarrow \left(V_l^{(z)} \right)^+ \sum_{t \in \text{Exp}_l^{(z)}} y_t a_{x_t}$
13: $\mathcal{X}_{l+1}^{(z)} \leftarrow \left\{ x \in \mathcal{X}_l^{(z)} : \max_{x' \in \mathcal{X}_l^{(z)}} (a_{x'} - a_x)^\top \hat{\theta}_l^{(z)} \leq 3\epsilon_l \right\}$
14: **else** for $t \in \text{Exp}_l^{(z)}$, sample empirical best action in $\mathcal{X}_l^{(z)}$
15: $t \leftarrow t + n_l^{(z)}$
16: **if** $\hat{z}_l^* = 0$ **then**
17: compute the $\hat{\Delta}^l$ -optimal design $\hat{\mu}_l$ and the corresponding regret $\kappa(\hat{\Delta}^l)$
18: **if** $\epsilon_l \leq \left(\kappa(\hat{\Delta}^l) \log(T)/T \right)^{1/3}$ **then** ▷ Recovery phase
19: Recovery $\leftarrow \{t, \dots, T\}$
20: sample empirical best action in $\mathcal{X}_{l+1}^{(-1)} \cup \mathcal{X}_{l+1}^{(1)}$ until the end of the budget, $t \leftarrow T$
21: **else** ▷ $\hat{\Delta}^l$ -optimal Exploration and Elimination
22: $\mu_l^{(0)}(x) \leftarrow \left\lceil \frac{2\hat{\mu}_l(x)}{\epsilon_l^2} \log \left(\frac{l(l+1)}{\delta} \right) \right\rceil$ for all $x \in \mathcal{X}$
23: $n_l^{(0)} \leftarrow \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x), \text{Exp}_l^{(0)} \leftarrow \{t, \dots, T \wedge (t+n_l^{(0)})\}$
24: **if** $t+n_l^{(0)} \leq T$ **then**
25: $\text{Explore}_l^{(0)} \leftarrow \text{True}$, choose each action $x \in \mathcal{X}$ exactly $\mu_l^{(0)}(x)$ times
26: $V_l^{(0)} \leftarrow \sum_{t \in \text{Exp}_l^{(0)}} a_{x_t} a_{x_t}^\top, \hat{\omega}_l^{(0)} \leftarrow e_{d+1}^\top \left(V_l^{(0)} \right)^+ \sum_{t \in \text{Exp}_l^{(0)}} y_t a_{x_t}$
27: **for** $x \in \mathcal{X}_{l+1}^{(-1)} \cup \mathcal{X}_{l+1}^{(1)}$ **do**
28: $\hat{m}_{l,x} \leftarrow a_x^\top \hat{\theta}_l^{(z_x)} - z_x \hat{\omega}_l^{(0)}$
29: $\hat{\Delta}_x^{l+1} \leftarrow \left(\max_{x' \in \mathcal{X}_{l+1}^{(-1)} \cup \mathcal{X}_{l+1}^{(1)}} \hat{m}_{l,x'} - \hat{m}_{l,x} + 4\epsilon_l \right) \wedge 2$
30: **for** $z \in \{-1, +1\}$ **do**
31: **if** $\max_{x \in \mathcal{X}_{l+1}^{(z)}} \hat{m}_{l,x} - 2\epsilon_l \geq \max_{x \in \mathcal{X}_{l+1}^{(-z)}} \hat{m}_{l,x} + 2\epsilon_l$ **then** $\hat{z}_{l+1}^* \leftarrow z$
32: **else** sample empirical best action in $\mathcal{X}_{l+1}^{(-1)} \cup \mathcal{X}_{l+1}^{(1)}$ until the end of the budget, $t \leftarrow T$
33: $t \leftarrow t + n_l^{(0)}$
34: $l \leftarrow l + 1$

C Proofs

For an event \mathcal{F} such that $\mathbb{P}(\mathcal{F}) > 0$, we denote by $\mathbb{E}_{|\mathcal{F}}$ (resp. $\mathbb{P}_{|\mathcal{F}}$) the expectation (resp. the probability) conditionally on \mathcal{F} .

C.1 Proof of Theorem 1

We begin by defining for $z \in \{-1, 0, +1\}$

$$L^{(z)} = \max \left\{ l \geq 1 : \text{Explore}_l^{(z)} = \text{True} \right\}$$

the largest integer l such that $\text{Explore}_l^{(z)} = \text{True}$. Recall that κ_* is the e_{d+1} -optimal variance. By definition of the algorithm, for all $l \leq L^{(0)} + 1$, $\hat{\Delta}^l \leq 2$, so $\kappa(\hat{\Delta}^l) \leq 2\kappa_*$. Now, let us also define

$$L_T = \max \left\{ l \geq 1 : \epsilon_l > \left(\frac{2\kappa_* \log(T)}{T} \right)^{1/3} \right\}.$$

Then, if $\text{Recovery} \neq \emptyset$, we must have $L^{(0)} \geq L_T$. Moreover, we see that since $\epsilon_{L_T} = 2^{2-L_T}$, we have $L_T \leq 2 + \frac{\log_2(T/(2\kappa_* \log(T)))}{3} \leq 3 \log_2(T)$ when $T > 1$.

We define a "bad" event \mathcal{F} , such that, on $\bar{\mathcal{F}}$, our estimators $\hat{\gamma}_l^{(z)}$ and $\hat{\omega}_l^{(z)}$ are close to the true parameters γ^* and ω^* for all rounds l . More precisely, let

$$\mathcal{F} = \bigcup_{l \geq 1} \mathcal{F}_l, \tag{7}$$

where for $l \geq 1$

$$\mathcal{F}_l = \left\{ \exists z \in \{-1, 1\} \text{ such that } \text{Explore}_l^{(z)} = \text{True}, \text{ and } x \in \mathcal{X}_l^{(z)} \text{ such that } \left| \begin{pmatrix} \hat{\gamma}_l^{(z)} - \gamma^* \\ \hat{\omega}_l^{(z)} - \omega^* \end{pmatrix}^\top \begin{pmatrix} x \\ z_x \end{pmatrix} \right| \geq \epsilon_l \right\} \\ \cup \left\{ \text{Explore}_l^{(0)} = \text{True} \text{ and } \left| \hat{\omega}_l^{(0)} - \omega^* \right| \geq \epsilon_l \right\}.$$

Then, the regret decomposes as

$$R_T \leq \sum_{t \leq T} \mathbb{E}_{|\bar{\mathcal{F}}} [(x^* - x_t)^\top \gamma^*] + 2T\mathbb{P}[\mathcal{F}]. \tag{8}$$

The following lemma relies on concentration of Gaussian variables to bound the probability of the event \mathcal{F} .

Lemma 10. $\mathbb{P}(\mathcal{F}) \leq 2\delta$.

Now, the first term of (8) can be decomposed as

$$\sum_{t \leq T} (x^* - x_t)^\top \gamma^* \leq \sum_{z \in \{-1, 0, +1\}} \sum_{l=1}^{L^{(z)}+1} \sum_{t \in \text{Exp}_l^{(z)}} (x^* - x_t)^\top \gamma^* + \sum_{t \in \text{Recovery}} (x^* - x_t)^\top \gamma^*,$$

where we use as convention that the sum over an empty set is null. Note that for $z \in \{-1, +1\}$, during the phase $\text{Exp}_l^{(z)}$ the algorithm only samples actions from $\mathcal{X}_l^{(z)}$. By contrast, during the phase $\text{Exp}_l^{(0)}$, even actions eliminated from the sets $\mathcal{X}_l^{(z)}$ can be sampled. Finally, if the algorithm stops during phase $\text{Exp}_{L^{(0)}+1}^{(0)}$, but does not have enough budget to complete the last $\hat{\Delta}^l$ -optimal Exploration and Elimination Phase, it samples the remaining actions in the set $\mathcal{X}_{L^{(0)}+2}^{(-1)} \cup \mathcal{X}_{L^{(0)}+2}^{(+1)}$. Hence, the first term of (8) can be upper-bounded

by

$$\begin{aligned}
\sum_{t \leq T} (x^* - x_t)^\top \gamma^* &\leq \sum_{z \in \{-1, +1\}} \sum_{l=1}^{L_T} \left(\sum_{x \in \mathcal{X}_l^{(z)}} \mu_l^{(z)}(x) \right) \max_{x \in \mathcal{X}_l^{(z)}} (x^* - x)^\top \gamma^* \\
&+ \sum_{z \in \{-1, +1\}} \sum_{l=L_T+1}^{L^{(z)+1}} \sum_{t \in \text{Exp}_l^{(z)}} (x^* - x_t)^\top \gamma^* + \sum_{t \in \text{Recovery}} (x^* - x_t)^\top \gamma^* \\
&+ \sum_{l=1}^{L^{(0)}} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \Delta_x + \mathbb{1} \left\{ \text{Explore}_{L^{(0)+1}^{(0)}} = \text{False} \right\} \sum_{t \in \text{Exp}_{L^{(0)+1}^{(0)}}} \max_{x \in \mathcal{X}_{L^{(0)+2}^{(-1)} \cup \mathcal{X}_{L^{(0)+2}^{(+1)}}} (x^* - x)^\top \gamma^*.
\end{aligned} \tag{9}$$

We begin by bounding the sum of the regret corresponding to the Recovery phase and to the phases $\text{Exp}_L^{(z)}$ for $z \in \{-1, +1\}$ and $l > L_T$ on the event $\overline{\mathcal{F}}$.

Bound on $\sum_{z \in \{-1, +1\}} \sum_{l=L_T+1}^{L^{(z)+1}} \sum_{t \in \text{Exp}_l^{(z)}} (x^* - x_t)^\top \gamma^* + \sum_{t \in \text{Recovery}} (x^* - x_t)^\top \gamma^*.$

Lemma 11. *Let $x^* \in \arg \max_{x \in \mathcal{X}} x^\top \gamma^*$ be an optimal action. Then, on the event $\overline{\mathcal{F}}$ defined in Equation (7), for $l \geq 1$ such that $\text{Explore}_l^{(z_{x^*})} = \text{True}$,*

$$\mathcal{X}_{l+1}^{(z_{x^*})} \subset \left\{ x \in \mathcal{X}_1^{(z_{x^*})} : (x^* - x)^\top \gamma^* < 10\epsilon_{l+1} \right\}. \tag{10}$$

Moreover, for $l \geq 1$ such that $\text{Explore}_l^{(-z_{x^*})} = \text{True}$,

$$\mathcal{X}_{l+1}^{(-z_{x^*})} \subset \left\{ x \in \mathcal{X}_1^{(-z_{x^*})} : (x^* - x)^\top \gamma^* < 42\epsilon_{l+1} \right\}.$$

Recall that if $\text{Recovery} \neq \emptyset$, $L^{(0)} \geq L_T$. Then, all actions sampled during the Recovery phase belong to $\mathcal{X}_{l+1}^{(-1)} \cup \mathcal{X}_{l+1}^{(+1)}$ for some $l \geq L_T$. Lemma 11 shows that, on $\overline{\mathcal{F}}$, for $l \geq L_T$, the actions in $\mathcal{X}_{l+1}^{(z)}$ are sub-optimal by at most $42\epsilon_{L_T+1}$. Then, we get that on the event $\overline{\mathcal{F}}$,

$$\begin{aligned}
\sum_{z \in \{-1, +1\}} \sum_{l=L_T+1}^{L^{(z)+1}} \sum_{t \in \text{Exp}_l^{(z)}} (x^* - x_t)^\top \gamma^* + \sum_{t \in \text{Recovery}} (x^* - x_t)^\top \gamma^* &\leq T \times 42\epsilon_{L_T+1} \\
&\leq 53\kappa_*^{1/3} T^{2/3} \log(T)^{1/3}.
\end{aligned} \tag{11}$$

Bound on $\sum_{l=1}^{L^{(0)}} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \Delta_x + \mathbb{1} \left\{ \text{Explore}_{L^{(0)+1}^{(0)}} = \text{False} \right\} \sum_{t \in \text{Exp}_{L^{(0)+1}^{(0)}}} \max_{x \in \mathcal{X}_{L^{(0)+2}^{(-1)} \cup \mathcal{X}_{L^{(0)+2}^{(+1)}}} (x^* - x)^\top \gamma^*.$

We begin by bounding $\mathbb{1} \left\{ \text{Explore}_{L^{(0)+1}^{(0)}} = \text{False} \right\} \sum_{t \in \text{Exp}_{L^{(0)+1}^{(0)}}} \max_{x \in \mathcal{X}_{L^{(0)+2}^{(-1)} \cup \mathcal{X}_{L^{(0)+2}^{(+1)}}} (x^* - x)^\top \gamma^*$. Recall that $n_{L^{(0)+1}^{(0)}} =$

$\sum_{x \in \mathcal{X}} \mu_{L^{(0)+1}^{(0)}}^{(0)}(x)$ is the budget that would be necessary to complete the $\widehat{\Delta}^l$ -optimal Exploration and Elimination phase at phase $L^{(0)} + 1$. On the one hand, Lemma 11 implies that on the event $\overline{\mathcal{F}}$,

$$\mathbb{1} \left\{ \text{Explore}_{L^{(0)+1}^{(0)}} = \text{False} \right\} \sum_{t \in \text{Exp}_{L^{(0)+1}^{(0)}}} \max_{x \in \mathcal{X}_{L^{(0)+2}^{(-1)} \cup \mathcal{X}_{L^{(0)+2}^{(+1)}}} (x^* - x)^\top \gamma^* \leq 42n_{L^{(0)+1}^{(0)}} \epsilon_{L^{(0)+2} \leq 21n_{L^{(0)+1}^{(0)}} \epsilon_{L^{(0)+1}.$$

On the other hand, for all $l \leq L^{(0)} + 1$, the definition of $\widehat{\Delta}^l$ implies that $\widehat{\Delta}_x^l \geq \epsilon_l$ for all $x \in \mathcal{X}$. Therefore, $21n_{L^{(0)}+1}^{(0)}\epsilon_{L^{(0)}+1} \leq 21n_{L^{(0)}+1}^{(0)} \min_x \widehat{\Delta}_x^{L^{(0)}+1}$. This implies that on $\overline{\mathcal{F}}$,

$$\mathbb{1} \left\{ \text{Explore}_{L^{(0)}+1}^{(0)} = \text{False} \right\} \sum_{t \in \text{Exp}_{L^{(0)}+1}^{(0)}} \max_{x \in \mathcal{X}_{L^{(0)}+2}^{(-1)} \cup \mathcal{X}_{L^{(0)}+2}^{(+1)}} (x^* - x)^\top \gamma^* \leq 21 \sum_{x \in \mathcal{X}} \mu_{L^{(0)}+1}^{(0)}(x) \widehat{\Delta}_x^{L^{(0)}+1}. \quad (12)$$

Next, to bound the remaining terms of Equation (9), we bound the regret $\sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \Delta_x$ of exploration phase $\text{Exp}_l^{(0)}$ using the following lemma.

Lemma 12. *For all $l > 0$, and $z \in \{-1, +1\}$, we have*

$$\sum_{x \in \mathcal{X}_l^{(z)}} \mu_l^{(z)}(x) \leq \frac{2(d+1)}{\epsilon_l^2} \log \left(\frac{kl(l+1)}{\delta} \right) + \frac{(d+1)(d+2)}{2}.$$

and on $\overline{\mathcal{F}}$, we have

$$\sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \Delta_x \leq \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \widehat{\Delta}_x^l \leq \frac{2\kappa(\widehat{\Delta}^l)}{\epsilon_l^2} \log \left(\frac{l(l+1)}{\delta} \right) + 2(d+1).$$

Then, Equation (12) and Lemma 12 imply that on $\overline{\mathcal{F}}$

$$\begin{aligned} \sum_{l=1}^{L^{(0)}} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \Delta_x + \mathbb{1} \left\{ \text{Explore}_{L^{(0)}+1}^{(0)} = \text{False} \right\} \sum_{t \in \text{Exp}_{L^{(0)}+1}^{(0)}} \max_{x \in \mathcal{X}_{L^{(0)}+2}^{(-1)} \cup \mathcal{X}_{L^{(0)}+2}^{(+1)}} (x^* - x)^\top \gamma^* \\ \leq 21 \sum_{l=1}^{L^{(0)}+1} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \widehat{\Delta}_x^l \\ \leq 42 \sum_{l=1}^{L^{(0)}+1} \frac{\kappa(\widehat{\Delta}^l)}{\epsilon_l^2} \log \left(\frac{l(l+1)}{\delta} \right) + 42(d+1)(L^{(0)}+1) \end{aligned} \quad (13)$$

We rely on the following Lemma to bound $\kappa(\widehat{\Delta}^l)$.

Lemma 13. *On $\overline{\mathcal{F}}$, we have for any $l \geq 1$ and any $\tau > 0$*

$$\kappa(\widehat{\Delta}^l) \leq 513 \left(1 + \frac{\epsilon_l}{\tau} \right) \kappa(\Delta \vee \tau).$$

and

$$\kappa(\widehat{\Delta}^l) \geq \kappa(\Delta \vee \epsilon_l).$$

Lemma 12 and Lemma 13 with $\tau = \epsilon_{L^{(0)}}$ imply that on $\overline{\mathcal{F}}$,

$$\begin{aligned} \sum_{l=1}^{L^{(0)}+1} \frac{\kappa(\widehat{\Delta}^l)}{\epsilon_l^2} \log \left(\frac{l(l+1)}{\delta} \right) &\leq 513 \kappa(\Delta \vee \epsilon_{L^{(0)}}) \log \left(\frac{(L^{(0)}+1)(L^{(0)}+2)}{\delta} \right) \left(\sum_{l=1}^{L^{(0)}+1} \frac{1}{\epsilon_l^2} + \sum_{l=1}^{L^{(0)}+1} \frac{1}{\epsilon_l \epsilon_{L^{(0)}}} \right) \\ &\leq 513 \kappa(\Delta \vee \epsilon_{L^{(0)}}) \log \left(\frac{6L^{(0)}}{\delta} \right) \left(\frac{16}{\epsilon_{L^{(0)}}^2} + \frac{4}{\epsilon_{L^{(0)}}^2} \right) \\ &\leq 10260 \log \left(\frac{6L^{(0)}}{\delta} \right) \frac{\kappa(\widehat{\Delta}^{L^{(0)}})}{\epsilon_{L^{(0)}}^2} \end{aligned} \quad (14)$$

where the last line follows from the second claim of Lemma 13. Now, by definition of $L^{(0)}$, $\epsilon_{L^{(0)}} \geq (\kappa(\widehat{\Delta}^{L^{(0)}}) \log(T)/T)^{1/3}$. Then, Equation (14) implies that

$$\sum_{l=1}^{L^{(0)}+1} \frac{\kappa(\widehat{\Delta}^l)}{\epsilon_l^2} \log\left(\frac{l(l+1)}{\delta}\right) \leq 10260 \log\left(\frac{6L^{(0)}}{\delta}\right) \kappa(\widehat{\Delta}^{L^{(0)}})^{1/3} \log(T)^{-2/3} T^{2/3}. \quad (15)$$

Moreover, we observe that during each phase l , but the last one, we sample at least

$$\max_{z \in \{-1, 1\}} \sum_{x \in \mathcal{X}_l^{(z)}} \tau_{l,x}^{(z)} \geq \frac{2(d+1)}{\delta_l^2} \log(kl(l+1)/\delta)$$

actions during the G-optimal explorations, so the number of phases $L^{(0)}$ is never larger than

$$\ell_T = 1 \vee \log_4(T).$$

Using this remark, together with Equations (13) and (15), we find that on $\overline{\mathcal{F}}$

$$\begin{aligned} \sum_{l=1}^{L^{(0)}} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \widehat{\Delta}_x^l + \mathbb{1} \left\{ \text{Explore}_{L^{(0)}+1}^{(0)} = \text{False} \right\} \sum_{t \in \text{Exp}_{L^{(0)}+1}^{(0)}} \max_{x \in \mathcal{X}_{L^{(0)}+2}^{(-1)} \cup \mathcal{X}_{L^{(0)}+2}^{(+1)}} (x^* - x)^\top \gamma^* \\ \leq 2^{19} \log\left(\frac{6L^{(0)}}{\delta}\right) \kappa(\widehat{\Delta}^{L^{(0)}}) T^{2/3} \log(T)^{-2/3} + 42\ell_T. \end{aligned} \quad (16)$$

Bound on $\sum_{z \in \{-1, +1\}} \sum_{l=1}^{L_T} \left(\sum_{x \in \mathcal{X}_l^{(z)}} \mu_l^{(z)}(x) \right) \max_{x \in \mathcal{X}_l^{(z)}} (x^* - x)^\top \gamma^*$. We bound the remaining term in Equation (9) using the first claim in Lemma 12 and Lemma 11. On $\overline{\mathcal{F}}$,

$$\begin{aligned} \sum_{z \in \{-1, +1\}} \sum_{l=1}^{L_T} \left(\sum_{x \in \mathcal{X}_l^{(z)}} \mu_l^{(z)}(x) \right) \max_{x \in \mathcal{X}_l^{(z)}} (x^* - x)^\top \gamma^* &\leq 2 \sum_{l=1}^{L_T} \left(\frac{2(d+1)}{\epsilon_l^2} \log\left(\frac{kl(l+1)}{\delta}\right) + \frac{(d+1)(d+2)}{2} \right) 42\epsilon_l \\ &\leq \frac{336(d+1)}{\epsilon_{L_T}} \log\left(\frac{kL_T(1+L_T)}{\delta}\right) + 168(d+1)(d+2) \\ &\leq 267(d+1) \kappa_*^{-1/3} T^{1/3} \log(T)^{-1/3} \log\left(\frac{kL_T(1+L_T)}{\delta}\right) \\ &\quad + 168(d+1)(d+2). \end{aligned} \quad (17)$$

Combing Equations (8), (9), (11), (16), and (17), and using $\delta = T^{-1}$, $\kappa(\widehat{\Delta}^{L^{(0)}}) \leq \kappa_*$ and $L_T \leq 4T/\log(2)$, we get for all $T \geq 1$

$$R_T \leq C \left(\kappa_*^{1/3} T^{2/3} \log(T)^{1/3} + (d \vee \kappa_*) \log(T) + d^2 + d \kappa_*^{-1/3} T^{1/3} \log(kT) \log(T)^{-1/3} \right)$$

for some absolute constant $C > 0$. Finally, for

$$T \geq \frac{((d \vee \kappa_*)^{3/2} \log(T)) \vee d^3}{\sqrt{\kappa_*}} \vee \frac{(d \log(kT))^3}{(\kappa_* \log(T))^2},$$

we get

$$R_T \leq C' \kappa_*^{1/3} T^{2/3} \log(T)^{1/3}.$$

C.2 Proof of Theorem 2

The beginning of the proof of Theorem 2 follows the same lines as the proof of Theorem 1. We begin by decomposing the regret as

$$R_T \leq \sum_{t \leq T} \mathbb{E}_{|\overline{\mathcal{F}}} [(x^* - x_t)^\top \gamma^*] + 2T\mathbb{P}[\mathcal{F}]. \quad (18)$$

where \mathcal{F} is defined in Equation (7). On the one hand, Lemma 10 implies $T\mathbb{P}[\mathcal{F}] \leq 2\delta T$. Then, Equation (18) implies

$$\begin{aligned} R_T \leq & 4\delta T + \mathbb{E}_{|\overline{\mathcal{F}}} \left[\sum_{z \in \{-1, +1\}} \sum_{l \geq 1}^{L^{(z)}+1} \sum_{t \in \text{Exp}_l^{(z)}} (x^* - x_t)^\top \gamma^* \right] + \mathbb{E}_{|\overline{\mathcal{F}}} \left[\sum_{t \in \text{Recovery}} (x^* - x_t)^\top \gamma^* \right] \\ & + \mathbb{E}_{|\overline{\mathcal{F}}} \left[\sum_{l=1}^{L^{(0)}} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \Delta_x \right] + \mathbb{E}_{|\overline{\mathcal{F}}} \left[\mathbf{1} \left\{ \text{Explore}_{L^{(0)}+1}^{(0)} = \text{False} \right\} \sum_{t \in \text{Exp}_{L^{(0)}+1}^{(0)}} \max_{x \in \mathcal{X}_{L^{(0)}+2}^{(-1)} \cup \mathcal{X}_{L^{(0)}+2}^{(+1)}} (x^* - x)^\top \gamma^* \right] \end{aligned} \quad (19)$$

where \mathcal{F} is defined in Equation (7), and where we used the convention that the sum over an empty set is null.

Bound on $\mathbf{1} \left\{ \text{Explore}_{L^{(0)}+1}^{(0)} = \text{False} \right\} \sum_{t \in \text{Exp}_{L^{(0)}+1}^{(0)}} \max_{x \in \mathcal{X}_{L^{(0)}+1}^{(z)}} (x^* - x)^\top \gamma^*.$

Similarly to the proof of Theorem 1, we use Lemma 11 and Lemma 13 to show that on $\overline{\mathcal{F}}$

$$\mathbf{1} \left\{ \text{Explore}_{L^{(0)}+1}^{(0)} = \text{False} \right\} \sum_{t \in \text{Exp}_{L^{(0)}+1}^{(0)}} \max_{x \in \mathcal{X}_{L^{(0)}+1}^{(z)}} (x^* - x)^\top \gamma^* \leq 21 \sum_{x \in \mathcal{X}} \mu_{L^{(0)}+1}^{(0)}(x) \widehat{\Delta}_x^{L^{(0)}+1}. \quad (20)$$

Bound on $\sum_{z \in \{-1, +1\}} \sum_{l \geq 1}^{L^{(z)}+1} \sum_{t \in \text{Exp}_l^{(z)}} (x^* - x_t)^\top \gamma^*.$

Lemma 11 shows that for $l \leq L^{(z)}$, the actions in $\mathcal{X}_{l+1}^{(z)}$ are sub-optimal by at most an additional factor at most $21\epsilon_l$. Let us set $l_{\Delta_{\min}} = \lceil -\log_2(\Delta_{\min}/21) \rceil$, so that

$$\frac{\Delta_{\min}}{42} \leq \epsilon_{l_{\Delta_{\min}}} \leq \frac{\Delta_{\min}}{21}.$$

For $l \geq l_{\Delta_{\min}}$, we have $\mathcal{X}_{l+1}^{(-1)} \cup \mathcal{X}_{l+1}^{(+1)} = \{x_{z^*}\}$. Thus, $l^{(-z_{x^*})} \leq l_{\Delta_{\min}}$, and for $l \geq l_{\Delta_{\min}}$, the algorithm selects only x^* during the phase $\text{Exp}_l^{(z^*)}$. Then, combining Lemmas 12 and 11, and the fact that $L^{(z)} + 1 \leq \ell_T$, we find that, on $\overline{\mathcal{F}}$,

$$\begin{aligned} \sum_{z \in \{-1, +1\}} \sum_{l=1}^{L^{(z)}+1} \sum_{t \in \text{Exp}_l^{(z)}} (x^* - x_t)^\top \gamma^* & \leq \sum_{z \in \{-1, +1\}} \sum_{l=1}^{l_{\Delta_{\min}}+1 \wedge \ell_T} \left(\sum_{x \in \mathcal{X}_l^{(z)}} \mu_l^{(z)}(x) \right) \max_{x \in \mathcal{X}_l^{(z)}} (x^* - x)^\top \gamma^* \\ & \leq 2 \sum_{l=1}^{l_{\Delta_{\min}}+1 \wedge \ell_T} \left(\frac{2(d+1)}{\epsilon_l^2} \log \left(\frac{kl(l+1)}{\delta} \right) + \frac{(d+1)(d+2)}{2} \right) 42\epsilon_l \\ & \leq 84(d+1)(d+2) + \epsilon_{l_{\Delta_{\min}}}^{-1} \times 672(d+1) \log \left(\frac{k(1+\ell_T)(2+\ell_T)}{\delta} \right) \\ & \leq 84(d+1)(d+2) + \frac{28224(d+1)}{\Delta_{\min}} \log \left(\frac{k(1+\ell_T)(2+\ell_T)}{\delta} \right). \end{aligned} \quad (21)$$

Bound on
$$\sum_{t \in \text{Recovery}} (x^* - x_t)^\top \gamma^* + \sum_{l=1}^{L^{(0)}} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \Delta_x + \sum_{x \in \mathcal{X}} \mu_{L^{(0)}+1}^{(0)}(x) \widehat{\Delta}_x^{L^{(0)}+1}.$$

We use the following lemma to bound the number of phases necessary to eliminate the sub-optimal group.

Lemma 14. *On the event $\overline{\mathcal{F}}$ defined in Equation (7), for $l \geq 1$ such that $\epsilon_l \leq \frac{\Delta_{\neq}}{8}$ and $\text{Explore}_L^{(0)} = \text{True}$, $\widehat{z}_{l+1}^* = z_{x^*}$.*

Let $l_{\Delta_{\neq}} = \lceil -\log(\Delta_{\neq}/8)/\log(2) \rceil$ be such that

$$\frac{\Delta_{\neq}}{16} \leq \epsilon_{l_{\Delta_{\neq}}} \leq \frac{\Delta_{\neq}}{8}. \quad (22)$$

Lemma 14 implies that on $\overline{\mathcal{F}}$, $L^{(0)} \leq l_{\Delta_{\neq}}$.

To bound the remaining terms, we consider two cases, corresponding to $\text{Recovery} = \emptyset$ and $\text{Recovery} \neq \emptyset$.

Case 1: $\text{Recovery} = \emptyset$. Our case assumption implies that

$$\sum_{t \in \text{Recovery}} (x^* - x_t)^\top \gamma^* = 0. \quad (23)$$

Lemma 13 implies that

$$\sum_{l=1}^{L^{(0)}} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \Delta_x + \sum_{x \in \mathcal{X}} \mu_{L^{(0)}+1}^{(0)}(x) \widehat{\Delta}_x^{L^{(0)}+1} \leq \sum_{l=1}^{L^{(0)}+1} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \widehat{\Delta}_x^l.$$

Moreover, $L^{(0)} \leq l_{\Delta_{\neq}} \wedge \ell_T$, so on $\overline{\mathcal{F}}$

$$\sum_{l=1}^{L^{(0)}+1} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \widehat{\Delta}_x^l \leq \sum_{l=1}^{(l_{\Delta_{\neq}} \wedge \ell_T)+1} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \widehat{\Delta}_x^l.$$

Using Lemma 12, we find that on $\overline{\mathcal{F}}$

$$\begin{aligned} \sum_{l=1}^{(l_{\Delta_{\neq}} \wedge \ell_T)+1} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \widehat{\Delta}_x^l &\leq \sum_{l=1}^{(l_{\Delta_{\neq}} \wedge \ell_T)+1} \frac{2\kappa(\widehat{\Delta}^l)}{\epsilon_l^2} \log\left(\frac{l(l+1)}{\delta}\right) + 2(d+1)(\ell_T+1) \\ &\leq 2 \log\left(\frac{(\ell_T+1)(\ell_T+2)}{\delta}\right) \sum_{l=1}^{l_{\Delta_{\neq}}+1} \frac{\kappa(\widehat{\Delta}^l)}{\epsilon_l^2} + 2(d+1)(\ell_T+1). \end{aligned}$$

Using Lemma 13 with $\tau = \Delta_{\neq}$ and (22), we have on $\overline{\mathcal{F}}$

$$\begin{aligned} \sum_{l=1}^{l_{\Delta_{\neq}}+1} \frac{\kappa(\widehat{\Delta}^l)}{\epsilon_l^2} &\leq 513\kappa(\Delta \vee \Delta_{\neq}) \sum_{l=1}^{l_{\Delta_{\neq}}+1} (\epsilon_l^{-2} + \epsilon_l^{-1}/\Delta_{\neq}) \\ &\leq \frac{2^{18}\kappa(\Delta \vee \Delta_{\neq})}{\Delta_{\neq}^2}. \end{aligned}$$

We obtain on $\overline{\mathcal{F}}$

$$\sum_{l=1}^{L^{(0)}+1} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \widehat{\Delta}_x^l \leq 2^{19} \log\left(\frac{(\ell_T+1)(\ell_T+2)}{\delta}\right) \frac{\kappa(\Delta \vee \Delta_{\neq})}{\Delta_{\neq}^2} + 2(d+1)(\ell_T+1). \quad (24)$$

Combining Equations (21), (20), (23), and (24), we find that on $\overline{\mathcal{F}}$, when $\text{Recovery} = \emptyset$, there exists an absolute constant $c > 0$ such that for $\delta = T^{-1}$,

$$\begin{aligned} & \sum_{z \in \{-1, +1\}} \sum_{l \geq 1}^{L^{(z)}+1} \sum_{t \in \text{Exp}_l^{(z)}} (x^* - x_t)^\top \gamma^* + \sum_{t \in \text{Recovery}} (x^* - x_t)^\top \gamma^* + \sum_{l=1}^{L^{(0)}} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \Delta_x \\ & + \mathbb{1}\{\text{Explore}_{L^{(0)}+1}^{(0)} = \text{False}\} \sum_{t \in \text{Exp}_{L^{(0)}+1}^{(0)}} \max_{x \in \mathcal{X}_{L^{(0)}+2}^{(-1)} \cup \mathcal{X}_{L^{(0)}+2}^{(+1)}} (x^* - x)^\top \gamma^* \\ & \leq c \left(d^2 + \left(\frac{d}{\Delta_{\min}} \vee \frac{\kappa(\Delta \vee \Delta_{\neq})}{\Delta_{\neq}^2} \right) \log(T) + \frac{d}{\Delta_{\min}} \log(k) \right). \end{aligned} \quad (25)$$

Case 2: $\text{Recovery} \neq \emptyset$. In this case, the algorithm enters Recovery at phase $L^{(0)}$, so $\text{Explore}_{L^{(0)}+1}^{(0)} = \text{False}$ and $\text{Exp}_{L^{(0)}+1}^{(0)} = \emptyset$, and

$$\mathbb{1}\{\text{Explore}_{L^{(0)}+1}^{(0)} = \text{False}\} \sum_{t \in \text{Exp}_{L^{(0)}+1}^{(0)}} \max_{x \in \mathcal{X}_{L^{(0)}+2}^{(-1)} \cup \mathcal{X}_{L^{(0)}+2}^{(+1)}} (x^* - x)^\top \gamma^* = 0. \quad (26)$$

Using Lemma 11, we see that

$$\sum_{t \in \text{Recovery}} (x^* - x_t)^\top \gamma^* \leq 21T \epsilon_{L^{(0)}+1}.$$

On the other hand, in the Recovery phase, $\epsilon_{L^{(0)}+1} \leq \left(\kappa(\widehat{\Delta}^{L^{(0)}+1}) \log(T)/T \right)^{1/3}$. Thus,

$$\sum_{t \in \text{Recovery}} (x^* - x_t)^\top \gamma^* \leq \frac{21\kappa(\widehat{\Delta}^{L^{(0)}+1}) \log(T)}{\epsilon_{L^{(0)}+1}^2}.$$

Now, Lemma 12 show that

$$\sum_{l=1}^{L^{(0)}} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \Delta_x \leq 4 \log(2L^{(0)}\delta^{-1}) \sum_{l=1}^{L^{(0)}} \frac{\kappa(\widehat{\Delta}^l)}{\epsilon_l^2} + 4dL^{(0)}.$$

Combining these results, and using $L^{(0)} \leq \ell_T$, we see that

$$\sum_{t \in \text{Recovery}} (x^* - x_t)^\top \gamma^* + \sum_{l=1}^{L^{(0)}} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \Delta_x \leq 4dL^{(0)} + (4 \log(2\ell_T\delta^{-1}) \vee 21 \log(T)) \sum_{l=1}^{L^{(0)}+1} \frac{\kappa(\widehat{\Delta}^l)}{\epsilon_l^2}. \quad (27)$$

Using Lemma 13 with $\tau = \epsilon_{L^{(0)}}$, we see that

$$\begin{aligned} \sum_{l=1}^{L^{(0)}+1} \frac{\kappa(\widehat{\Delta}^l)}{\epsilon_l^2} & \leq 513 \sum_{l=1}^{L^{(0)}+1} \frac{\kappa(\Delta \vee \epsilon_{L^{(0)}})}{\epsilon_l^2} + 513 \sum_{l=1}^{L^{(0)}+1} \frac{\kappa(\Delta \vee \epsilon_{L^{(0)}})}{\epsilon_{L^{(0)}} \epsilon_l} \\ & \leq 10260 \frac{\kappa(\Delta \vee \epsilon_{L^{(0)}})}{\epsilon_{L^{(0)}}^2}. \end{aligned}$$

Now, the algorithm enters the Recovery phase before finding the best group, so we must have $L^{(0)} \leq l_{\Delta_{\neq}}$. This implies that

$$\sum_{l=1}^{L^{(0)}+1} \frac{\kappa(\widehat{\Delta}^l)}{\epsilon_l^2} \leq 2^{18} \frac{\kappa(\Delta \vee \epsilon_{L^{(0)}})}{\Delta_{\neq}^2}.$$

Finally, note that $L^{(0)} \geq L_T$, so $\epsilon_{L^{(0)}} \leq \epsilon_{L_T} = \epsilon_T$, and

$$\sum_{l=1}^{L^{(0)}+1} \frac{\kappa(\widehat{\Delta}^l)}{\epsilon_l^2} \leq 2^{18} \frac{\kappa(\Delta \vee \epsilon_T)}{\Delta_{\neq}^2}. \quad (28)$$

Combining Equations (21), (26), (27), and (28), we find that on $\overline{\mathcal{F}}$, when $\text{Recovery} \neq \emptyset$, there exists an absolute constant $c > 0$ such that for $\delta = T^{-1}$,

$$\begin{aligned} & \sum_{z \in \{-1, +1\}} \sum_{l \geq 1}^{L^{(z)}+1} \sum_{t \in \text{Exp}_l^{(z)}} (x^* - x_t)^\top \gamma^* + \sum_{t \in \text{Recovery}} (x^* - x_t)^\top \gamma^* + \sum_{l=1}^{L^{(0)}} \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \Delta_x \\ & + \mathbf{1}\{\text{Explore}_{L^{(0)}+1}^{(0)} = \text{False}\} \sum_{t \in \text{Exp}_{L^{(0)}+1}^{(0)}} \max_{x \in \mathcal{X}_{L^{(0)}+2}^{(-1)} \cup \mathcal{X}_{L^{(0)}+2}^{(+1)}} (x^* - x)^\top \gamma^* \\ & \leq c \left(d^2 + \left(\frac{d}{\Delta_{\min}} \vee \frac{\kappa(\Delta \vee \epsilon_T)}{\Delta_{\neq}^2} \right) \log(T) + \frac{d \log(k)}{\Delta_{\min}} \right). \end{aligned} \quad (29)$$

Conclusion We conclude the proof of Theorem 2 by combining Equations (19), (25) and (29).

C.3 Proof of Theorem 3

Consider the actions \mathcal{A} defined in the following lemma.

Lemma 15. *Let the action set be given by $\mathcal{A} = \left\{ \begin{pmatrix} x_1 \\ z_{x_1} \end{pmatrix}, \dots, \begin{pmatrix} x_{d+1} \\ z_{x_{d+1}} \end{pmatrix} \right\}$, where $\begin{pmatrix} x_1 \\ z_{x_1} \end{pmatrix} = e_1 + e_{d+1}$, $\begin{pmatrix} x_i \\ z_{x_i} \end{pmatrix} = e_i - e_{d+1}$ for $i \in \{2, \dots, d\}$, and $\begin{pmatrix} x_{d+1} \\ z_{x_{d+1}} \end{pmatrix} = -\left(1 - \frac{2}{\sqrt{\kappa_*}+1}\right) e_1 - e_{d+1}$. It holds that*

$$\min_{\pi \in \mathcal{P}^{\mathcal{A}}} \left\{ e_{d+1}^\top \left(\sum_{\begin{pmatrix} x \\ z \end{pmatrix} \in \mathcal{A}} \pi_x \begin{pmatrix} x \\ z_x \end{pmatrix} \begin{pmatrix} x \\ z_x \end{pmatrix}^\top \right)^+ e_{d+1} \right\} = \kappa.$$

By Lemma 15, $\mathcal{A} \in \mathbf{A}_{\kappa_*, d}$. We will introduce two bandit problems characterized by two parameters $\theta_T^{(1)}$ and $\theta_T^{(2)}$ - assuming that the noise ξ_t is Gaussian and i.i.d. - and we prove that for any algorithm, the regret for one of those two problems must be of larger order than $\kappa_*^{1/3} T^{2/3}$.

We also consider the following two alternative problems. For a small $1/4 > \rho_T > 0$ where $\rho_T = T^{-1/3} \kappa_*^{1/3}$ (satisfied since $T > 4^3 \kappa_*$), the two alternative action parameters are defined as:

$$\begin{aligned} \gamma_T^{(1)} &= \frac{1 + \rho_T}{2} e_1 + \frac{1 - \rho_T}{2} e_2 - \frac{\rho_T}{2} \left(\sum_{3 \leq j \leq d} e_j \right) \\ \gamma_T^{(2)} &= \frac{1 - \rho_T}{2} e_1 + \frac{1 + \rho_T}{2} e_2 + \frac{\rho_T}{2} \left(\sum_{3 \leq j \leq d} e_j \right). \end{aligned}$$

On top of this, two bias parameters are defined as $\omega_T^{(1)} = -\frac{\rho_T}{2}$ and $\omega_T^{(2)} = \frac{\rho_T}{2}$. Through this, we define the two bandit problems of the sketch of proof of Lemma 15 characterized by $\theta_T^{(1)} = \begin{pmatrix} \gamma_T^{(1)} \\ \omega_T^{(1)} \end{pmatrix}$ and $\theta_T^{(2)} = \begin{pmatrix} \gamma_T^{(2)} \\ \omega_T^{(2)} \end{pmatrix}$ - and where the distribution of the noise ξ_t is supposed to be Gaussian and i.i.d. We refer to these two problems respectively as **Problem 1** and **Problem 2**. We write $R_T^{(1)}$, $\mathbb{P}^{(1)}$ and $\mathbb{E}^{(1)}$ (respectively $R_T^{(2)}$, $\mathbb{P}^{(2)}$)

and $\mathbb{E}^{(2)}$) for the regret, probability and expectation for the first bandit problem, when the parameter is $\theta_T^{(1)}$ (respectively the second bandit problem with $\theta_T^{(2)}$). We also write $\mathbb{P}_j^{(i)}$ for the distribution of a sample received in **Problem i** when sampling action x_j at any given time t - note that by definition of the bandit problems, this distribution does not depend on t and on the past samples given that action x_j is sampled.

The three following facts hold on these two bandit problems:

Fact 1 The parameters $\gamma_T^{(1)}$ and $\gamma_T^{(2)}$ are chosen so that x_1 is the unique best action for **Problem 1**, and x_2 is the unique best action for **Problem 2**. Choosing any sub-optimal action induces an instantaneous regret of at least ρ_T , and choosing the very sub-optimal action x_{d+1} induces an instantaneous regret of at least $1/2$.

Fact 2 Because of the chosen bias parameters, the distributions of the evaluations of all actions but x_{d+1} are exactly the same under the two bandit problems characterized by $\theta^{(1)}$ and $\theta^{(2)}$ - i.e. exactly the same data is observed under the two alternative bandit problems defined by the two alternative parameters for all actions but x_{d+1} . More precisely, for $i \in \{1, 2\}$, in **Problem i** and at any time t , when sampling action x_i where $i \leq 2$, we observe a sample distributed according to $\mathcal{N}(1/2, 1)$ - i.e. $\mathbb{P}_j^{(i)}$ is $\mathcal{N}(1/2, 1)$ - and when sampling action x_i where $2 < i \leq d+1$, we observe a sample distributed according to $\mathcal{N}(0, 1)$ - i.e. $\mathbb{P}_j^{(i)}$ is $\mathcal{N}(0, 1)$.

Fact 3 The distributions of the outcomes of the evaluation of action x_{d+1} differs in the two bandit problems. Set $\alpha = 2/(\sqrt{\kappa_*} + 1)$. In **Problem 1**, $\mathbb{P}_{d+1}^{(1)}$ is $\mathcal{N}(-\frac{1-\alpha-\rho_T\alpha}{2}, 1)$. In **Problem 2**, $\mathbb{P}_{d+1}^{(2)}$ is $\mathcal{N}(-\frac{1-\alpha+\rho_T\alpha}{2}, 1)$. So that the difference between the means of the evaluations of action x_{d+1} in the two bandit problems is $\Delta = \rho_T\alpha = \frac{2\rho_T}{\sqrt{\kappa_*}+1} \leq \frac{2\rho_T}{\sqrt{\kappa_*}}$.

For $i \leq d+1$, we write $N_i(T)$ for the number of times that action x_i has been selected before time T . In **Problem 1**, choosing the action x_{d+1} leads to an instantaneous regret larger than $\frac{1}{2}$ (**Fact 1**), so that

$$R_T^{(1)} \geq \frac{\mathbb{E}^{(1)} [N_{x_{d+1}}(T)]}{2}.$$

If $\mathbb{E}^{(1)} [N_{d+1}(T)] \geq \frac{T^{2/3}\kappa_*^{1/3}}{2}$, then Theorem 1 follows immediately; we therefore consider from now on the case when

$$\mathbb{E}^{(1)} [N_{d+1}(T)] \leq \frac{T^{2/3}\kappa_*^{1/3}}{2}. \quad (30)$$

Now, let us define the event

$$F = \left\{ N_1(T) \geq \frac{T}{2}\kappa_*^{1/3} \right\}.$$

Note that action x_1 is optimal for **Problem 1** and that action x_2 is optimal for **Problem 2** (**Fact 1**). Since choosing an action that is sub-optimal leads to an instantaneous regret larger than ρ_T (**Fact 1**), we also have

$$R_T^{(1)} \geq \frac{T\rho_T}{2} \mathbb{P}^{(1)}(\bar{F})$$

and

$$R_T^{(2)} \geq \frac{T\rho_T}{2} \mathbb{P}^{(2)}(F).$$

Then, Bretagnolle-Huber inequality (see, e.g., Theorem 14.2 in [21]) implies that

$$R_T^{(1)} + R_T^{(2)} \geq \frac{T\rho_T}{4} \exp\left(-KL\left(\mathbb{P}^{(1)}, \mathbb{P}^{(2)}\right)\right).$$

For the choice $\rho_T = T^{-1/3}\kappa_*^{1/3}$, this implies that

$$R_T^{(1)} + R_T^{(2)} \geq \frac{T^{2/3}\kappa_*^{1/3}}{4} \exp\left(-KL\left(\mathbb{P}^{(1)}, \mathbb{P}^{(2)}\right)\right). \quad (31)$$

Now, the Kullback-Leibler divergence between $\mathbb{P}^{(1)}$ and $\mathbb{P}^{(2)}$ can be rewritten as follows (see, e.g., Lemma 15.1 in [21]) :

$$KL(\mathbb{P}^{(1)}, \mathbb{P}^{(2)}) = \frac{1}{2} \sum_{j \leq d+1} \mathbb{E}^{(1)} [N_j(T)] KL(\mathbb{P}_j^{(1)}, \mathbb{P}_j^{(2)}).$$

By **Fact 2**, we have that for any $j \leq d$, $\mathbb{P}_j^{(1)} = \mathbb{P}_j^{(2)}$. So that

$$KL(\mathbb{P}^{(1)}, \mathbb{P}^{(2)}) = \frac{1}{2} \mathbb{E}^{(1)} [N_{d+1}(T)] KL(\mathbb{P}_{d+1}^{(1)}, \mathbb{P}_{d+1}^{(2)}).$$

By the characterization of $\mathbb{P}_{d+1}^{(1)}, \mathbb{P}_{d+1}^{(2)}$ in **Fact 3**, and recalling that the Kullback-Leibler divergence between two normalized Gaussian distributions is given by the squared distance between their means, we find that

$$KL(\mathbb{P}^{(1)}, \mathbb{P}^{(2)}) = \frac{1}{2} \mathbb{E}^{(1)} [N_{d+1}(T)] \bar{\Delta}^2.$$

Thus, by the definition of $\bar{\Delta}$ in **Fact 3** and by Equation (30)

$$KL(\mathbb{P}^{(1)}, \mathbb{P}^{(2)}) = \frac{1}{2} \mathbb{E}^{(1)} [N_{d+1}(T)] \left(\frac{2\rho_T}{\sqrt{\kappa_* + 1}} \right)^2 \leq \frac{T^{2/3} \kappa_*^{1/3}}{4} \times \frac{4\rho_T^2}{\kappa_*} = 1, \quad (32)$$

reminding that $\rho_T = T^{-1/3} \kappa_*^{1/3}$.

Combining Equations (31) and (32) implies that

$$\max \{ R_T^{(1)}, R_T^{(2)} \} \geq \frac{T^{2/3} \kappa_*^{1/3}}{8} \exp(-1),$$

which concludes the proof of Theorem 3.

C.4 Proof of Theorems 4

Theorems 4 follows directly from the next Theorem.

Theorem 6. *For all $\kappa_* \geq 1$ and all $d \geq 4$, there exists an action set $\mathcal{A} \in \mathbf{A}_{\kappa_*, d}$, such that for all bandit algorithms, for all $(\Delta_{\min}, \Delta_{\neq}) \in (0, 1/8)^2$ with $\Delta_{\min} \leq \Delta_{\neq}$, and for all budget $T \geq 2$, there exists a problem characterized by $\theta \in \Theta_{\Delta_{\min}, \Delta_{\neq}}^{\mathcal{A}}$ such that the regret of the algorithm on the problem satisfies*

$$\begin{aligned} R_T^\theta &\geq \left[\frac{d}{10\Delta_{\min}} \log(T) \left[1 - \frac{\log\left(\frac{8d \log(T)}{\Delta_{\min}^2}\right)}{\log(T)} \right] \right] \vee \left[\frac{\kappa_* + 1}{4\Delta_{\neq}^2} \log(T) \left[1 - \frac{\log\left(\frac{8\kappa_* \log(T)}{\Delta_{\neq}^3}\right)}{\log(T)} \right] \right] \\ &\vee \left[\frac{\kappa_*}{4\Delta_{\neq}^2} \left[1 \wedge \log\left(\frac{T\Delta_{\neq}^3}{8\kappa_*}\right) \right] \right]. \end{aligned} \quad (33)$$

Moreover, on this problem, $\kappa(\Delta) \in [\kappa_*/8, 2\kappa_*]$.

Remark 1. *Note that Theorem 6 allows us to recover a lower bound similar to that of Theorem 3 by choosing Δ_{\neq} and Δ_{\min} of the order $\kappa_*^{1/3} T^{-1/3}$, however this bound only holds for d larger than 4.*

We prove Theorem 6 for the following set of actions \mathcal{A} : $\mathcal{A} = \left\{ \begin{pmatrix} x_1 \\ z_{x_1} \end{pmatrix}, \dots, \begin{pmatrix} x_{d+1} \\ z_{x_{d+1}} \end{pmatrix} \right\}$, where $\begin{pmatrix} x_i \\ z_{x_i} \end{pmatrix} = e_i + e_{d+1}$, for $i \in \{2, \dots, \lfloor d/2 \rfloor\}$, $\begin{pmatrix} x_i \\ z_{x_i} \end{pmatrix} = e_i - e_{d+1}$ for $i \in \{\lfloor d/2 \rfloor + 1, \dots, d\}$, and $\begin{pmatrix} x_{d+1} \\ z_{x_{d+1}} \end{pmatrix} = -\left(1 - \frac{2}{\sqrt{\kappa_* + 1}}\right) e_1 - e_{d+1}$. Then, by Lemma 16, for this choice of action set, we have $\mathcal{A} \in \mathbf{A}_{\kappa_*, d}$.

We consider the following set of bandit problems: for $i \in \{1, \dots, \lfloor d/2 \rfloor + 1\}$ **Problem i** is characterized by the parameter $\theta^{(i)}$, where $\theta^{(i)} = \begin{pmatrix} \gamma^{(i)} \\ \omega^{(i)} \end{pmatrix}$ is defined as:

$$\begin{aligned} \gamma^{(1)} &= \frac{1 + \Delta_{\neq} - \Delta_{\min}}{2} \left(\sum_{1 \leq j \leq \lfloor d/2 \rfloor} e_j \right) + \frac{1 - \Delta_{\neq} - \Delta_{\min}}{2} \left(\sum_{\lfloor d/2 \rfloor + 1 \leq j \leq d} e_j \right) + \Delta_{\min} e_1 + \Delta_{\min} e_{\lfloor d/2 \rfloor + 1} \\ \gamma^{(i)} &= \gamma^{(1)} + 2\Delta_{\min} e_i + 2\Delta_{\min} e_{\lfloor d/2 \rfloor + i} \quad \forall i \in \{2, \dots, \lfloor d/2 \rfloor\} \\ \gamma^{(\lfloor d/2 \rfloor + 1)} &= \frac{1 - \Delta_{\neq} - \Delta_{\min}}{2} \left(\sum_{1 \leq j \leq \lfloor d/2 \rfloor} e_j \right) + \frac{1 + \Delta_{\neq} - \Delta_{\min}}{2} \left(\sum_{\lfloor d/2 \rfloor + 1 \leq j \leq d} e_j \right) + \Delta_{\min} e_1 + \Delta_{\min} e_{\lfloor d/2 \rfloor + 1}, \end{aligned}$$

and the bias parameters are defined as $\omega^{(i)} = -\frac{\Delta_{\neq}}{2} \quad \forall i \in \{1, \dots, \lfloor d/2 \rfloor\}$, and otherwise $\omega^{(\lfloor d/2 \rfloor + 1)} = \frac{\Delta_{\neq}}{2}$. We write $\mathbb{E}^{(i)}, \mathbb{P}^{(i)}, R_T^{(i)}$ for resp. the probability, expectation, and regret, in **Problem i**. Note that this choice of parameters ensures that $\forall i \in \{1, \dots, \lfloor d/2 \rfloor + 1\}$, $\theta^{(i)} \in \Theta_{\Delta_{\min}, \Delta_{\neq}}^A$.

Set $\mathcal{A} = \left\{ \begin{pmatrix} x_1 \\ z_{x_1} \end{pmatrix}, \dots, \begin{pmatrix} x_{d+1} \\ z_{x_{d+1}} \end{pmatrix} \right\}$, where $\begin{pmatrix} x_i \\ z_{x_i} \end{pmatrix} = e_i + e_{d+1}$, for $i \in \{2, \dots, \lfloor d/2 \rfloor\}$, $\begin{pmatrix} x_i \\ z_{x_i} \end{pmatrix} = e_i - e_{d+1}$ for $i \in \{\lfloor d/2 \rfloor + 1, \dots, d\}$, and $\begin{pmatrix} x_{d+1} \\ z_{x_{d+1}} \end{pmatrix} = -\left(1 - \frac{2}{\sqrt{\kappa_*} + 1}\right) e_1 - e_{d+1}$. Then, Lemma 16 shows that $\mathcal{A} \in \mathbf{A}_{\kappa_*, d}$.

Lemma 16. *It holds that*

$$\min_{\pi \in \mathcal{P}_{e_{d+1}}^A} \left\{ e_{d+1}^\top \left(\sum_{\begin{pmatrix} x \\ z_x \end{pmatrix} \in \mathcal{A}} \pi(x) \begin{pmatrix} x \\ z_x \end{pmatrix} \begin{pmatrix} x \\ z_x \end{pmatrix}^\top \right)^+ e_{d+1} \right\} = \kappa_*.$$

The following facts hold:

- Fact 1** For any $i \in \{1, \dots, \lfloor d/2 \rfloor + 1\}$, action x_i is the unique optimal action in **Problem i**. Since $1/2 \geq \Delta_{\neq} \geq \Delta_{\min}$, sampling any other (sub-optimal) action leads to an instantaneous regret of at least Δ_{\min} . Moreover, choosing an action in the group $-z_i$ leads to an instantaneous regret of at least Δ_{\neq} .
- Fact 2** In **Problem i** for any $i \in \{1, \dots, \lfloor d/2 \rfloor + 1\}$, action $d + 1$ is very sub-optimal and sampling it leads to an instantaneous regret higher than $(1 - 2/(\sqrt{\kappa_*} + 1))(1 - \Delta_{\neq} + \Delta_{\min}) + (1 + \Delta_{\neq} + \Delta_{\min})/2 \geq 1/2$, since $\kappa_* \geq 1$ and $1/2 \geq \Delta_{\neq} \geq \Delta_{\min}$.
- Fact 3** In **Problem i**, for $i \in \{1, \dots, \lfloor d/2 \rfloor + 1\}$, when sampling action x_j at time, t the distribution of the observation does not depend on t or on the past (except through the choice of x_j) and is $\mathbb{P}_j^{(i)}$. It is characterized as:

$$\begin{aligned} \forall i \in \{1, \dots, \lfloor d/2 \rfloor + 1\}, \mathbb{P}_1^{(i)}, \mathbb{P}_{\lfloor d/2 \rfloor + 1}^{(i)} &\text{ are } \mathcal{N}((1 + \Delta_{\min})/2, 1) \\ \forall i \in \{1, \dots, \lfloor d/2 \rfloor + 1\}, \forall j \in \{2, \dots, d\} \setminus \{\lfloor d/2 \rfloor + 1, i, \lfloor d/2 \rfloor + i\}, \mathbb{P}_j^{(i)} &\text{ is } \mathcal{N}((1 - \Delta_{\min})/2, 1), \\ \forall i \in \{2, \lfloor d/2 \rfloor\}, \mathbb{P}_i^{(i)} &\text{ is } \mathcal{N}((1 + 3\Delta_{\min})/2, 1) \quad \mathbb{P}_{\lfloor d/2 \rfloor + i}^{(i)} \text{ is } \mathcal{N}((1 + 3\Delta_{\min})/2, 1) \\ \forall i \in \{1, \lfloor d/2 \rfloor\}, \mathbb{P}_{d+1}^{(i)} &\text{ is } \mathcal{N}(-(1 - \alpha)(1 + \Delta_{\neq} + \Delta_{\min})/2 + \Delta_{\neq}/2, 1), \\ \mathbb{P}_{d+1}^{(\lfloor d/2 \rfloor + 1)} &\text{ is } \mathcal{N}(-(1 - \alpha)(1 - \Delta_{\neq} + \Delta_{\min})/2 - \Delta_{\neq}/2, 1) \quad \text{where } \alpha = 2/(\sqrt{\kappa_*} + 1). \end{aligned}$$

So that:

- Fact 3.1** For any $i \in \{2, \dots, \lfloor d/2 \rfloor\}$, between **Problem 1** and **Problem i**, the only actions that provide different evaluations when sampled are action i and action $\lfloor d/2 \rfloor + i$, and the mean gaps in both cases is $2\Delta_{\min}$.
- Fact 3.2** Between **Problem 1** and **Problem** $\lfloor d/2 \rfloor + 1$, the only action that provide different evaluation when sampled is action $d + 1$, and the mean gap in this case is $\alpha\Delta_{\neq}$.

For $j \leq d+1$, we write $N_j(T)$ for the total number of times action x_j has been selected before time T . Then, for $j \in \{1, \dots, \lfloor d/2 \rfloor\}$, let $E^{(j)} = \{N_i(T) \leq T/2\}$. Note that for $i \in \{1, \dots, \lfloor d/2 \rfloor\}$, in **Problem 1** the action x_i is the optimal action. Therefore, for any efficient algorithm, for all $i \in \{1, \dots, \lfloor d/2 \rfloor\}$ the event $E^{(i)}$ should have a low probability under $\mathbb{P}^{(i)}$. Indeed, for $i \in \{1, \dots, \lfloor d/2 \rfloor\}$, the regret of the algorithm under **Problem 1** can be lower-bounded as follows - see **Facts 1 and 2**:

$$R_T^{(i)} \geq \sum_{j \leq \lfloor d/2 \rfloor, j \neq i} \mathbb{E}^{(i)} [N_j(T)] \Delta_{\min} + \sum_{\lfloor d/2 \rfloor + 1 \leq j \leq d} \mathbb{E}^{(i)} [N_j(T)] \Delta_{\neq} + \frac{\mathbb{E}^{(i)} [N_{d+1}(T)]}{2}. \quad (34)$$

Since $\sum_j \mathbb{E}^{(i)} [N_j(T)] = T$ and $\Delta_{\min} \leq \Delta_{\neq} \leq \frac{1}{2}$, this implies together with **Facts 1**:

$$R_T^{(i)} \geq \left(T - \mathbb{E}^{(i)} [N_i(T)] \right) \Delta_{\min}$$

Using the definition of $E^{(i)}$, we find that

$$R_T^{(i)} \geq \frac{T \Delta_{\min}}{2} \mathbb{P}^{(i)} \left(E^{(i)} \right). \quad (35)$$

In particular for **Problem 1**, for any $i \in \{1, \dots, \lfloor d/2 \rfloor\}$,

$$R_T^{(1)} \geq \frac{T \Delta_{\min}}{2} \mathbb{P}^{(1)} \left(\overline{E^{(i)}} \right). \quad (36)$$

since $E^{(1)} \supset \overline{E^{(i)}}$.

Similarly, let us also define the event $F = \left\{ \sum_{i \leq \lfloor d/2 \rfloor} N_i(T) \geq T/2 \right\}$. Then, in **Problem 1**, the group 1 contains the optimal action, and so for any efficient algorithm, the event F should have a low probability under $\mathbb{P}^{(1)}$. Indeed, Equation (34) also implies

$$R_T^{(1)} \geq \left(T - \mathbb{E}^{(1)} \left[\sum_{i \leq \lfloor d/2 \rfloor} N_i(T) \right] \right) \Delta_{\neq} \geq \frac{T \Delta_{\neq}}{2} \mathbb{P}^{(1)} \left(\overline{F} \right). \quad (37)$$

On the other hand, for any efficient algorithm, the event F should have high probability under $\mathbb{P}^{(\lfloor d/2 \rfloor + 1)}$. Indeed, under problem **Problem $\lfloor d/2 \rfloor + 1$** , the regret can be lower-bounded as follows - see **Facts 1 and 2**:

$$R_T^{(\lfloor d/2 \rfloor + 1)} \geq \sum_{j \leq \lfloor d/2 \rfloor} \mathbb{E}^{(\lfloor d/2 \rfloor + 1)} [N_j(T)] \Delta_{\neq} + \sum_{\lfloor d/2 \rfloor + 2 \leq j \leq d} \mathbb{E}^{(\lfloor d/2 \rfloor + 1)} [N_j(T)] \Delta_{\min} + \frac{\mathbb{E}^{(\lfloor d/2 \rfloor + 1)} [N_{d+1}(T)]}{2}.$$

which implies that

$$R_T^{(\lfloor d/2 \rfloor + 1)} \geq \sum_{j \leq \lfloor d/2 \rfloor} \mathbb{E}^{(\lfloor d/2 \rfloor + 1)} [N_j(T)] \Delta_{\neq} \geq \frac{T \Delta_{\neq}}{2} \mathbb{P}^{(\lfloor d/2 \rfloor + 1)} (F). \quad (38)$$

Now, Bretagnolle-Huber inequality (see, e.g., Theorem 14.2 in [21]) implies that for all $i \in \{2, \dots, \lfloor d/2 \rfloor\}$,

$$\frac{1}{2} \exp \left(-KL \left(\mathbb{P}^{(1)}, \mathbb{P}^{(i)} \right) \right) \leq \mathbb{P}^{(i)} \left(E^{(i)} \right) + \mathbb{P}^{(1)} \left(\overline{E^{(i)}} \right) \quad (39)$$

and that

$$\frac{1}{2} \exp \left(-KL \left(\mathbb{P}^{(1)}, \mathbb{P}^{(\lfloor d/2 \rfloor + 1)} \right) \right) \leq \mathbb{P}^{(\lfloor d/2 \rfloor + 1)} (F) + \mathbb{P}^{(1)} (\overline{F}). \quad (40)$$

On the one hand, Equation (39) implies that for any $i \in \{2, \dots, \lfloor d/2 \rfloor\}$,

$$\begin{aligned} KL \left(\mathbb{P}^{(1)}, \mathbb{P}^{(i)} \right) &\geq -\log \left(2\mathbb{P}^{(i)} \left(E^{(i)} \right) + 2\mathbb{P}^{(1)} \left(\overline{E^{(i)}} \right) \right) \\ &\geq \log(T) - \log \left(2T\mathbb{P}^{(i)} \left(E^{(i)} \right) + 2T\mathbb{P}^{(1)} \left(\overline{E^{(i)}} \right) \right). \end{aligned} \quad (41)$$

Combining Equations (35), (36), and (41), we find that

$$KL(\mathbb{P}^{(1)}, \mathbb{P}^{(i)}) \geq \log(T) - \log\left(\frac{4(R_T^{(i)} + R_T^{(1)})}{\Delta_{\min}}\right). \quad (42)$$

On the other hand, Equation (40) implies that

$$\begin{aligned} KL(\mathbb{P}^{(1)}, \mathbb{P}^{(\lfloor d/2 \rfloor + 1)}) &\geq -\log\left(2\mathbb{P}^{(\lfloor d/2 \rfloor + 1)}(F) + 2\mathbb{P}^{(1)}(\bar{F})\right) \\ &\geq \log(T) - \log\left(2T\mathbb{P}^{(\lfloor d/2 \rfloor + 1)}(F) + 2T\mathbb{P}^{(1)}(\bar{F})\right). \end{aligned} \quad (43)$$

Combining Equations (35), (36), and (43), we find that

$$KL(\mathbb{P}^{(1)}, \mathbb{P}^{(\lfloor d/2 \rfloor + 1)}) \geq \log(T) - \log\left(\frac{4(R_T^{(\lfloor d/2 \rfloor + 1)} + R_T^{(1)})}{\Delta_{\neq}}\right). \quad (44)$$

Also, note that for all $i \in \{2, \dots, \lfloor d/2 \rfloor + 1\}$, the Kullback-Leibler divergence between $\mathbb{P}^{(1)}$ and $\mathbb{P}^{(i)}$ can be decomposed as follows (see, e.g., Lemma 15.1 in [21]) :

$$KL(\mathbb{P}^{(1)}, \mathbb{P}^{(i)}) = \sum_{j \leq d+1} \mathbb{E}^{(1)}[N_j(T)] KL(\mathbb{P}_j^{(1)}, \mathbb{P}_j^{(i)}). \quad (45)$$

Lower bound in $d\Delta_{\min}^{-1} \log T$. By design, for $i \in \{2, \dots, \lfloor d/2 \rfloor\}$, all actions but x_i and $x_{\lfloor d \rfloor + i}$ have the same distribution under $\mathbb{P}^{(1)}$ and $\mathbb{P}^{(i)}$ - see **Fact 3.1**. Then, Equation (45) becomes from **Fact 3.1** and from the expression of KL divergence between standard Gaussian distributions:

$$KL(\mathbb{P}^{(1)}, \mathbb{P}^{(i)}) = \frac{4\Delta_{\min}^2}{2} \mathbb{E}^{(1)}[N_i(T)] + \frac{4\Delta_{\min}^2}{2} \mathbb{E}^{(1)}[N_{\lfloor d \rfloor + i}(T)].$$

So that, summing over $i \in \{2, \dots, \lfloor d/2 \rfloor\}$, and by **Fact 1**:

$$\sum_{i \in \{2, \dots, \lfloor d/2 \rfloor\}} KL(\mathbb{P}^{(1)}, \mathbb{P}^{(i)}) \leq 2\Delta_{\min} R_T^{(1)}.$$

So that by Equation (42) (summing over $i \in \{2, \dots, \lfloor d/2 \rfloor\}$):

$$\begin{aligned} 2\Delta_{\min} R_T^{(1)} &\geq \sum_{i \in \{2, \dots, \lfloor d/2 \rfloor\}} \left[\log(T) - \log\left(\frac{4(R_T^{(i)} + R_T^{(1)})}{\Delta_{\min}}\right) \right] \\ &= (\lfloor d/2 \rfloor - 1) \log(T) - \sum_{i \in \{2, \dots, \lfloor d/2 \rfloor\}} \log\left(\frac{4(R_T^{(i)} + R_T^{(1)})}{\Delta_{\min}}\right). \end{aligned}$$

Let us assume that our algorithm satisfies $\max_{i \leq \lfloor d/2 \rfloor} R_T^{(i)} \leq \frac{d \log(T)}{\Delta_{\min}}$ - otherwise the bound immediately follows for this algorithm. Then

$$\begin{aligned} R_T^{(1)} &\geq \frac{1}{2\Delta_{\min}} (\lfloor d/2 \rfloor - 1) \log(T) - \frac{1}{2\Delta_{\min}} \sum_{i \in \{2, \dots, \lfloor d/2 \rfloor\}} \log\left(\frac{8d \log T}{\Delta_{\min}^2}\right) \\ &\geq \frac{1}{2\Delta_{\min}} (\lfloor d/2 \rfloor - 1) \left[\log(T) - \log\left(\frac{8d \log(T)}{\Delta_{\min}^2}\right) \right]. \end{aligned} \quad (46)$$

Sine $d \geq 4$, we note that $\lfloor d/2 \rfloor - 1 \geq d/5$. This concludes the proof for this part of the bound.

Lower bound in $\kappa_* \Delta_{\neq}^{-2} \log T$. By design, all actions but x_{d+1} have the same evaluation under **Problem 1** and **Problem $\lfloor d/2 \rfloor + 1$** - see **Fact 3.2**. Then, by **Fact 3.2** and the expression between the KL divergence of standard Gaussians, Equation (45) becomes

$$KL(\mathbb{P}^{(1)}, \mathbb{P}^{(\lfloor d/2 \rfloor + 1)}) = \mathbb{E}^{(1)} [N_{d+1}(T)] \frac{(\alpha \Delta_{\neq})^2}{2} = \frac{1}{2} \mathbb{E}^{(1)} [N_{d+1}(T)] \left(\frac{2\Delta_{\neq}}{\sqrt{\kappa_*} + 1} \right)^2.$$

Combined with equation (44), this implies that

$$\frac{1}{2} \mathbb{E}^{(1)} [N_{d+1}(T)] \left(\frac{2\Delta_{\neq}}{\sqrt{\kappa_*} + 1} \right)^2 \geq \log(T) - \log \left(\frac{4(R_T^{(\lfloor d/2 \rfloor + 1)} + R_T^{(1)})}{\Delta_{\neq}} \right). \quad (47)$$

Let us assume that our algorithm satisfies $\max_{i \leq \lfloor d/2 \rfloor + 1} R_T^{(i)} \leq \frac{\kappa_* \log(T)}{\Delta_{\neq}^2}$ - otherwise the bound immediately follows for this algorithm. We then have

$$\frac{1}{2} \mathbb{E}^{(1)} [N_{d+1}(T)] \left(\frac{2\Delta_{\neq}}{\sqrt{\kappa_*} + 1} \right)^2 \geq \log(T) - \log \left(\frac{8\kappa_* \log(T)}{\Delta_{\neq}^3} \right).$$

Using Equation (34), we find that

$$R_T^{(1)} \geq \frac{\kappa_* + 1}{4\Delta_{\neq}^2} \left[\log(T) - \log \left(\frac{8\kappa_* \log(T)}{\Delta_{\neq}^3} \right) \right]. \quad (48)$$

Lower bound in $\kappa_* \Delta_{\neq}^{-2}$. Let us assume that our algorithm satisfies $\max_{i \leq \lfloor d/2 \rfloor + 1} R_T^{(i)} \leq \frac{\kappa_*}{\Delta_{\neq}^2}$ - otherwise the bound immediately follows for this algorithm. Then, Equation (47) implies

$$\frac{1}{2} \mathbb{E}^{(1)} [N_{d+1}(T)] \left(\frac{2\Delta_{\neq}}{\sqrt{\kappa_*}} \right)^2 \geq \log(T) - \log \left(\frac{8\kappa_*}{\Delta_{\neq}^3} \right).$$

Using again Equation (34), we find that

$$R_T^{(1)} \geq \frac{\kappa_* + 1}{4\Delta_{\neq}^2} \log \left(\frac{T\Delta_{\neq}^3}{8\kappa_*} \right). \quad (49)$$

We conclude the proof of Theorem 6 by combining Equations (46), (48) and (49).

Bounds on $\kappa(\Delta)$ Finally, the following lemma allows to express $\kappa(\Delta)$ as a function of κ_* .

Lemma 17. *For any $i \in \{1, \dots, \lfloor d/2 \rfloor + 1\}$, the gap vector Δ verifies*

$$\kappa(\Delta) = \frac{(1 + \sqrt{\kappa_*})^2 \Delta_{d+1}}{4}$$

where $\Delta_{d+1} = \max_i (x_i - x_{d+1})^\top \gamma^{(i)}$.

On the one hand, since $\kappa_* \geq 1$, we see that $\kappa_* \leq (1 + \sqrt{\kappa_*})^2 \leq 4\kappa_*$. On the other hand, $1/2 \leq \Delta_{d+1} \leq 2$, so $\kappa(\Delta) \in [\frac{\kappa_*}{8}, 2\kappa_*]$.

C.5 Extension of the gap-dependent lower bounds to $d = 2, 3$

Theorem 4 can be extended to $d \in \{2, 3\}$ by considering separately the cases $\frac{d}{\Delta_{\min}} \geq \frac{\kappa}{\Delta_{\neq}^2}$ and $\frac{d}{\Delta_{\min}} < \frac{\kappa}{\Delta_{\neq}^2}$.

Case 1 : $\frac{d}{\Delta_{\min}} \geq \frac{\kappa}{\Delta_{\neq}^2}$ Let us consider the set of actions defined by $\mathcal{A} = \left\{ \begin{pmatrix} x_1 \\ z_{x_1} \end{pmatrix}, \dots, \begin{pmatrix} x_{d+1} \\ z_{x_{d+1}} \end{pmatrix} \right\}$, where $\begin{pmatrix} x_i \\ z_{x_i} \end{pmatrix} = e_1 + e_{d+1}$ for $i \in \{1, \dots, d\}$, and $\begin{pmatrix} x_{d+1} \\ z_{x_{d+1}} \end{pmatrix} = -\left(1 - \frac{2}{\sqrt{\kappa_* + 1}}\right) e_1 - e_{d+1}$. Using the same proof as in Lemma 15, we see that

$$\min_{\pi \in \mathcal{P}^{\mathcal{A}}} \left\{ e_{d+1}^\top \left(\sum_{\begin{pmatrix} x \\ z \end{pmatrix} \in \mathcal{A}} \pi_x \begin{pmatrix} x \\ z_x \end{pmatrix} \begin{pmatrix} x \\ z_x \end{pmatrix}^\top \right)^+ e_{d+1} \right\} = \kappa.$$

Then, we consider the following problems : for $i \leq d$, **Problem i** is characterized by the parameter $\theta^{(i)}$, where $\theta^{(i)} = \begin{pmatrix} \gamma^{(i)} \\ \omega^{(i)} \end{pmatrix}$ is defined as:

$$\begin{aligned} \gamma^{(1)} &= \frac{1 - \Delta_{\min}}{2} \sum_{i \leq d} e_i + \Delta_{\min} e_1 \\ \gamma^{(i)} &= \frac{1 - \Delta_{\min}}{2} \sum_{i \leq d} e_i + \Delta_{\min} e_1 + \Delta_{\min} e_i \quad \text{for } i > 1 \end{aligned}$$

and the bias parameters are defined as $\omega^{(i)} = 0$ for $i \leq d$. The following facts hold:

Fact 1 For any $i \in \{1, \dots, d\}$, action x_i is the unique optimal action in **Problem i**. Sampling any other (sub-optimal) action leads to an instantaneous regret of at least Δ_{\min} .

Fact 2 In **Problem i**, for $i \in \{1, \dots, d\}$, when sampling action x_j at time, t the distribution of the observation does not depend on t or on the past (except through the choice of x_j) and is $\mathbb{P}_j^{(i)}$. It is characterized as:

$$\begin{aligned} \forall i \in \{1, \dots, d\}, \mathbb{P}_1^{(i)} &\text{ is } \mathcal{N}((1 + \Delta_{\min})/2, 1) \\ \forall i \in \{1, \dots, d\}, \mathbb{P}_{d+1}^{(1)} &\text{ is } \mathcal{N}\left(-\left(1 - \frac{2}{\sqrt{\kappa_* + 1}}\right)(1 + \Delta_{\min})/2, 1\right) \\ \forall i \in \{2, \dots, d\}, \mathbb{P}_i^{(i)} &\text{ is } \mathcal{N}((1 + 3\Delta_{\min})/2, 1) \\ \forall i, j \in \{2, \dots, d\}, i \neq j : \mathbb{P}_j^{(i)} &\text{ is } \mathcal{N}((1 - \Delta_{\min})/2, 1) \end{aligned}$$

So that for any $i \in \{2, \dots, d\}$, between **Problem 1** and **Problem i**, the only action that provides different evaluations when sampled is action i , and the mean gap is $2\Delta_{\min}$.

Since $\Delta_{\neq} \leq \frac{1}{8}$, this choice of parameters ensures that $\forall i \in \{1, \dots, d\}$, $\theta^{(i)} \in \Theta_{\Delta_{\min}, \Delta_{\neq}, \kappa_*}^{\mathcal{A}}$. Adapting the proof of Lemma 15, we note that the minimal variance of bias estimation is at least κ_* . This proves that $\mathcal{A} \in \Theta_{\Delta_{\min}, \Delta_{\neq}, \kappa_*}^{\mathcal{A}}$. Now, the lower bound

$$R_T \geq \frac{d-1}{2\Delta_{\min}} \left[\log(T) - \log\left(\frac{8d \log(T)}{\Delta_{\min}^2}\right) \right]$$

follows directly using arguments from the proof of Theorem 6.

Case 2 : $\frac{d}{\Delta_{\min}} > \frac{\kappa}{\Delta_{\neq}^2}$ Let the action set be given by $\mathcal{A} = \left\{ \begin{pmatrix} x_1 \\ z_{x_1} \end{pmatrix}, \dots, \begin{pmatrix} x_{d+1} \\ z_{x_{d+1}} \end{pmatrix} \right\}$, where $\begin{pmatrix} x_1 \\ z_{x_1} \end{pmatrix} = e_1 + e_{d+1}$, $\begin{pmatrix} x_i \\ z_{x_i} \end{pmatrix} = e_i - e_{d+1}$ for $i \in \{2, \dots, d\}$, and $\begin{pmatrix} x_{d+1} \\ z_{x_{d+1}} \end{pmatrix} = -\left(1 - \frac{2}{\sqrt{\kappa_* + 1}}\right) e_1 - e_{d+1}$. By Lemma 15, $\mathcal{A} \in \mathbf{A}_{\kappa_*, d}$. We consider two bandit problems characterized by two parameters $\theta^{(1)}$ and $\theta^{(2)}$, defined as:

$$\begin{aligned} \gamma^{(1)} &= \frac{1 + \Delta_{\neq}}{2} e_1 + \frac{1 - \Delta_{\neq}}{2} e_2 - \frac{\Delta_{\neq}}{2} e_3 \\ \gamma^{(2)} &= \frac{1 - \Delta_{\neq}}{2} e_1 + \frac{1 + \Delta_{\neq}}{2} e_2 + \frac{\Delta_{\neq}}{2} e_3. \end{aligned}$$

On top of this, two bias parameters are defined as $\omega^{(1)} = -\frac{\Delta_{\neq}}{2}$ and $\omega^{(2)} = \frac{\Delta_{\neq}}{2}$.

The following facts hold:

Fact 1 For any $i \in \{1, 2\}$, action x_i is the unique optimal action in **Problem i**. Since $1/2 \geq \Delta_{\neq}$, sampling any other (sub-optimal) action leads to an instantaneous regret of at least Δ_{\neq} .

Fact 2 In **Problem i**, for $i \in \{1, \dots, d\}$, when sampling action x_j at time, t the distribution of the observation does not depend on t or on the past (except through the choice of x_j) and is $\mathbb{P}_j^{(i)}$. It is characterized as:

$$\begin{aligned} \forall i \in \{1, 2\}, \forall j \in \{1, 2\}, \mathbb{P}_j^{(i)} & \text{ is } \mathcal{N}(1/2, 1) \\ \forall i \in \{1, 2\}, \mathbb{P}_3^{(1)} & \text{ is } \mathcal{N}(0, 1) \\ \mathbb{P}_{d+1}^{(1)} & \text{ is } \mathcal{N}\left(\left(1 - \frac{2}{\sqrt{\kappa_*} + 1}\right) \left(\frac{1 + \Delta_{\neq}}{2}\right) + \frac{\Delta_{\neq}}{2}, 1\right) \\ \mathbb{P}_{d+1}^{(2)} & \text{ is } \mathcal{N}\left(\left(1 - \frac{2}{\sqrt{\kappa_*} + 1}\right) \left(\frac{1 - \Delta_{\neq}}{2}\right) - \frac{\Delta_{\neq}}{2}, 1\right) \end{aligned}$$

So that, between **Problem 1** and **Problem 2**, the only action that provides different evaluations when sampled is action 1, and the mean gaps in both cases is $\frac{2\Delta_{\neq}}{\sqrt{\kappa_*} + 1}$.

Note that the minimum gap for these parameters is $\Delta_{\neq} \geq \Delta_{\min}$. Thus, this choice of parameters ensures that $\forall i \in \{1, \dots, d\}, \theta^{(i)} \in \Theta_{\Delta_{\min}, \Delta_{\neq}, \kappa_*}^A$. Adapting the proof of Lemma 15, we note that the minimal variance of bias estimation is at least κ_* . This proves that $\mathcal{A} \in \Theta_{\Delta_{\min}, \Delta_{\neq}, \kappa_*}^A$. Then, the lower bound

$$R_T \geq \frac{\kappa_* + 1}{4\Delta_{\neq}^2} \left[\log(T) - \log\left(\frac{8\kappa_* \log(T)}{\Delta_{\neq}^3}\right) \right].$$

follows directly using arguments from the proof of Theorem 6.

C.6 Auxiliary Lemmas

C.6.1 Proof of Lemma 1

Lemma 1 follows from the characterization of κ_* given in Lemma 5. We begin by proving the first statement. Assume that $\kappa_* > 1$ (otherwise the first statement is void). Note that for all $u \in \mathbb{R}^d$, $\lim_{\lambda \rightarrow +\infty} (\max_{x \in \mathcal{X}} (x^\top (\lambda u) + z_x)^2)^{-1} = 0$, so the minimum over $u \in \mathbb{R}^d$ of $(\max_{x \in \mathcal{X}} (x^\top (\lambda u) + z_x)^2)^{-1}$ is attained for some vector $\tilde{u} \in \mathbb{R}^d$. Since $\kappa_* > 1$, \tilde{u} is not null. Moreover, $\max_{x \in \mathcal{X}} (1 + z_x x^\top \tilde{u})^2 < 1$, so $\max_{x \in \mathcal{X}} z_x x^\top \tilde{u} < 0$. Thus, for all $x \in \mathcal{X}$, $x^\top \tilde{u}$ and z_x are of opposite sign, and $x^\top \tilde{u} \neq 0$. This implies that the hyperplane containing 0 with normal vector \tilde{u} contains no action, and separates the two groups. Moreover,

$$\kappa_*^{-1/2} = \max_{x \in \mathcal{X}} |z_x x^\top \tilde{u} + 1|.$$

We denote $x^{(1)} \in \operatorname{argmax}_{x \in \mathcal{X}} z_x x^\top \tilde{u}$, and $x^{(2)} \in \operatorname{argmin}_{x \in \mathcal{X}} z_x x^\top \tilde{u}$. Let us show that $(z_{x^{(1)}} x^{(1)\top} \tilde{u} + 1) = -(1 + z_{x^{(2)}} x^{(2)\top} \tilde{u})$, i.e that $z_{x^{(1)}} x^{(1)\top} \tilde{u} + z_{x^{(2)}} x^{(2)\top} \tilde{u} = -2$. Indeed, note that

$$\kappa_*^{-1/2} = (z_{x^{(1)}} x^{(1)\top} \tilde{u} + 1) \vee -(1 + z_{x^{(2)}} x^{(2)\top} \tilde{u}).$$

Then, for $u' = \frac{-2}{(z_{x^{(1)}} x^{(1)\top} \tilde{u} + z_{x^{(2)}} x^{(2)\top} \tilde{u})} \tilde{u}$, we see that

$$z_{x^{(1)}} x^{(1)\top} u' + 1 = -(1 + z_{x^{(2)}} x^{(2)\top} u') = \max_{x \in \mathcal{X}} |z_x x^\top u' + 1|.$$

By contradiction, let us first assume that $z_{x^{(1)}} x^{(1)\top} \tilde{u} + z_{x^{(2)}} x^{(2)\top} \tilde{u} < -2$. Then,

$$\max_{x \in \mathcal{X}} |z_x x^\top u' + 1| = z_{x^{(1)}} x^{(1)\top} u' + 1 < z_{x^{(1)}} x^{(1)\top} \tilde{u} + 1 = \kappa_*^{-1/2}$$

which contradicts the definition of κ_* .

Similarly, if we assume that $z_{x^{(1)}}x^{(1)\top}\tilde{u} + z_{x^{(2)}}x^{(2)\top}\tilde{u} > -2$, then

$$\max_{x \in \mathcal{X}} |z_x x^\top u' + 1| = -(z_{x^{(2)}}x^{(2)\top}u' + 1) < -(z_{x^{(2)}}x^{(2)\top}\tilde{u} + 1) = \kappa_*^{-1/2}$$

which contradicts again the definition of κ_* . Therefore,

$$(z_{x^{(1)}}x^{(1)\top}\tilde{u} + 1) = -\left(1 + z_{x^{(2)}}x^{(2)\top}\tilde{u}\right) = \kappa_*^{-1/2}.$$

Then, the hyperplane containing 0 with normal vector \tilde{u} separates the actions of the two groups. Moreover, the margin is $-z_{x^{(1)}}x^{(1)\top}\tilde{u} = 1 - \kappa_*^{-1/2}$, while the maximum distance of all points is $-z_{x^{(2)}}x^{(2)\top}\tilde{u} = 1 + \kappa_*^{-1/2}$. Thus, there exists \tilde{u} such that the hyperplane containing 0 with normal vector \tilde{u} separates the actions of the two groups, with margin equal to $\frac{\sqrt{\kappa_*}-1}{\sqrt{\kappa_*}+1}$ times the maximum distance of all points to the hyperplane.

Conversely, assume that there exists $\kappa > \kappa_*$ such that there exists $u \in \mathbb{R}^d$ such that the hyperplane containing 0 with normal vector u separates the actions of the two groups, with margin equal to $\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1} = \frac{1-\kappa^{-1/2}}{1+\kappa^{-1/2}}$ times the maximum distance of all points to the hyperplane, denoted hereafter d . Since the hyperplane separates the points, we can assume without loss of generality that for all $x \in \mathcal{X}$, $z_x x^\top u < 0$. Similarly, up to a renormalization, we can assume without loss of generality that $d = 1 + \kappa^{-1/2}$. Then,

$$\begin{aligned} \max_{x \in \mathcal{X}} |z_x x^\top u + 1| &= \left(\max_{x \in \mathcal{X}} z_x x^\top u + 1\right) \vee -\left(\min_{x \in \mathcal{X}} z_x x^\top u + 1\right) \\ &= \left(-\frac{1 - \kappa^{-1/2}}{1 + \kappa^{-1/2}} \times (1 + \kappa^{-1/2}) + 1\right) \vee -(1 - \kappa^{-1/2} - 1) = \kappa^{-1/2} < \kappa_*^{-1/2} \end{aligned}$$

which contradicts the definition of κ_* . This concludes the proof of the first statement.

To prove the second statement, let us assume that no separating hyperplane containing zero exists. Then, for all $u \in \mathbb{R}^d$, there exists $x \in \mathcal{X}$ such that $z_x x^\top u \geq 0$. This implies that $\min_{u \in \mathbb{R}^d} \max_{x \in \mathcal{X}} (z_x x^\top u + 1) \geq 1$, so $\kappa_* \leq 1$. Choosing $u = 0$, we see that $\kappa_* \geq 1$, which implies that $\kappa_* = 1$.

C.6.2 Proof of Lemma 2

Since for all $\gamma \in \mathcal{X}$ and all $x \in \mathcal{X}$, $|x^\top \gamma| \leq 1$, it is easy to see that the gaps are bounded by 2, and that $\tilde{\kappa} \leq 2\kappa_*$.

Let us now show that $\tilde{\kappa} \geq \kappa_*/2$.

$$\begin{aligned} \left(x^{(1)}, x^{(2)}, \tilde{\gamma}\right) &\in \operatorname{argmax}_{(x, x') \in \mathcal{X}, \gamma \in \mathcal{C}(\mathcal{X})} (x - x')^\top \gamma \\ \bar{x} &= \frac{1}{2}(x^{(1)} + x^{(2)}) \\ \tilde{n} &= \sum_{x \in \mathcal{X}} \tilde{\mu}(x) \\ \text{and } \tilde{x} &= \frac{1}{\tilde{n}} \sum_{x \in \mathcal{X}} \tilde{\mu}(x)x. \end{aligned}$$

Recall that κ_* can equivalently be defined as the budget necessary to estimate the bias with a variance smaller than 1. Therefore, we have

$$\tilde{n} \geq \kappa_*. \tag{50}$$

Let us define Δ_{\max} as $\Delta_{\max} = (x^{(1)} - x^{(2)})^\top \tilde{\gamma} = \max_{(x, x') \in \mathcal{X}, \gamma \in \mathcal{C}(\mathcal{X})} (x - x')^\top \gamma$. By definition of $\tilde{\kappa}$ and $\tilde{\mu}$,

$$\begin{aligned} \tilde{\kappa} &\geq \sum_{x \in \mathcal{X}} \tilde{\mu}(x)(x^{(1)} - x)^\top \tilde{\gamma} \\ &= \tilde{n}(x^{(1)} - \tilde{x})^\top \tilde{\gamma}. \end{aligned}$$

Using Equation (50), we find that

$$\begin{aligned}\frac{\tilde{\kappa}}{\kappa_*} &\geq (x^{(1)} - \bar{x})^\top \tilde{\gamma} + (\bar{x} - \tilde{x})^\top \tilde{\gamma} \\ &= \frac{\Delta_{\max}}{2} + (\bar{x} - \tilde{x})^\top \tilde{\gamma}.\end{aligned}\tag{51}$$

Now, since $\tilde{\gamma} \in \mathcal{C}(\mathcal{X})$, we also have $-\tilde{\gamma} \in \mathcal{C}(\mathcal{X})$, and therefore

$$\begin{aligned}\tilde{\kappa} &\geq \sum_{x \in \mathcal{X}} \tilde{\mu}(x) (x^{(2)} - x)^\top (-\tilde{\gamma}) \\ &= \tilde{n}(\tilde{x} - x^{(2)})^\top \tilde{\gamma}\end{aligned}$$

Using again Equation (50), we find that

$$\begin{aligned}\frac{\tilde{\kappa}}{\kappa_*} &\geq (\tilde{x} - \bar{x})^\top \tilde{\gamma} + (\bar{x} - x^{(2)})^\top \tilde{\gamma} \\ &= (\tilde{x} - \bar{x})^\top \tilde{\gamma} + \frac{\Delta_{\max}}{2}.\end{aligned}\tag{52}$$

Combining Equations (51) and (52), we find that

$$\frac{\tilde{\kappa}}{\kappa_*} \geq \frac{\Delta_{\max}}{2} + |(\bar{x} - \tilde{x})^\top \tilde{\gamma}|.$$

This implies in particular that $\tilde{\kappa} \geq \frac{\Delta_{\max} \kappa_*}{2}$.

To conclude the proof of the Lemma, we show that $\Delta_{\max} \geq 1$. By contradiction, assume that $\Delta_{\max} < 1$.

For all non-zero vector $u \in \mathbb{R}^d$, let us denote $x_u = \operatorname{argmax}_{x \in \mathcal{X}} |x^\top u|$. Since \mathcal{X} spans \mathbb{R}^d , we necessarily have $|x_u^\top u| > 0$, so we can define the normalized vector $\tilde{u} = u/|x_u^\top u|$ such that \tilde{u} belongs to the set $\mathcal{C}(\mathcal{X})$. Finally, denote $x_u^{(1)}, x_u^{(2)} \in \operatorname{argmax}_{x, x' \in \mathcal{X}} (x_u^{(1)} - x_u^{(2)})^\top \tilde{u}$. Note that by definition of Δ_{\max} , we always have $(x_u^{(1)} - x_u^{(2)})^\top \tilde{u} \leq \Delta_{\max} < 1$.

Case 1 : $x_u^\top \tilde{u} > 0$ Then, by definition of x_u and $x_u^{(1)}$, we see that $x_u^{(1)\top} \tilde{u} = x_u^\top \tilde{u} = 1$. Then, $(x_u^{(1)} - x_u^{(2)})^\top \tilde{u} < 1$ implies that $1 - x_u^{(2)\top} \tilde{u} < 1$, so $x_u^{(2)\top} \tilde{u} > 0$, and in particular $x_u^{(2)\top} u > 0$.

Case 2 : $x_u^\top \tilde{u} < 0$ Then, by definition of x_u and $x_u^{(2)}$, we see that $x_u^{(2)\top} \tilde{u} = x_u^\top \tilde{u} = -1$. Then $(x_u^{(1)} - x_u^{(2)})^\top \tilde{u} < 1$ implies that $x_u^{(1)\top} \tilde{u} + 1 < 1$, so $x_u^{(1)\top} \tilde{u} < 0$, and in particular $x_u^{(1)\top} u < 0$.

Putting together Case 1 and Case 2, we see that $x_u^{(1)\top} u$ and $x_u^{(2)\top} u$ are of the same sign and are not null. By definition of $x_u^{(1)}$ and $x_u^{(2)}$, we conclude that for all $x \in \mathcal{X}$, the sign of $x^\top u$ is the same, and that $x^\top u$ is not 0. Since this is true for all non-zero vector u , this implies in particular that no hyperplane containing the origin can separate the actions, which contradicts the assumption that \mathcal{X} spans \mathbb{R}^d .

C.6.3 Proof of Lemmas 3 and 4

We begin by proving Lemma 4. Recall that π is a G-optimal design for the set $\{a_x : x \in \mathcal{X}\}$, and that μ is defined as $\mu(x) = \lceil m\pi(x) \rceil$ for all $x \in \mathcal{X}$.

We first observe that $V(\pi) = A_\pi^\top A_\pi$, where A_π is the matrix with lines given by $[\sqrt{\pi(x)} a_x^\top]_{x \in \mathcal{X}}$. Since the supports of μ and π are the same, we get that $\operatorname{Range}(A_\pi^\top) = \operatorname{Range}(A_\mu^\top)$. As a consequence

$$\operatorname{Range}(V(\pi)) = \operatorname{Range}(A_\pi^\top) = \operatorname{Range}(A_\mu^\top) = \operatorname{Range}(V(\mu)),$$

and $x \in \operatorname{Range}(V(\mu))$ for all $x \in \mathcal{X}$. This ensures that $a_x^\top \hat{\theta}_\mu$ is an unbiased estimator of $a_x^\top \theta^*$.

Furthermore $V(\mu) \succcurlyeq mV(\pi)$, so the variance $a_x^\top V(\mu) + a_x$ of $a_x^\top \hat{\theta}_\mu$ is upper-bounded by $a_x^\top V(\mu) + a_x \leq m^{-1} a_x^\top V(\pi) + a_x$. Now, the General Equivalence Theorem of Kiefer and Pukelshein shows that $\max_{x \in \mathcal{X}} a_x^\top V(\pi) + a_x \leq d + 1$. Thus, $a_x^\top V(\pi) + a_x \leq m^{-1}(d + 1)$.

We now prove Lemma 3. Recall that $\pi \in \mathcal{M}_{e_{d+1}}^{\mathcal{X}}$ is such that $e_{d+1} \in \text{Range } V(\pi)$, and that μ is defined as $\mu(x) = \lceil m\pi(x) \rceil$ for all $x \in \mathcal{X}$. Using similar arguments, we can show that $e_{d+1} \in \text{Range}(V(\mu))$, which ensures that $e_{d+1}^\top \hat{\theta}_\mu$ is an unbiased estimator of $e_{d+1}^\top \theta^*$. The second part of the Lemma follows directly using that $V(\mu) \succcurlyeq mV(\pi)$.

C.6.4 Proof of Lemma 5

Elfving's set \mathcal{S} for estimating the bias in the biased linear bandit problem is given by

$$\mathcal{S} = \text{convex hull} \left\{ \begin{pmatrix} x \\ z_x \end{pmatrix}, \begin{pmatrix} -x \\ -z_x \end{pmatrix} : x \in \mathcal{X} \right\},$$

or equivalently by

$$\mathcal{S} = \text{convex hull} \left\{ \pm \begin{pmatrix} z_x x \\ 1 \end{pmatrix} : x \in \mathcal{X} \right\}.$$

Now, Theorem 5 indicates that $\kappa_*^{-1/2} e_{d+1}$ belongs to a supporting hyperplane of \mathcal{S} . We first show that when \mathcal{A} spans \mathbb{R}^{d+1} , any normal vector $w \in \mathbb{R}^{d+1}$ to this hyperplane is such that $w^\top e_{d+1} \neq 0$.

By contradiction, let us assume that $\kappa_*^{-1/2} e_{d+1}$ belongs to some supporting hyperplane \mathcal{H} of \mathcal{S} parametrized as $\mathcal{H} = \{a \in \mathbb{R}^{d+1} : a^\top w = b\}$, where the normal vector w is of the form $w = \begin{pmatrix} u \\ 0 \end{pmatrix}$. Then, $\kappa_*^{-1/2} e_{d+1} \in \mathcal{H}$, so $\kappa_*^{-1/2} e_{d+1}^\top w = b$, and thus $b = 0$. Now, \mathcal{H} is a supporting hyperplane of \mathcal{S} , so for all $a \in \mathcal{S}$ we see that $a^\top w \leq b$. In particular, for all $x \in \mathcal{X}$, $x^\top u \leq 0$ and $-x^\top u \leq 0$, so $x^\top u = 0$. This implies that \mathcal{X} is supported by an hyperplane in \mathbb{R}^d with normal vector u , which contradicts our assumption that \mathcal{A} spans \mathbb{R}^{d+1} . Thus, the supporting hyperplane of \mathcal{S} containing $\kappa_*^{-1/2} e_{d+1}$ has a normal vector $w \in \mathbb{R}^{d+1}$ such that $w^\top e_{d+1} \neq 0$. In particular, we can parameterize this hyperplane as $\mathcal{H}_{u,b} = \{a \in \mathbb{R}^{d+1} : a^\top \begin{pmatrix} u \\ 1 \end{pmatrix} = b\}$ for some $b \in \mathbb{R}$ and $u \in \mathbb{R}^d$.

Now, if $\mathcal{H}_{u,b}$ is a supporting hyperplane of \mathcal{S} , then, by definition, \mathcal{S} is contained in the half space $\{a \in \mathbb{R}^{d+1} : a^\top \begin{pmatrix} u \\ 1 \end{pmatrix} \leq b\}$. In particular, for all $x \in \mathcal{X}$, one must have $z_x x^\top u + 1 \leq b$ and $-z_x x^\top u - 1 \leq b$: therefore, for all $x \in \mathcal{X}$, $|z_x x^\top u + 1| \leq b$. Moreover, $\mathcal{H}_{u,b}$ is a supporting hyperplane of \mathcal{S} , so there exists an extreme point $a \in \mathcal{S}$ such that $a \in \mathcal{H}_{u,b}$. Note that \mathcal{S} is the convex hull of $\{\pm \begin{pmatrix} z_x x \\ 1 \end{pmatrix} : x \in \mathcal{X}\}$, so the extreme points of \mathcal{S} are in $\{\pm \begin{pmatrix} z_x x \\ 1 \end{pmatrix} : x \in \mathcal{X}\}$. In particular, this implies that $b = \max\{|z_x x^\top u + 1| : x \in \mathcal{X}\}$. Thus, the supporting hyperplane of \mathcal{S} containing $\kappa_*^{-1/2} e_{d+1}$ is necessarily of the form $\mathcal{H}_{u, \max\{|z_x x^\top u + 1| : x \in \mathcal{X}\}}$.

On the one hand, $\kappa_*^{-1/2} e_{d+1}$ belongs to the boundary of \mathcal{S} and therefore to a supporting hyperplane $\mathcal{H}_{u, \max\{|z_x x^\top u + 1| : x \in \mathcal{X}\}}$ of \mathcal{S} . Then, there exists $u \in \mathbb{R}^d$ such that $\kappa_*^{-1/2} = \max\{|z_x x^\top u + 1| : x \in \mathcal{X}\}$.

On the other hand, it is easy to verify that for all $u \in \mathbb{R}^d$, $\mathcal{H}_{u, \max\{|z_x x^\top u + 1| : x \in \mathcal{X}\}}$ is a supporting hyperplane of \mathcal{S} . Now, $\kappa_*^{-1/2} e_{d+1}$ belongs to \mathcal{S} , so $\kappa_*^{-1/2} e_{d+1}^\top \begin{pmatrix} u \\ 1 \end{pmatrix} \leq \max\{|z_x x^\top u + 1| : x \in \mathcal{X}\}$.

These two results imply that

$$\kappa_*^{-1/2} = \min_{u \in \mathbb{R}^d} \max_{x \in \mathcal{X}} |z_x x^\top u + 1|$$

which proves the Lemma.

C.6.5 Proof of Lemma 6

We prove that $2(\sqrt{\kappa_*} - 1)^2 \vee 1 \leq \alpha \leq 8(\kappa_* + 1)$. Lemma 6 follows directly by noticing that $\alpha \geq 1$ and $\kappa_* \geq 1$.

Let us begin by proving that $2(\sqrt{\kappa_*} - 1)^2 \leq \alpha$ for $\kappa_* > 1$ (otherwise this inequality is automatically verified). Note that for all $u \in \mathbb{R}^d$, $\lim_{\lambda \rightarrow +\infty} \frac{1}{\max_{x \in \mathcal{X}} (x^\top (\lambda u) + z_x)^2} = 0$, so the minimum over $u \in \mathbb{R}^d$ of $\frac{1}{\max_{x \in \mathcal{X}} (x^\top u + z_x)^2} = 0$ is attained for some vector $\tilde{u} \in \mathbb{R}^d$. Let us also denote $\tilde{x} \in \arg\max_{x \in \mathcal{X}} (z_x x^\top \tilde{u} + 1)^2$, such that

$$\kappa_* = \frac{1}{(z_{\tilde{x}} \tilde{x}^\top \tilde{u} + 1)^2}.$$

With these notations, we see that for all $x \in \mathcal{X}$,

$$(z_x x^\top \tilde{u} + 1)^2 \leq (z_{\tilde{x}} \tilde{x}^\top \tilde{u} + 1)^2 = \kappa_*^{-1} < 1.$$

This implies that for all $x \in \mathcal{X}$,

$$z_x x^\top \tilde{u} \leq -1 + \kappa_*^{-1/2} < 0.$$

Now, let us denote $x^{(1)}, x^{(2)} \in \operatorname{argmax}_{x, x' \in \mathcal{X}} (x - x')^\top \tilde{u}$. By definition of α , we see that

$$\alpha \geq \frac{\left((x^{(1)} - x^{(2)})^\top \tilde{u} \right)^2}{(z_{\tilde{x}} \tilde{x}^\top \tilde{u} + 1)^2} = \left((x^{(1)} - x^{(2)})^\top \tilde{u} \right)^2 \times \kappa_*.$$

Since $z_x x^\top \tilde{u} < 0$ for all $x \in \mathcal{X}$, and since no group is empty, we can conclude that there exists $x, x' \in \mathcal{X}$ such that $x^\top \tilde{u} > 0$ and $x'^\top \tilde{u} < 0$. In particular, by definition of $x^{(1)}$ and $x^{(2)}$, we see that $(x^{(1)})^\top \tilde{u} > 0$ and $(x^{(2)})^\top \tilde{u} < 0$. Then,

$$\left((x^{(1)} - x^{(2)})^\top \tilde{u} \right)^2 \geq \left((x^{(1)})^\top \tilde{u} \right)^2 + \left((x^{(2)})^\top \tilde{u} \right)^2 \geq 2(1 - \kappa_*^{-1/2})^2.$$

This implies that

$$\alpha \geq 2(1 - \kappa_*^{-1/2})^2 \times \kappa_* = 2(\sqrt{\kappa_*} - 1)^2.$$

Let us now prove that $\alpha \geq 1$. Note that by assumption, \mathcal{X} spans \mathbb{R}^d , and in particular there exists $\tilde{u} \in \mathbb{R}^d$ and $x, x' \in \mathcal{X}$ such that $\max_{x \in \mathcal{X}} x^\top \tilde{u} > 0$ and $\min_{x \in \mathcal{X}} x^\top \tilde{u} \leq 0$. Thus, $\max_{x, x' \in \mathcal{X}} \left((x - x')^\top \tilde{u} \right)^2 \geq \max_{x \in \mathcal{X}} (x^\top \tilde{u})^2$. For any $\lambda > 0$, choosing $u = \lambda \tilde{u}$ in the definition of α implies that

$$\alpha \geq \frac{\lambda^2 \max_{x \in \mathcal{X}} (x^\top u)^2}{\max_{x \in \mathcal{X}} (\lambda z_x x^\top u + 1)^2}.$$

Letting λ go to infinity, we find that $\alpha \geq 1$.

Finally, we prove that $\alpha \leq 8(\kappa_* + 1)$. For all $u \in \mathbb{R}^d$, we see that

$$\frac{\max_{x, x' \in \mathcal{X}} \left((x - x')^\top u \right)^2}{\max_{x \in \mathcal{X}} (z_x x^\top u + 1)^2} \leq \frac{4 \max_{x \in \mathcal{X}} (z_x x^\top u)^2}{\max_{x \in \mathcal{X}} (z_x x^\top u + 1)^2}.$$

Now, we see that

$$\frac{\max_{x \in \mathcal{X}} (z_x x^\top u)^2}{\max_{x \in \mathcal{X}} (z_x x^\top u + 1)^2} \leq \frac{2 \max_{x \in \mathcal{X}} (z_x x^\top u + 1)^2 + 2}{\max_{x \in \mathcal{X}} (z_x x^\top u + 1)^2} \leq 2 + \frac{2}{\max_{x \in \mathcal{X}} (z_x x^\top u + 1)^2}.$$

This in turn implies that for all $u \in \mathbb{R}^d$,

$$\frac{\max_{x, x' \in \mathcal{X}} \left((x - x')^\top u \right)^2}{\max_{x \in \mathcal{X}} (z_x x^\top u + 1)^2} \leq 8(1 + \kappa_*),$$

which finally implies that $\alpha \leq 8(1 + \kappa_*)$.

C.6.6 Proof of Lemma 8

Proof of Claim i) The proof of the first claim is immediate by definition of κ . Indeed, let $\widetilde{\mathcal{M}} = \left\{ \mu \in \mathcal{M}_{e_{d+1}}^{\mathcal{X}} : e_{d+1}^\top V(\mu)^+ e_{d+1} \leq 1 \right\}$ be the set of measures μ admissible for estimating ω^* with a precision level 1. Then,

$$\kappa(c\Delta) = \min_{\mu \in \widetilde{\mathcal{M}}} \sum_x \mu(x) c \Delta_x = c \min_{\mu \in \widetilde{\mathcal{M}}} \sum_x \mu(x) \Delta_x = c\kappa(\Delta).$$

Proof of Claim ii) The proof of the second claim is also straightforward. If $\Delta \leq \Delta'$, then for all $\mu \in \widetilde{\mathcal{M}}$, $\sum_x \mu(x) \Delta_x \leq \sum_x \mu(x) \Delta'_x$. Recall that $\mu^{\Delta'} = \operatorname{argmin}_{\mu \in \widetilde{\mathcal{M}}} \sum_x \mu(x) \Delta'_x$. Then,

$$\kappa(\Delta') = \sum_x \mu^{\Delta'}(x) \Delta'_x \geq \sum_x \mu^{\Delta'}(x) \Delta_x \geq \min_{\mu \in \widetilde{\mathcal{M}}} \sum_x \mu(x) \Delta_x = \kappa(\Delta).$$

Proof of Claim iii) To prove the third claim, note that

$$\begin{aligned}
\kappa(\Delta \vee \Delta') &= \min_{\mu \in \widetilde{\mathcal{M}}} \sum_x \mu(x) (\Delta_x \vee \Delta'_x) \\
&\geq \min_{\mu \in \widetilde{\mathcal{M}}} \left(\sum_x \mu(x) \Delta_x \vee \sum_x \mu(x) \Delta'_x \right) \\
&\geq \left(\min_{\mu \in \widetilde{\mathcal{M}}} \sum_x \mu(x) \Delta_x \right) \vee \left(\min_{\mu \in \widetilde{\mathcal{M}}} \sum_x \mu(x) \Delta'_x \right) \\
&\geq \kappa(\Delta) \vee \kappa(\Delta').
\end{aligned}$$

Proof of Claim iv) Recall that

$$\kappa(\Delta) = \min_{\mu \in \widetilde{\mathcal{M}}} \sum_x \mu(x) \Delta_x.$$

Let us define a sequence $(\mu_n)_{n \in \mathbb{N}} \in \widetilde{\mathcal{M}}^{\mathbb{N}}$ such that $\sum_x \mu_n(x) \Delta_x \xrightarrow{n \rightarrow \infty} \kappa(\Delta)$, and let us denote $\kappa_n = \sum_x \mu_n(x) \Delta_x$. According to Claim ii), we have

$$\kappa(\Delta) \leq \kappa(\Delta \vee \epsilon) = \min_{\mu \in \widetilde{\mathcal{M}}} \sum_x \mu(x) (\Delta_x \vee \epsilon) \leq \sum_x \mu_n(x) \Delta_x + \epsilon \sum_x \mu_n(x).$$

It follows that for all n ,

$$\kappa(\Delta) \leq \liminf_{\epsilon \rightarrow 0^+} \kappa(\Delta \vee \epsilon) \leq \limsup_{\epsilon \rightarrow 0^+} \kappa(\Delta \vee \epsilon) \leq \kappa_n.$$

Letting n go to infinity, we get that $\lim_{\epsilon \rightarrow 0^+} \kappa(\Delta \vee \epsilon) = \kappa(\Delta)$.

C.6.7 Proof of Lemma 9

Setting $\mu \cdot \Delta = (\mu(x) \Delta_x)_{x \in \mathcal{X}}$ and

$$V_{\Delta}(\lambda) = \sum_{x \in \mathcal{X}} \lambda_x \begin{pmatrix} \Delta_x^{-1/2} x \\ \Delta_x^{-1/2} z_x \end{pmatrix} \begin{pmatrix} \Delta_x^{-1/2} x \\ \Delta_x^{-1/2} z_x \end{pmatrix}^{\top},$$

we observe that $V_{\Delta}(\mu \cdot \Delta) = V(\mu)$. Hence,

$$\kappa(\Delta) = \min_{\substack{\mu \in \mathcal{M}^+ \\ e_{d+1}^{\top} V_{\Delta}(\mu \cdot \Delta) + e_{d+1} \leq 1}} \sum_{x \in \mathcal{X}} (\mu \cdot \Delta)_x.$$

We observe that $e_{d+1} \in \text{Range}(V(\mu))$ is equivalent to $e_{d+1} \in \text{Range}(V_{\Delta}(\mu \cdot \Delta))$. Hence, $\mu^{\Delta} \cdot \Delta = \lambda^{\Delta}$ where

$$\lambda^{\Delta} \in \underset{\substack{\lambda \in \mathbb{R}_+^{\mathcal{X}} \\ e_{d+1} \in \text{Range}(V_{\Delta}(\lambda)) \\ e_{d+1}^{\top} V_{\Delta}(\lambda) + e_{d+1} \leq 1}}{\text{argmin}} \sum_{x \in \mathcal{X}} \lambda_x.$$

The conclusion then follows by noticing that by homogeneity, $\lambda^{\Delta} = \kappa^{\Delta} \pi^{\Delta}$.

C.6.8 Proof of Lemma 10

Lemma 10 follows directly from Lemmas 18 and 19.

Lemma 18.

$$\mathbb{P} \left(\exists l \geq 1, z \in \{-1, 1\} \text{ such that } \text{Explore}_l^{(z)} = \text{True}, \text{ and } x \in \mathcal{X}_l^{(z)} \text{ such that } \left| \begin{pmatrix} \widehat{\gamma}_l^{(z)} - \gamma^* \\ \widehat{\omega}_l^{(z)} - \omega^* \end{pmatrix}^{\top} \begin{pmatrix} x \\ z_x \end{pmatrix} \right| \geq \epsilon_l \right) \leq \delta.$$

Lemma 19.

$$\mathbb{P} \left(\exists l \geq 1 \text{ such that } \text{Explore}_l^{(0)} = \text{True} \text{ and } \left| \widehat{\omega}_l^{(0)} - \omega^* \right| \geq \epsilon_l \right) \leq \delta.$$

C.6.9 Proof of Lemma 11

To prove Lemma 11, we rely on the following key lemma. This lemma proves that on $\overline{\mathcal{F}}$, i.e. when the error bounds hold, the algorithm never eliminates the best action or the best group.

Lemma 20. *On the event $\overline{\mathcal{F}}$, for all $x^* \in \operatorname{argmax}_{x \in \mathcal{X}} x^\top \gamma^*$ and all l such that $\text{Explore}_l^{(z_{x^*})} = \text{True}$, $x^* \in \mathcal{X}_{l+1}^{(z_{x^*})}$. Moreover, on the event $\overline{\mathcal{F}}$, for all l such that $\text{Explore}_l^{(0)} = \text{True}$, there exists $x^* \in \operatorname{argmax}_{x \in \mathcal{X}} x^\top \gamma^*$ such that $\widehat{z}_{l+1}^* \neq -z_{x^*}$.*

Let $l \geq 1$ be such that $\text{Explore}_l^{(z_{x^*})} = \text{True}$. Then, on $\overline{\mathcal{F}}$, $x^* \in \mathcal{X}_{l+1}^{(z_{x^*})}$ by Lemma 20. Moreover, for all $x \in \mathcal{X}_{l+1}^{(z_{x^*})}$, by definition of $\mathcal{X}_{l+1}^{(z_{x^*})}$, we have that on $\overline{\mathcal{F}}$

$$\left(\begin{pmatrix} x^* \\ z_{x^*} \end{pmatrix} - \begin{pmatrix} x \\ z_{x^*} \end{pmatrix} \right)^\top \begin{pmatrix} \widehat{\gamma}_l^{(z)} \\ \widehat{\omega}_l^{(z)} \end{pmatrix} \leq 3\epsilon_l.$$

which implies that

$$\left(\begin{pmatrix} x^* \\ z_{x^*} \end{pmatrix} - \begin{pmatrix} x \\ z_{x^*} \end{pmatrix} \right)^\top \begin{pmatrix} \gamma^* \\ \omega^* \end{pmatrix} \leq 3\epsilon_l + \left| \begin{pmatrix} x^* \\ z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_l^{(z)} - \gamma^* \\ \widehat{\omega}_l^{(z)} - \omega^* \end{pmatrix} \right| + \left| \begin{pmatrix} x \\ z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_l^{(z)} - \gamma^* \\ \widehat{\omega}_l^{(z)} - \omega^* \end{pmatrix} \right|.$$

Thus, on the event $\overline{\mathcal{F}}$, for all $x \in \mathcal{X}_{l+1}^{(z_{x^*})}$

$$(x^* - x)^\top \gamma^* < 5\epsilon_l,$$

which proves Equation (10). To prove the second claim of Lemma 11, assume that for all $x' \in \operatorname{argmax}_{x \in \mathcal{X}} x^\top \gamma^*$, $z_{x'} = z_{x^*}$ (when this does not hold, the second claim follows from Equation (10)). Now, let $l \geq 1$ be such that $\text{Explore}_l^{(-z_{x^*})} = \text{True}$. By Lemma 20, on $\overline{\mathcal{F}}$, $x^* \in \mathcal{X}_l^{(z_{x^*})}$ and $\widehat{z}_l^* = 0$. Then, the algorithm is unable to determine the group containing the best set during the phase $\text{Exp}_{l-1}^{(0)}$, so there must exist $x' \in \mathcal{X}_l^{(-z_{x^*})}$ such that

$$\begin{pmatrix} x^* \\ z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_{l-1}^{(z_{x^*})} \\ \widehat{\omega}_{l-1}^{(z_{x^*})} \end{pmatrix} \leq \begin{pmatrix} x' \\ -z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_{l-1}^{(-z_{x^*})} \\ \widehat{\omega}_{l-1}^{(-z_{x^*})} \end{pmatrix} + 2z_{x^*} \widehat{\omega}_{l-1}^{(0)} + 4\epsilon_{l-1}.$$

It follows that

$$\begin{pmatrix} x^* - x' \\ 2z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \gamma^* \\ \omega^* \end{pmatrix} \leq \begin{pmatrix} x^* \\ z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \gamma^* - \widehat{\gamma}_{l-1}^{(z_{x^*})} \\ \omega^* - \widehat{\omega}_{l-1}^{(z_{x^*})} \end{pmatrix} + \begin{pmatrix} x' \\ -z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_{l-1}^{(-z_{x^*})} - \gamma^* \\ \widehat{\omega}_{l-1}^{(-z_{x^*})} - \omega^* \end{pmatrix} + 2z_{x^*} \widehat{\omega}_{l-1}^{(0)} + 4\epsilon_{l-1}.$$

On $\overline{\mathcal{F}}$, this implies that

$$\begin{pmatrix} x^* - x' \\ 2z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \gamma^* \\ \omega^* \end{pmatrix} < 2z_{x^*} \widehat{\omega}_{l-1}^{(0)} + 6\epsilon_{l-1}$$

so

$$(x^* - x')^\top \gamma^* \leq 2z_{x^*} \left(\widehat{\omega}_{l-1}^{(0)} - \omega^* \right) + 6\epsilon_{l-1} < 8\epsilon_{l-1} = 16\epsilon_l. \quad (53)$$

Moreover, for all $x \in \mathcal{X}_{l+1}^{(-z_{x^*})}$ we have $(a_{x'} - a_x)^\top \widehat{\theta}_l^{(-z_{x^*})} \leq 3\epsilon_l$, so following the same lines as for the first claim, we get $(x' - x)^\top \gamma^* < 5\epsilon_l$. Combining this bound with (53), we get

$$\max_{x \in \mathcal{X}_{l+1}^{(-z_{x^*})}} (x^* - x)^\top \gamma^* < 21\epsilon_l.$$

This concludes the proof of Lemma 11.

C.6.10 Proof of Lemma 12

For $z \in \{-1, +1\}$ and $l > 0$,

$$\sum_x \mu_l^{(z)}(x) \leq \sum_x \frac{2(d+1)\pi_l^{(z)}(x)}{\epsilon_l^2} \log\left(\frac{kl(l+1)}{\delta}\right) + |\text{supp}(\pi_l^{(z)})|.$$

Now, $\text{supp}(\pi_l^{(z)}) \leq \frac{(d+1)(d+2)}{2}$ and $\sum_x \pi_l^{(z)}(x) = 1$, so

$$\sum_x \mu_l^{(z)}(x) \leq \frac{2(d+1)}{\epsilon_l^2} \log\left(\frac{kl(l+1)}{\delta}\right) + \frac{(d+1)(d+2)}{2}$$

which proves the first claim of Lemma 12.

To prove the second claim, we bound the regret for bias estimation at stage l as follows. On $\overline{\mathcal{F}}$, we have $\Delta_x \leq \widehat{\Delta}_x^l$ for all $x \in \mathcal{X}$ and $l \geq 1$, so

$$\sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \Delta_x \leq \sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \widehat{\Delta}_x^l.$$

Recall that $\hat{\mu}_l$ is the $\widehat{\Delta}^l$ -optimal design, and that for all $x \in \mathcal{X}$, $\mu_l^{(0)}(x) = \lceil \frac{2\hat{\mu}_l(x)}{\epsilon_l^2} \log\left(\frac{l(l+1)}{\delta}\right) \rceil$. Since $\widehat{\Delta}_x^l \leq 2$ for all $x \in \mathcal{X}$, we have

$$\sum_{x \in \mathcal{X}} \mu_l^{(0)}(x) \widehat{\Delta}_x^l \leq \sum_{x \in \mathcal{X}} \frac{2\hat{\mu}_l(x)}{\epsilon_l^2} \log\left(\frac{l(l+1)}{\delta}\right) \widehat{\Delta}_x^l + 2|\text{supp}(\mu_l^{(0)})|$$

and $|\text{supp}(\mu_l^{(0)})| \leq d+1$, so

$$\sum_x \mu_l^{(0)}(x) \Delta_x \leq \frac{2}{\epsilon_l^2} \log\left(\frac{l(l+1)}{\delta}\right) \sum_{x \in \mathcal{X}} \hat{\mu}_l(x) \widehat{\Delta}_x^l + 2(d+1).$$

By definition of $\hat{\mu}_l(x)$, we have that

$$\sum_{x \in \mathcal{X}} \hat{\mu}_l(x) \widehat{\Delta}_x^l = \kappa(\widehat{\Delta}^l).$$

It follows that, on $\overline{\mathcal{F}}$,

$$\sum_x \mu_l^{(0)}(x) \Delta_x \leq \sum_x \mu_l^{(0)}(x) \widehat{\Delta}_x^l \leq \frac{2}{\epsilon_l^2} \log\left(\frac{l(l+1)}{\delta}\right) \kappa(\widehat{\Delta}^l) + 2(d+1).$$

C.6.11 Proof of Lemma 13

For the first claim, we rely on the next lemma.

Lemma 21. *Let us set $\ell_x = \max\{l \geq 1 : x \in \mathcal{X}_l^{(-1)} \cup \mathcal{X}_l^{(1)}\}$. On $\overline{\mathcal{F}}$, we have for any $l \geq 1$*

1. $\widehat{\Delta}_x^l \leq \Delta_x + 16\epsilon_l$ for all $x \in \mathcal{X}_l^{(-1)} \cup \mathcal{X}_l^{(1)}$ (i.e. for all x such that $l \leq \ell_x$);
2. if $\Delta_x \geq 21\epsilon_l$ then $\ell_x \leq l$;
3. $\epsilon_{\ell_x} < \Delta_x$ for all $x \in \mathcal{X}$.

Lemma 13 relies on the following remarks : if Δ, Δ' are such that $\Delta_x \leq \Delta'_x$ for all $x \in \mathcal{X}$, then by Lemma 8 (ii), $\kappa(\Delta) \leq \kappa(\Delta')$. Let us now prove that for all $l \geq 1$ and all $x \in \mathcal{X}$, $\widehat{\Delta}_x^l \leq 513(\Delta \vee \epsilon_l)$.

Case $\epsilon_l \geq \Delta_x$. On $\overline{\mathcal{F}}$, we have $l \leq \ell_x - 1$ according to the third claim of Lemma 21. So, on $\overline{\mathcal{F}}$,

$$\widehat{\Delta}_x^l \leq \Delta_x + 16\epsilon_l \leq 17(\Delta_x \vee \epsilon_l).$$

Case $\epsilon_l < \Delta_x$. Then, on $\overline{\mathcal{F}}$, we have $32\epsilon_{l+5} < \Delta_x$ and so $l + 5 \geq \ell_x$ according to the second claim of Lemma 21. Hence, on $\overline{\mathcal{F}}$, according to Lemma 21, we have

$$\begin{aligned}\widehat{\Delta}_x^l &\leq \max_{k=0,\dots,5} \widehat{\Delta}_x^{\ell_x-k} \leq \Delta_x + 16\epsilon_{\ell_x-5} \\ &\leq \Delta_x + 512\epsilon_{\ell_x} \leq 513\Delta_x.\end{aligned}$$

Thus, for all $l \geq 1$ and all $x \in \mathcal{X}$,

$$\widehat{\Delta}_x^l \leq 513(\Delta \vee \epsilon_l).$$

Now, let $\widetilde{\mathcal{M}} = \left\{ \mu \in \mathcal{M}_{e_{d+1}}^{\mathcal{X}} : e_{d+1}^\top V(\mu)^+ e_{d+1} \geq 1 \right\}$ the measures μ admissible for estimating ω^* with a precision level 1. Note that for all $a, b, c > 0$,

$$(1 + ab^{-1})(c \vee b) = (c + cab^{-1}) \vee (a + b) \geq c \vee (a + b) \geq c \vee a. \quad (54)$$

Using Equation (54) with $a = \Delta_x$, $b = \tau$ and $c = \epsilon$, we see that

$$\kappa(\Delta \vee \epsilon) = \min_{\mu \in \widetilde{\mathcal{M}}} \sum_x \mu(x)(\Delta_x \vee \epsilon) \leq (1 + \epsilon/\tau) \min_{\mu \in \widetilde{\mathcal{M}}} \sum_x \mu(x)(\Delta_x \vee \tau) = (1 + \epsilon/\tau)\kappa(\Delta \vee \tau).$$

Using Lemma 8 together with $\widehat{\Delta}_x^l \leq 513(\Delta \vee \epsilon_l)$, we find that

$$\kappa(\widehat{\Delta}_x^l) \leq 513\kappa(\Delta \vee \epsilon_l) \leq 513(1 + \epsilon_l/\tau)\kappa(\Delta \vee \tau).$$

This proves the first claim of Lemma 13.

To prove the second claim, we use Lemma 8 and the fact that for all x , $\widehat{\Delta}_x^l \geq \epsilon_l$. Moreover, on $\overline{\mathcal{F}}$, $\widehat{\Delta}_x^l \geq \Delta_x$ for all $x \in \mathcal{X}$. Then, $\kappa(\widehat{\Delta}) \geq \kappa(\epsilon_l \vee \Delta)$ by Lemma 8 (iii).

C.6.12 Proof of Lemmas 14

To prove Lemma 14, let us consider l such that $\epsilon_l \leq \frac{\Delta_{\neq}}{8}$. According to Lemma 20, on $\overline{\mathcal{F}}$ we know that $\widehat{z}_l^* \neq -z_{x^*}$. When $\widehat{z}_l^* = z_{x^*}$, then we also have $\widehat{z}_{l+1}^* = z_{x^*}$ and the conclusion follows immediately. Let us consider now the case where $\widehat{z}_l^* = 0$. By definition of Δ_{\neq} , for all $x' \in \mathcal{X}_{l+1}^{(-z_{x^*})}$,

$$(x^* - x')^\top \gamma^* \geq \Delta_{\neq}.$$

This implies that

$$\begin{aligned}\begin{pmatrix} x^* \\ z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_l^{(z_{x^*})} \\ \widehat{\omega}_l^{(z_{x^*})} \end{pmatrix} - z_{x^*} \widehat{\omega}_l^{(0)} &\geq \max_{x \in \mathcal{X}_{l+1}^{(-z_{x^*})}} \begin{pmatrix} x \\ -z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_l^{(-z_{x^*})} \\ \widehat{\omega}_l^{(-z_{x^*})} \end{pmatrix} + z_{x^*} \widehat{\omega}_l^{(0)} \\ &+ \begin{pmatrix} x^* \\ z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_l^{(z_{x^*})} - \gamma^* \\ \widehat{\omega}_l^{(z_{x^*})} - \omega^* \end{pmatrix} + \min_{x \in \mathcal{X}_{l+1}^{(-z_{x^*})}} \begin{pmatrix} x \\ -z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \gamma^* - \widehat{\gamma}_l^{(-z_{x^*})} \\ \omega^* - \widehat{\omega}_l^{(-z_{x^*})} \end{pmatrix} \\ &+ \Delta_{\neq} + 2z_{x^*} (\omega^* - \widehat{\omega}_l^{(0)}).\end{aligned}$$

On $\overline{\mathcal{F}}$, it follows that

$$\begin{pmatrix} x^* \\ z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_l^{(z_{x^*})} \\ \widehat{\omega}_l^{(z_{x^*})} \end{pmatrix} - z_{x^*} \widehat{\omega}_l^{(0)} - 2\epsilon_l \geq \max_{x \in \mathcal{X}_{l+1}^{(-z_{x^*})}} \begin{pmatrix} x \\ -z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_l^{(-z_{x^*})} \\ \widehat{\omega}_l^{(-z_{x^*})} \end{pmatrix} + z_{x^*} \widehat{\omega}_l^{(0)} - 6\epsilon_l + \Delta_{\neq}.$$

When $\Delta_{\neq} \geq 8\epsilon_l$, this implies that $\widehat{z}_{l+1}^* = z_{x^*}$.

C.6.13 Proof of Lemmas 16 and 15

We prove Lemma 16. The proof of Lemma 15 follows by noticing that the two actions sets are equal up to a permutation of the direction of some basis vectors. To prove Lemma 15, we rely on Elfving's characterization of c -optimal design, given in Theorem 5. Theorem 5 shows that for $\pi \in \mathcal{P}^{\{1, \dots, d+1\}}$ to be e_{d+1} -optimal, there must exist $t > 0$ and $\zeta \in \{-1, +1\}^{d+1}$ such that

$$\begin{aligned} \sum_{1 \leq i \leq d+1} \pi_i &= 1 \\ 0 &= \pi_1 \zeta_1 - \left(1 - \frac{2}{\sqrt{\kappa_*} + 1}\right) \pi_{d+1} \zeta_{d+1} \\ \forall i \in \{2, \dots, d\}, 0 &= \pi_i \zeta_i \\ t &= \sum_{1 \leq i \leq \lfloor d/2 \rfloor} \pi_i \zeta_i - \sum_{\lfloor d/2 \rfloor + 1 \leq i \leq d+1} \pi_i \zeta_i. \end{aligned}$$

Solving this system, we find that $t^{-2} = \kappa_*$. Note that the unicity of the solution for the corresponding probability measure π guarantees that te_{d+1} belongs to the boundary of \mathcal{S} .

C.6.14 Proof of Lemma 17

For a given parameter γ^* , let us denote by Δ_i the gap corresponding to the action i . To compute $\kappa(\Delta)$, we could want to rely on Lemma 9 to find the Δ -optimal design, corresponding to the e_{d+1} -optimal design on the rescaled features $\Delta_x^{-1/2} \begin{pmatrix} x \\ z_x \end{pmatrix}$. Theorem 5 indeed allows us to compute such a design, as seen in the proof of Lemma 16. Unfortunately, we cannot rescale the features using the true gaps, since $\Delta_{x^*} = 0$. To circumvent this problem, we rely on the following reasoning :

1. We use Lemma 9 and Theorem 5 to compute the design $\mu^{\Delta \vee \epsilon}$ for $\epsilon \in (0, \Delta_{\min})$; and the corresponding regret $\kappa(\Delta \vee \epsilon)$;
2. We find the value of $\kappa(\Delta)$ by noticing that $\epsilon \mapsto \kappa(\Delta \vee \epsilon)$ is continuous at 0.

For $\epsilon \in (0, \Delta_{\min})$, define $\bar{\Delta} = \Delta \vee \epsilon$, and $\bar{x} = \bar{\Delta}_x^{-1/2} x$. Let $\bar{\pi}$ denote the e_{d+1} -optimal design for the rescaled features \bar{x} , and let $\bar{\kappa}_*$ denote its variance. Then, Lemma 9 ensures that $\kappa(\bar{\Delta}) = \bar{\kappa}_*$.

Now, Theorem 5 shows that there exists $\zeta \in \{-1, +1\}^{d+1}$ such that

$$\begin{aligned} \sum_{1 \leq i \leq d+1} \bar{\pi}_i &= 1 \\ 0 &= \bar{\pi}_1 \zeta_1 \bar{\Delta}_1^{-1/2} - \left(1 - \frac{2}{\sqrt{\bar{\kappa}_*} + 1}\right) \bar{\pi}_{d+1} \zeta_{d+1} \bar{\Delta}_{d+1}^{-1/2} \\ \forall i \in \{2, \dots, d\}, 0 &= \bar{\pi}_i \zeta_i \bar{\Delta}_i^{-1/2} \\ \bar{\kappa}_*^{-1/2} &= \sum_{1 \leq i \leq \lfloor d/2 \rfloor} \bar{\pi}_i \zeta_i \bar{\Delta}_i^{-1/2} - \sum_{\lfloor d/2 \rfloor + 1 \leq i \leq d+1} \bar{\pi}_i \zeta_i \bar{\Delta}_i^{-1/2} \end{aligned}$$

and $\bar{\kappa}_*^{-1/2} e_{d+1}$ belongs to the boundary of \mathcal{S} . Solving this system, we find that

$$\kappa(\bar{\Delta})^{-1/2} = \bar{\kappa}_*^{-1/2} = \frac{\left(\frac{2}{\sqrt{\bar{\kappa}_*} + 1}\right) \bar{\Delta}_{d+1}^{-1/2}}{1 + \left(1 - \frac{2}{\sqrt{\bar{\kappa}_*} + 1}\right) \bar{\Delta}_{d+1}^{-1/2} \bar{\Delta}_1^{-1/2}}.$$

As in Lemma 16, the unicity of the solution for the corresponding probability measure $\bar{\pi}$ guarantees that $\bar{\kappa}_*^{-1/2} e_{d+1}$ belongs to the boundary of the Elfving's set. Now, $\epsilon \leq \Delta_{\min}$, so

$$\kappa(\bar{\Delta})^{-1/2} = \kappa(\Delta \vee \epsilon)^{-1/2} = \frac{\left(\frac{2}{\sqrt{\kappa_*} + 1}\right) \Delta_{d+1}^{-1/2}}{1 + \left(1 - \frac{2}{\sqrt{\kappa_*} + 1}\right) \Delta_{d+1}^{-1/2} \epsilon^{1/2}}.$$

The fourth claim of Lemma 8 ensures that $\kappa(\Delta \vee \epsilon) \xrightarrow{\epsilon \rightarrow 0} \kappa(\Delta)$. Therefore,

$$\kappa(\Delta) = \lim_{\epsilon \rightarrow 0} \left(\frac{\left(\frac{2}{\sqrt{\kappa_*} + 1} \right) \Delta_{d+1}^{-1/2}}{1 + \left(1 - \frac{2}{\sqrt{\kappa_*} + 1} \right) \Delta_{d+1}^{-1/2} \epsilon^{1/2}} \right)^{-2} = \frac{(\sqrt{\kappa_*} + 1)^2 \Delta_{d+1}}{4}.$$

C.6.15 Proof of Lemma 18

Recall that $\xi_t = y_t - x_t^\top \gamma^* - z_{x_t} \omega^*$. For $l \geq 0$ and $z \in \{-1, +1\}$, when $\text{Explore}_l^{(z)} = \text{True}$, the least square estimator $\begin{pmatrix} \widehat{\gamma}_l^{(z)} \\ \widehat{\omega}_l^{(z)} \end{pmatrix}$ is given by

$$\begin{aligned} \begin{pmatrix} \widehat{\gamma}_l^{(z)} \\ \widehat{\omega}_l^{(z)} \end{pmatrix} &= \left(V_l^{(z)} \right)^+ \sum_{t \in \text{Exp}_l^{(z)}} \left(\begin{pmatrix} x_t \\ z_{x_t} \end{pmatrix}^\top \begin{pmatrix} \gamma^* \\ \omega^* \end{pmatrix} + \xi_t \right) \begin{pmatrix} x_t \\ z_{x_t} \end{pmatrix} \\ &= \left(V_l^{(z)} \right)^+ \left(V_l^{(z)} \right) \begin{pmatrix} \gamma^* \\ \omega^* \end{pmatrix} + \left(V_l^{(z)} \right)^+ \sum_{t \in \text{Exp}_l^{(z)}} \xi_t \begin{pmatrix} x_t \\ z_{x_t} \end{pmatrix}, \end{aligned}$$

where $\left(V_l^{(z)} \right)^+$ is a generalized inverse of $V_l^{(z)}$. Since $V_l^{(z)} \left(V_l^{(z)} \right)^+ V_l^{(z)} = V_l^{(z)}$, multiplying the left and right hand side of the last equation by $V_l^{(z)}$, we find that

$$V_l^{(z)} \begin{pmatrix} \widehat{\gamma}_l^{(z)} - \gamma^* \\ \widehat{\omega}_l^{(z)} - \omega^* \end{pmatrix} = V_l^{(z)} \left(V_l^{(z)} \right)^+ \sum_{t \in \text{Exp}_l^{(z)}} \xi_t \begin{pmatrix} x_t \\ z_{x_t} \end{pmatrix}. \quad (55)$$

By Lemma 4, for all $x \in \mathcal{X}_l^{(z)}$, $\begin{pmatrix} x \\ z_x \end{pmatrix} \in \text{Range} \left(V_l^{(z)} \right)$, so

$$V_l^{(z)} \left(V_l^{(z)} \right)^+ \begin{pmatrix} x \\ z_x \end{pmatrix} = \begin{pmatrix} x \\ z_x \end{pmatrix}. \quad (56)$$

Then,

$$\begin{aligned} \begin{pmatrix} \widehat{\gamma}_l^{(z)} - \gamma^* \\ \widehat{\omega}_l^{(z)} - \omega^* \end{pmatrix}^\top \begin{pmatrix} x \\ z_x \end{pmatrix} &= \begin{pmatrix} \widehat{\gamma}_l^{(z)} - \gamma^* \\ \widehat{\omega}_l^{(z)} - \omega^* \end{pmatrix}^\top V_l^{(z)} \left(V_l^{(z)} \right)^+ \begin{pmatrix} x \\ z_x \end{pmatrix} \\ &= \sum_{t \in \text{Exp}_l^{(z)}} \begin{pmatrix} x_t \\ z_{x_t} \end{pmatrix}^\top \left(V_l^{(z)} \right)^+ V_l^{(z)} \left(V_l^{(z)} \right)^+ \begin{pmatrix} x \\ z_x \end{pmatrix} \xi_t \\ &= \sum_{t \in \text{Exp}_l^{(z)}} \begin{pmatrix} x_t \\ z_{x_t} \end{pmatrix}^\top \left(V_l^{(z)} \right)^+ \begin{pmatrix} x \\ z_x \end{pmatrix} \xi_t, \end{aligned}$$

where the first and third lines follow from Equation (56), and the second line follows from Equation (55). By definition of our algorithm, conditionally on $\mathcal{X}_l^{(z)}$ and $\text{Explore}_l^{(z)} = \text{True}$, the variables $(\xi_t)_{t \in \text{Exp}_l^{(z)}}$ are independent centered normal gaussian variables. Then,

$$\mathbb{P}_{|\mathcal{X}_l^{(z)}, \text{Explore}_l^{(z)} = \text{True}} \left(\left| \begin{pmatrix} \widehat{\gamma}_l^{(z)} - \gamma^* \\ \widehat{\omega}_l^{(z)} - \omega^* \end{pmatrix}^\top \begin{pmatrix} x \\ z_x \end{pmatrix} \right| \geq \sqrt{2 \sum_{t \in \text{Exp}_l^{(z)}} \left(\begin{pmatrix} x_t \\ z_{x_t} \end{pmatrix}^\top \left(V_l^{(z)} \right)^+ \begin{pmatrix} x \\ z_x \end{pmatrix} \right)^2 \log \left(\frac{kl(l+1)}{\delta} \right)} \right) \leq \frac{\delta}{kl(l+1)}.$$

Expanding $\left(\begin{pmatrix} x_t \\ z_{x_t} \end{pmatrix}^\top (V_l^{(z)})^+ \begin{pmatrix} x \\ z_x \end{pmatrix}\right)^2 = \begin{pmatrix} x \\ z_x \end{pmatrix}^\top (V_l^{(z)})^+ \begin{pmatrix} x_t \\ z_{x_t} \end{pmatrix} \begin{pmatrix} x_t \\ z_{x_t} \end{pmatrix}^\top (V_l^{(z)})^+ \begin{pmatrix} x \\ z_x \end{pmatrix}$, and using the definition of $V_l^{(z)}$, we find that

$$\mathbb{P}_{|\mathcal{X}_l^{(z)}, \text{Explore}_l^{(z)}=\text{True}} \left(\left| \begin{pmatrix} \widehat{\gamma}_l^{(z)} - \gamma^* \\ \widehat{\omega}_l^{(z)} - \omega^* \end{pmatrix}^\top \begin{pmatrix} x \\ z_x \end{pmatrix} \right| \geq \sqrt{2 \begin{pmatrix} x \\ z_x \end{pmatrix}^\top (V_l^{(z)})^+ V_l^{(z)} (V_l^{(z)})^+ \begin{pmatrix} x \\ z_x \end{pmatrix} \log \left(\frac{kl(l+1)}{\delta} \right)} \right) \leq \frac{\delta}{kl(l+1)}$$

which in turn implies (using Equation (56))

$$\mathbb{P}_{|\mathcal{X}_l^{(z)}, \text{Explore}_l^{(z)}=\text{True}} \left(\left| \begin{pmatrix} \widehat{\gamma}_l^{(z)} - \gamma^* \\ \widehat{\omega}_l^{(z)} - \omega^* \end{pmatrix}^\top \begin{pmatrix} x \\ z_x \end{pmatrix} \right| \geq \sqrt{2 \left\| \begin{pmatrix} x \\ z_x \end{pmatrix} \right\|_{(V_l^{(z)})^+}^2 \log \left(\frac{kl(l+1)}{\delta} \right)} \right) \leq \frac{\delta}{kl(l+1)}$$

Now, using Lemma 4 and the definition of μ_l^z , we see that for all $x \in \mathcal{X}_l^{(z)}$,

$$\begin{pmatrix} x \\ z_x \end{pmatrix}^\top (V_l^{(z)})^+ \begin{pmatrix} x \\ z_x \end{pmatrix} \leq \frac{\epsilon_l^2}{2 \log(kl(l+1)/\delta)}.$$

Finally, for all $x \in \mathcal{X}_l^{(z)}$,

$$\begin{aligned} & \mathbb{P}_{|\mathcal{X}_l^{(z)}, \text{Explore}_l^{(z)}=\text{True}} \left(\left| \begin{pmatrix} \widehat{\gamma}_l^{(z)} - \gamma^* \\ \widehat{\omega}_l^{(z)} - \omega^* \end{pmatrix}^\top \begin{pmatrix} x \\ z_x \end{pmatrix} \right| \geq \epsilon_l \right) \\ & \leq \mathbb{P}_{|\mathcal{X}_l^{(z)}, \text{Explore}_l^{(z)}=\text{True}} \left(\left| \begin{pmatrix} \widehat{\gamma}_l^{(z)} - \gamma^* \\ \widehat{\omega}_l^{(z)} - \omega^* \end{pmatrix}^\top \begin{pmatrix} x \\ z_x \end{pmatrix} \right| \geq \sqrt{2 \left\| \begin{pmatrix} x \\ z_x \end{pmatrix} \right\|_{(V_l^{(z)})^+}^2 \log \left(\frac{kl(l+1)}{\delta} \right)} \right) \leq \frac{\delta}{kl(l+1)}. \end{aligned}$$

Integrating out the conditioning on the value of $\mathcal{X}_l^{(z)}$ and $\text{Explore}_l^{(z)}$ and using a union bound yields the desired result.

C.6.16 Proof of Lemma 19

The proof is similar to that of Lemma 18. If $\text{Explore}_l^{(0)} = \text{True}$, then $\widehat{\omega}_l$ is defined as

$$\widehat{\omega}_l^{(0)} = e_{d+1}^\top (V_l^{(0)})^+ \sum_{t \in \text{Exp}_l^{(0)}} \left(\begin{pmatrix} x_t \\ z_{x_t} \end{pmatrix}^\top \begin{pmatrix} \gamma^* \\ \omega^* \end{pmatrix} + \xi_t \right) \begin{pmatrix} x_t \\ z_{x_t} \end{pmatrix}.$$

Since $\begin{pmatrix} x \\ z_x \end{pmatrix}_{x \in \mathcal{X}}$ spans \mathbb{R}^{d+1} , μ is finite and $e_{d+1} \in \text{Range}(V(\hat{\mu}_l))$. Then, according to Lemma 3, for every round l , we have $e_{d+1} \in \text{Range}(V_l^{(0)})$, so $V_l^{(0)} (V_l^{(0)})^+ e_{d+1} = e_{d+1}$. This implies that

$$\widehat{\omega}_l^{(0)} - \omega^* = \sum_{t \in \text{Exp}_l^{(0)}} e_{d+1}^\top (V_l^{(0)})^+ \begin{pmatrix} x_t \\ z_{x_t} \end{pmatrix} \xi_t.$$

By definition of our algorithm, conditionally on $\text{Explore}_l^{(0)} = \text{True}$, the variables $(\xi_t)_{t \in \text{Exp}_l^{(0)}}$ are independent centered normal gaussian variables. Then,

$$\mathbb{P}_{|\text{Explore}_l^{(0)}=\text{True}} \left(\left| \widehat{\omega}_l^{(0)} - \omega^* \right| \geq \sqrt{2 \sum_{t \in \text{Exp}_l^{(0)}} \left(e_{d+1}^\top (V_l^{(0)})^+ \begin{pmatrix} x_t \\ z_{x_t} \end{pmatrix} \right)^2 \log \left(\frac{l(l+1)}{\delta} \right)} \right) \leq \frac{\delta}{l(l+1)}.$$

Using again $V_l^{(0)} \left(V_l^{(0)} \right)^+ e_{d+1} = e_{d+1}$ and the definition of $V_l^{(0)}$, we find that

$$\mathbb{P}_{|\text{Explore}_l^{(0)}=\text{True}} \left(\left| \widehat{\omega}_l^{(0)} - \omega^* \right| \geq \sqrt{2e_{d+1}^\top \left(V_l^{(0)} \right)^+ e_{d+1} \log \left(\frac{l(l+1)}{\delta} \right)} \right) \leq \frac{\delta}{l(l+1)}. \quad (57)$$

Now, Lemma 3 and the definition of $\mu_l^{(0)}$ imply that

$$e_{d+1}^\top \left(V_l^{(0)} \right)^+ e_{d+1} \leq \frac{\epsilon_l^2}{2 \log(l(l+1)/\delta)}.$$

Finally, Equation (57) implies that

$$\mathbb{P}_{|\text{Explore}_l^{(0)}=\text{True}} \left(\left| \widehat{\omega}_l^{(0)} - \omega^* \right| \geq \epsilon_l \right) \leq \frac{\delta}{l(l+1)}.$$

Using a union bound over the phases $\text{Exp}_l^{(0)}$ yields the result.

C.6.17 Proof of Lemma 20

To prove Lemma 20, we begin by showing that it is enough to prove that for $l \geq 1$,

$$\begin{aligned} \mathcal{F}_l \supset & \left\{ \exists x^* \in \underset{x \in \mathcal{X}}{\text{argmax}} x^\top \gamma^* : \text{Explore}_l^{(z_{x^*})} = \text{True} \text{ and } x^* \notin \mathcal{X}_{l+1}^{(z_{x^*})} \right\} \\ \cup & \left\{ \overline{\bigcap_{l' \leq l} \left\{ \exists x^* \in \underset{x \in \mathcal{X}}{\text{argmax}} x^\top \gamma^* : \text{Explore}_{l'}^{(z_{x^*})} = \text{True} \text{ and } x^* \notin \mathcal{X}_{l'+1}^{(z_{x^*})} \right\}} \right. \\ & \left. \bigcap \left\{ \text{Explore}_l^{(0)} = \text{True} \text{ and } \forall x^* \in \underset{x \in \mathcal{X}}{\text{argmax}} x^\top \gamma^*, \widehat{z}_{l+1}^* = -z_{x^*} \right\} \right\}. \end{aligned} \quad (58)$$

Indeed, denoting $\mathcal{F}_l^{(1)} = \left\{ \exists x^* \in \underset{x \in \mathcal{X}}{\text{argmax}} x^\top \gamma^* : \text{Explore}_l^{(z_{x^*})} = \text{True} \text{ and } x^* \notin \mathcal{X}_{l+1}^{(z_{x^*})} \right\}$ and $\mathcal{F}_l^{(2)} = \left\{ \text{Explore}_l^{(0)} = \text{True} \text{ and } \forall x^* \in \underset{x \in \mathcal{X}}{\text{argmax}} x^\top \gamma^*, \widehat{z}_{l+1}^* = -z_{x^*} \right\}$, we see that Equation (58) would then be rewritten as

$$\mathcal{F}_l \supset \mathcal{F}_l^{(1)} \cup \left\{ \bigcap_{l' \leq l} \overline{\mathcal{F}_{l'}^{(1)}} \cap \mathcal{F}_l^{(2)} \right\}$$

which implies

$$\bigcup_{l \geq 1} \mathcal{F}_l \supset \bigcup_{l \geq 1} \left\{ \mathcal{F}_l^{(1)} \cup \left\{ \bigcap_{l' \leq l} \overline{\mathcal{F}_{l'}^{(1)}} \cap \mathcal{F}_l^{(2)} \right\} \cup \mathcal{F}_{l'}^{(1)} \right\} \supset \bigcup_{l \geq 1} \left\{ \mathcal{F}_l^{(1)} \cup \mathcal{F}_l^{(2)} \right\}.$$

Then, Equation (58) would imply that

$$\overline{\mathcal{F}} = \overline{\bigcup_{l \geq 1} \mathcal{F}_l} \subset \overline{\bigcup_{l \geq 1} \left\{ \mathcal{F}_l^{(1)} \cup \mathcal{F}_l^{(2)} \right\}} = \bigcap_{l \geq 1} \left\{ \overline{\mathcal{F}_l^{(1)}} \cap \overline{\mathcal{F}_l^{(2)}} \right\},$$

thus proving Lemma 20. To prove Equation (58), we show that both $\mathcal{F}_l^{(1)}$ and $\bigcap_{l' \leq l} \overline{\mathcal{F}_{l'}^{(1)}} \cap \mathcal{F}_l^{(2)}$ imply \mathcal{F}_l .

If $\mathcal{F}_l^{(1)}$ is true: then $\exists x^* \in \underset{x \in \mathcal{X}}{\text{argmax}} x^\top \gamma^* : \text{Explore}_l^{(z_{x^*})} = \text{True}$ and $x^* \notin \mathcal{X}_{l+1}^{(z_{x^*})}$.

Without loss of generality, assume that $l > 1$ is the smallest integer such that $\text{Explore}_l^{(z_{x^*})} = \text{True}$ and

$x^* \notin \mathcal{X}_{l+1}^{(z_{x^*})}$. Then, necessarily $x^* \in \mathcal{X}_l^{(z_{x^*})}$ (because either $l = 1$, or $\text{Explore}_{l-1}^{(z_{x^*})} = \text{True}$). Now, because $x^* \in \mathcal{X}_l^{(z_{x^*})} \setminus \mathcal{X}_{l+1}^{(z_{x^*})}$, there exists $x \in \mathcal{X}_l^{(z_{x^*})}$ such that

$$(x - x^*)^\top \widehat{\gamma}_l^{(z_{x^*})} \geq 3\epsilon_l$$

and in particular

$$x^\top \widehat{\gamma}_l^{(z_{x^*})} - \epsilon_l > (x^*)^\top \widehat{\gamma}_l^{(z_{x^*})} + \epsilon_l.$$

Recall that by definition of x^* , $(\gamma^*)^\top (x^* - x) \geq 0$. This in turn implies that

$$\begin{pmatrix} x \\ z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_l^{(z_{x^*})} - \gamma^* \\ \widehat{\omega}_l^{(z_{x^*})} - \omega^* \end{pmatrix} - \epsilon_l > \begin{pmatrix} x^* \\ z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_l^{(z_{x^*})} - \gamma^* \\ \widehat{\omega}_l^{(z_{x^*})} - \omega^* \end{pmatrix} + \epsilon_l.$$

The last equation implies that either $\begin{pmatrix} x \\ z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \gamma_l^{(z)} - \gamma^* \\ \widehat{\omega}_l^{(z)} - \omega^* \end{pmatrix} > \epsilon_l$ or $\begin{pmatrix} x^* \\ z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \gamma_l^{(z)} - \gamma^* \\ \widehat{\omega}_l^{(z)} - \omega^* \end{pmatrix} < -\epsilon_l$, which in turn implies \mathcal{F}_l .

If $\bigcap_{l' \leq l} \overline{\mathcal{F}_{l'}^{(1)}} \cap \mathcal{F}_l^{(2)}$ is true: then $\text{Explore}_l^{(0)} = \text{True}$ and $\forall x^* \in \text{argmax}_{x \in \mathcal{X}} x^\top \gamma^*, \widehat{z}_{l+1}^* = -z_{x^*}$. Moreover, for all $l' \leq l$, $\text{Explore}_{l'}^{(z_{x^*})} = \text{False}$ or $x^* \in \mathcal{X}_{l'+1}^{(z_{x^*})}$.

Note that this case can only hold if all optimal actions x^* belong to the same group z_{x^*} . Without loss of generality, assume that $l > 1$ is the smallest integer such that $\text{Explore}_l^{(0)} = \text{True}$ and $\widehat{z}_{l+1}^* = -z_{x^*}$, and for all $l' \leq l$, $\text{Explore}_{l'}^{(z_{x^*})} = \text{False}$ or $x^* \in \mathcal{X}_{l'+1}^{(z_{x^*})}$. Note that because $\text{Explore}_l^{(0)} = \text{True}$, necessarily $\text{Explore}_{l'}^{(z_{x^*})} = \text{True}$ for all $l' \leq l$, and in particular $x^* \in \mathcal{X}_{l+1}^{(z_{x^*})}$.

Then, there exists $x \in \mathcal{X}_{l+1}^{(-z_{x^*})}$ such that

$$\begin{pmatrix} x \\ -z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_l^{(-z_{x^*})} \\ \widehat{\omega}_l^{(-z_{x^*})} \end{pmatrix} - \begin{pmatrix} x^* \\ z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_l^{(z_{x^*})} \\ \widehat{\omega}_l^{(z_{x^*})} \end{pmatrix} + 2z_{x^*} \widehat{\omega}_l^{(0)} \geq 4\epsilon_l.$$

Recall that all optimal actions x^* are in the same group z_{x^*} , so $(\gamma^*)^\top (x^* - x) > 0$. This in turn implies that

$$\begin{pmatrix} x \\ -z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_l^{(-z_{x^*})} - \gamma^* \\ \widehat{\omega}_l^{(-z_{x^*})} - \omega^* \end{pmatrix} - \begin{pmatrix} x^* \\ z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_l^{(z_{x^*})} - \gamma^* \\ \widehat{\omega}_l^{(z_{x^*})} - \omega^* \end{pmatrix} + 2z_{x^*} (\widehat{\omega}_l^{(0)} - \omega^*) \geq 4\epsilon_l.$$

The last equation implies that either $\begin{pmatrix} x \\ -z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_l^{(-z_{x^*})} - \gamma^* \\ \widehat{\omega}_l^{(-z_{x^*})} - \omega^* \end{pmatrix} \geq \epsilon_l$, or $\begin{pmatrix} x^* \\ z_{x^*} \end{pmatrix}^\top \begin{pmatrix} \widehat{\gamma}_l^{(z_{x^*})} - \gamma^* \\ \widehat{\omega}_l^{(z_{x^*})} - \omega^* \end{pmatrix} \leq -\epsilon_l$, or $z_{x^*} (\widehat{\omega}_l^{(0)} - \omega^*) \geq \epsilon_l$, which in turn implies \mathcal{F}_l .

C.6.18 Proof of Lemma 21

The first claim holds for $l = 1$. For $l \geq 1$, for any $x \in \mathcal{X}_{l+1}^{(-1)} \cup \mathcal{X}_{l+1}^{(1)}$, we have $\widehat{\Delta}_x^{l+1} \leq \Delta_x + 8\epsilon_l$ on $\overline{\mathcal{F}}$ according to the definition of $\widehat{\Delta}^{l+1}$ and \mathcal{F} . The first claim then follows.

For the second claim, Lemma 11 gives that, on $\overline{\mathcal{F}}$, $\Delta_x < 21\epsilon_l$ for any $x \in \mathcal{X}_{l+1}^{(-1)} \cup \mathcal{X}_{l+1}^{(1)}$. So $\Delta_x \geq 21\epsilon_l$ implies $x \notin \mathcal{X}_{l+1}^{(-1)} \cup \mathcal{X}_{l+1}^{(1)}$ and hence $l \geq \ell_x$ on $\overline{\mathcal{F}}$.

For the third claim, we notice that

$$\max_{x' \in \mathcal{X}_{\ell_x}^{(z_x)}} (a_{x'} - a_x)^\top \widehat{\theta}_{\ell_x}^{(z_x)} > 3\epsilon_{\ell_x},$$

since $x \notin \mathcal{X}_{\ell_x+1}$. Since the left-hand side is smaller than $\Delta_x + 2\epsilon_{\ell_x}$ on $\overline{\mathcal{F}}$, we get $\Delta_x > \epsilon_{\ell_x}$.