

---

# Learning to Explore with Lagrangians for Bandits under Unknown Constraints

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Pure exploration in bandits can model eclectic real-world decision making prob-  
2 lems, such as tuning hyper-parameters or conducting user studies, where sample  
3 frugality is desired. Thus, considering different safety, resource, and fairness  
4 constraints on the decision space has gained increasing attention. In this paper,  
5 we study generalisation of these problems as pure exploration in multi-armed  
6 bandits with unknown linear constraints. First, we propose a Lagrangian relaxation  
7 of the sample complexity lower bound for pure exploration. We further derive  
8 how this lower bound converges to the existing lower bound for pure exploration  
9 under known constraints, and how the hardness of the problem changes with the  
10 geometry induced by the constraint estimation procedure. We further leverage the  
11 Lagrangian lower bound and properties of convex optimisation to propose two  
12 computationally efficient extensions of Track-and-Stop and Gamified Explorations,  
13 namely LATS and LAGEX. Designing these algorithms require us to propose a  
14 new constraint-adaptive stopping rule, and also at each step, using pessimistic  
15 estimates of constraints in the Lagrangian lower bound. We show that these algo-  
16 rithms asymptotically achieve the desired sample complexity bounds. Finally, we  
17 conduct numerical experiments with different reward distributions and constraints  
18 that validate efficient performance of LAGEX and LATS with respect to baselines.

## 19 1 Introduction

20 Decision-making under uncertainty is a ubiquitous challenge encountered across various domains,  
21 including clinical trials [VBW15], recommendation systems [ZY24] and more. Multi-Armed Bandit  
22 (MAB) serves as an archetypal framework for sequential decision-making under uncertainty and  
23 allows to study the involved information-utility trade-offs [LS20]. In MAB, at each step, an agent  
24 interacts with an environment consisting of  $K$  decisions (aka arms) corresponding to  $K$  noisy  
25 feedback distribution (aka reward distribution). At each step, the agent takes a decision, and obtains  
26 a reward from the corresponding reward distribution. The goal of the agent is to compute a *policy*,  
27 i.e. a distribution over the decisions, that maximises a certain utility metric (e.g. accumulated  
28 rewards [ACBF02], probability of identifying the best arm [KCG16] etc.) over time.

29 In this paper, we focus on the *pure exploration* problem of MABs, where the agent interacts by  
30 sequentially realising a sequence of policies (or experiments) with the goal of *answering a query*  
31 *as correctly as possible*. A well-studied pure exploration problem is Best-Arm Identification (BAI),  
32 where the agent aims to identify the arm with the highest expected reward [EDMM02a, BMS09, JN14,  
33 KCG16]. BAI has been increasingly applied for hyper-parameter tuning [LJD<sup>+</sup>17], communication  
34 networks [LPJ22], influenza mitigation [LVR<sup>+</sup>19], finding the optimal dose of a drug [AKR21a] etc.  
35 However, real-world scenarios often involve constraints on the arms that must be satisfied [CBJD24].

36 **Example 1** (Optimal treatment plan [KCJ23, CGS22a]). *We want to identify the optimal treatment of*  
37 *a patient with rheumatoid arthritis, when the first and second-line of treatments have failed [KCJ23].*

There is large variability in the choice of next treatment, and several drugs are a priori considered to be equally good but might not work equally well together [KCJ23]. Let us assume that we have  $K$  drugs for efficacies and side-effects are unknown. We assume that efficacy of each drug comes from a reward distribution with unknown mean  $\mu_a$  and also  $d$  side-effects (e.g. drop in heart- and liver-function scores) are due to unknown drug-specific constraints  $\mathbf{A}_a$ . The constraint  $\mathbf{A}_a$  represent scores that are deemed to be unsafe, and thus, we want  $\mathbf{A}\pi \leq \mathbf{0}$ . Thus, given a treatment plan  $\pi$  of mixing the drugs, the mean efficacy is  $\langle \mu, \pi \rangle$ , and the extent of side-effects is  $\mathbf{A}\pi$ . These scores and the efficacy serves can be measured after each application of a drug, with the stochasticity due to inter-patient variability. Additionally, changing a treatment has additional human involvement and cost, which is better to be minimised. We aim to reliably find the most effective drug cocktail  $\pi^*$  with minimum number of interactions with patients, while retaining the side-effects in the safe zone.

**Pure exploration under constraints.** In recent years, real-life applications like above naturally motivated study of pure exploration under a set of known and unknown constraints [KSS18, WWJ21, LZYL23, WZZ23, CBJD24]. Specifically, we aim to find the optimal policy that maximises the expected rewards over the set of arms and also satisfies the true constraints. The agent, at each step  $t$ , selects an action according to a chosen policy observes associated reward and the cost incurred with respect to the constraints. The agent uses the feedbacks to update the estimates of the expected rewards and constraints. Using these estimates, the agent chooses the next policy and action till the optimal policy satisfying constraints is identified with confidence  $1 - \delta$ . This is known as the *fixed-confidence setting* of pure exploration [WWJ21, CBJD24], while there is also *fixed-budget setting* which is of independent interest [KSS18, LZYL23, FN22]. Existing literature has studied either the general linear constraints when they are known [WWJ21, CBJD24, CWM<sup>+</sup>22a], or very specific type of unknown constraints, e.g. safety [WWJ21], knapsack [LZYL23], fairness [WZZ23], preferences [LTHK22] etc. Here, we study the *pure exploration problem in the fixed-confidence setting subject to unknown linear constraints on the policy*, which generalises all these settings (Section 2). Further discussions on related works is in Appendix B.1.

Recently, [CBJD23] show if the constraints are known, i.e the feasible set is deterministic, a bandit instance may become harder or easier depending on the geometry of the constraints. They state that *studying similar phenomenon for unknown constraints as an open problem* as the feasible policy space is not deterministic anymore. As we have to construct an estimate of the feasible policy set, we have to simultaneously control concentration of the means of the unknown reward distributions, and those of the feasible sets, simultaneously. This leads us to two questions

- How does the hardness of the pure exploration under unknown constraints change if the constraints are estimated sequentially?
- How can we design a generic algorithmic scheme to track both the constraints and the optimal policy with sample- and computational-efficiency?

**Our contributions** address these questions as follows:

1. *Lagrangian relaxation of the lower bound.* The minimum number of samples required to conduct pure exploration with fixed confidence and the corresponding interactive policy are expressed through lower bound, which is an optimisation problem under known constraints. For unknown constraints, we propose a novel Lagrangian relaxation of this optimisation problem in the lower bound (Sec. 3). At every step, we use pessimistic estimates of the constraints, while the Lagrangian multipliers help us to track the interaction between the objective function and the structure of the estimated constraints. We use multiple results from convex analysis to show that the Lagrangian relaxation with pessimistic constraints preserves all the continuity properties of the lower bound under known constraints. Additionally, it satisfies strong duality to yield a unique optimal policy for interactions, and also bounds on the Lagrangian multipliers at every step. Further, we characterise the Lagrangian lower bound for Gaussian rewards, which connects with the lower bound for known constraints.

2. *A generic algorithm design.* Now, we leverage this Lagrangian lower bound with pessimistic estimates of constraints to propose two algorithms, namely LATS (Lagrangian Track and Stop) and LAGEX (Lagrangian Gamified EXplorer). First, we develop a new stopping rule for the unknown constraint setting as we have to ensure that both the mean estimates and pessimistic estimates of constraints concentrate close to their true values before reliably recommend an optimal policy using them. Then, we extend the Track-and-Stop [GK16] and gamified explorer [DKM19b] approaches for the Lagrangian lower bound to design LATS and LAGEX, respectively (Sec. 4).

3. *Upper bound on sample complexities.* We provide upper bounds on sample complexities of LATS and LAGEX (Sec. 3). This requires proving a novel concentration of the constraints, and also

consequent concentration of optimal policies under constraints. We show that due to constraint LATS achieve an upper bound, which is  $(1 + \mathfrak{s})$  times the upper bound of TS under known constraints.  $\mathfrak{s}$  is the shadow price of the true constraint and quantifies its stability under perturbation. In contrast, LAGEX leads to an upper bound that has only an additive  $\mathfrak{s}$  factor with the known constraint lower bound. This suggests that LAGEX should be more sample-efficient than LATS. Our experimental results (Sec. 5) across multiple settings validate that LAGEX requires the least samples among competing algorithms and it can exactly follow the hardness due to constraints across environments.

## 2 Pure exploration under unknown constraints

**Notation.**  $x, \mathbf{x}, \mathbf{X}$ , adn  $\mathcal{X}$  denote a scalar, a vector, a matrix, and a set respectively. For a positive semi-definite matrix  $\mathbf{A}$  and vector  $\mathbf{z}$ ,  $\|\mathbf{z}\|_{\mathbf{A}}^2 = \langle \mathbf{z}, \mathbf{A}\mathbf{z} \rangle$ . Also, in  $\mathbb{R}_+^d$ , we include  $\mathbf{0}^d$ .  $[K]$  refers to  $\{1, \dots, K\}$ .  $\text{Supp}(P)$  denotes the support of a distribution  $P$ .  $\Delta_K$  is the simplex over  $[K]$ .

**Problem formulation.** We work with a MAB instance with  $K \in \mathbb{N}$  arms. Each arm  $a \in [K]$  has a reward distribution  $P_a$  with unknown means  $\mu_a \in \mathbb{R}$ . The agent, at each time step  $t \in \mathbb{N}$ , chooses an action  $A_t \in [K]$ , and observes a stochastic reward  $R_t \sim P_{A_t}$ . A feasible policy  $\pi \in \Delta_K$  satisfies  $\mathbf{A}\pi \leq \mathbf{0}$  with respect to the set of  $d$  linear constraints  $\mathbf{A} \in \mathbb{R}^{d \times K^1}$ .

If we have known  $\mathbf{A}$ , the agent would have access to the non-empty and compact set of feasible policies  $\mathcal{F} \triangleq \{\pi \in \Delta_K : \mathbf{A}\pi \leq \mathbf{0}\}$  and the agent would aim to identify the optimal feasible policy, i.e. the one yielding the highest expected reward while maintaining the constraints,

$$\pi_{\mathcal{F}}^* \triangleq \arg \max_{\pi \in \mathcal{F}} \mu^T \pi. \quad (1)$$

In our setting, we do not have access to the true set of constraints. At any time  $t \in \mathbb{N}$ , we rather produce  $\hat{\mathbf{A}}_t$  as an estimate of  $\mathbf{A}$ . Then, the agent has access to an estimated feasible set  $\hat{\mathcal{F}} \triangleq \{\pi \in \Delta_K : \hat{\mathbf{A}}\pi \leq \mathbf{0}\}$  and can identify the optimal feasible policy to be  $\pi_{\hat{\mathcal{F}}}^* \triangleq \arg \max_{\pi \in \hat{\mathcal{F}}} \mu^T \pi$ .

Using the estimated constraints, we want to design an algorithm that finally recommends policy that is  $(1 - \delta)$ -correct and  $(1 - \delta)$ -feasible.

**Definition 1** ( $(1 - \delta)$ -correct and  $(1 - \delta)$ -feasible recommended policy). *For  $\delta \in [0, 1)$ , a policy recommended by a pure exploration algorithm is  $(1 - \delta)$ -correct and  $(1 - \delta)$ -feasible if  $\Pr[\pi_{\hat{\mathcal{F}}}^* \neq \pi_{\mathcal{F}}^*] \leq \delta$  and  $\Pr[\mathbf{A}\pi_{\hat{\mathcal{F}}}^* \geq \mathbf{0}] \leq \delta$ .*

For the algorithm to yield correct and feasible recommendations at anytime it stops, we want such an estimate of  $\mathbf{A}$  that  $\hat{\mathcal{F}}$  is a superset of  $\mathcal{F}$ . Otherwise, we cannot ensure that the true optimal and feasible policy  $\pi_{\mathcal{F}}^* \in \hat{\mathcal{F}}$ , i.e. the set where we are searching for the recommended policy. Additionally, we might hit a degenerate case where  $\pi_{\hat{\mathcal{F}}}^*$  might not exist. To ensure these properties, in Section 3, we design pessimistic estimates of  $\mathbf{A}$  and use them further. This property echoes the spirit of the optimistic-pessimistic algorithms designed for regret-minimisation setting of bandits under unknown constraints [PGBJ20, MAAT20, LLSY21].

**Goal.** In addition to recommending a  $(1 - \delta)$ -correct and  $(1 - \delta)$ -feasible policy, we want to use minimum number of interactions to identify it. For this, pure exploration algorithms use a stopping rule which stops at a random stopping time  $\tau_\delta$ . Here, we aim to design an algorithm that recommends a  $(1 - \delta)$ -correct and  $(1 - \delta)$ -feasible policy while keeping  $\mathbb{E}[\tau_\delta]$  as small as possible.

**Motivation: Extension of multiple prior problems.** First, we clarify our motivation by showing how different interesting problems are special cases of our setting.

**a. Thresholding Bandits.** Thresholding bandits [AKR21a] are motivated from the safe dose finding problem, where one wants to identify the highest dose of a drug below a known safety level. This has also motivated the safe arm identification problem [WWJ21]. Our setting generalises it further to detect the optimal combination of doses of a set of drugs yielding highest efficacy while staying below the safety threshold. Formally, we want to identify  $\pi^* = \arg \max \mu^T \pi$ , such that  $\mathbf{I}\pi \leq \mathbf{I}\theta$ . Additionally, we allow different different thresholds for different drugs.

**b. Optimal policy under knapsack.** Bandits under knapsack constraints are studied both in BAI [LZYL23, TTCRJ12, LSY21] and regret minimisation literature [BKS18, AD16, ISSS22,

<sup>1</sup>We assume that simplex constraints are augmented in  $\mathbf{A}$  for ease of notation.

ADL16, SS18]. Detecting an optimal arm might have additional resource constraints than the number of required samples. This led to study of BAI with knapsacks under fixed-budget setting [LZYL23]. But as in regret-minimisation [SS18], one might finally want a policy that maximises utility and satisfies knapsack constraints. For example, we want to manage caches where the recommended memory allocation should satisfy a certain resource budget but can violate them during exploration. Formally,  $\pi_\tau^* = \arg \max_{\pi \in C_A} \hat{\mu}_{\tau\delta}^T \pi$ , where  $C_A \triangleq \{\mathbf{A}\pi_{\tau\delta} \leq c\}$ .

**c. BAI with fairness across sub-populations.** BAI with fairness constraints on sub-populations (BAICS) [WZZ23] aims to select an arm that must be fair across all the  $l$  sub-populations rather than the whole population as in standard BAI. Let us think of a problem where there are  $l$  sub-groups of patients and we have  $K$  number of drugs to administer with reward means  $\mu_k$ . Then we are looking for a combination of drugs rather than a single drug to administer as  $\pi^* = \arg \max_{\pi \in \Delta_K} \mu^T \pi$ , such that  $\mathbb{1}_{\mu_m \geq 0}^T \pi = 1, \forall m \in [l]$ . Hence, our setting solve the BAICS as a special case.

### 3 Lagrangian Relaxation of the Lower Bound and its Properties

Now, we discuss the Lagrangian relaxation of the lower bound and its properties that are necessary to ensure design of a correct and feasible pure exploration algorithm under unknown constraints. We require two structural assumptions to establish our approach.

**Assumption 1** (Structural assumptions on means, policy, and constraints). (a) *The mean vector  $\mu$  belongs to a bounded subset  $\mathcal{D}$  of  $\mathbb{R}^K$ .* (b) *There exists a unique optimal feasibly policy (Equation (1)).* (c) *For the true constraint  $\mathbf{A}$ , there exists a non-zero slack vector  $\Gamma$ , such that  $\max_{\pi \in \Delta_K} (-\mathbf{A}\pi) = \Gamma$ .*

The unique optimal and feasible policy assumption is used following [CBJD23] to ensure that solution of Equation (1) is an extreme point of the polytope  $\mathcal{F}$ . The assumption on slack is analogous of using the assumption of existence of a safe-arm [PGBJ20], or existence of Slater’s condition for the constraint optimisation problem [LLSY21]. Standing on these assumptions, we further prove that  $\pi_{\mathcal{F}}^*$  is unique, i.e  $\pi_{\mathcal{F}}^*$  is an extreme point in the polytope  $\hat{\mathcal{F}}$ .

#### 3.1 Information acquisition and estimates of constraints

The agent acquires new information at every step  $t \in \mathbb{N}$  by sampling an action  $A_t \sim \omega_t$ .  $\omega_t \in \Delta_K$  is called the allocation policy, and is used for interaction at step  $t$ . This yields a noisy reward  $R_t \in \mathbb{R}$  and cost vector  $\mathbf{A}_t \in \mathbb{R}^d$ . We remind that in MAB, due to independence of arms, we can represent the  $a$ -th arm as the  $a$ -th basis of  $\mathbb{R}^K$ . Thus, using the observations obtained till  $t$ , we estimate the mean vector as  $\hat{\mu}_t \triangleq \Sigma_t^{-1} \left( \sum_{s=1}^{t-1} R_s A_s \right)$ . Here,  $\Sigma_t \triangleq \sum_{s=1}^t A_s A_s^\top$  is the Gram matrix or the design matrix at time  $t$ , and  $M > 0$  is the regularisation parameter. Similarly, the estimate of the  $i$ -th row of the constraint matrix is  $\hat{\mathbf{A}}_t^i \triangleq \Sigma_t^{-1} \left( \sum_{s=1}^{t-1} \mathbf{A}_s^{i, A_s} A_s \right)$ . But we observe that using  $\hat{\mathbf{A}}_t$  to define the feasible policy set does not ensure that for any  $t$ , it is a superset of  $\mathcal{F}$ .

Thus, we define a confidence ellipsoid around  $\hat{\mathbf{A}}_t$  that always includes  $\mathbf{A}$  with probability  $1 - \delta$ , and further allows us to define a pessimistic estimate of  $\mathbf{A}$ . Formally, the confidence ellipsoid for each row  $i \in [d]$  of  $\hat{\mathbf{A}}^i$  is

$$\mathcal{C}_t \triangleq \{ \mathbf{A}' \in \mathbb{R}^{d \times K} \mid \| \mathbf{A}'^i - \hat{\mathbf{A}}_t^i \|_{\Sigma_t} \leq f(t, \delta) \forall i \in [d] \} \quad (2)$$

Here,  $f(\delta, t) \triangleq 1 + \sqrt{\frac{1}{2} \log \frac{K}{\delta} + \frac{1}{4} \log \det \Sigma_t}$  is a monotonically non-decreasing function of  $t$ .

**Lemma 1** (Pessimistic estimate of constraints). *If we use the pessimistic estimate of  $\mathbf{A}$  from the confidence ellipsoid  $\mathcal{C}_t$  to define the feasible policy set at time  $t$  as*

$$\hat{\mathcal{F}}_t \triangleq \{ \pi \in \Delta_K : \min_{\mathbf{A}' \in \mathcal{C}_t} \mathbf{A}' \pi \leq \mathbf{0} \}, \quad (3)$$

*we observe that  $\mathcal{F} \subseteq \hat{\mathcal{F}}_t$  for all  $t \in \mathbb{N}$ .*

In Figure 1, we visualise this result using the numerical values obtained from our algorithms. We observe that as we acquire more samples, our estimated feasible policy set  $\hat{\mathcal{F}}_t \rightarrow \mathcal{F}$ .

**Remark.** Our design principle of the estimators echos the spirit of the optimistic-pessimistic algorithms from the regret-minimisation under constraints literature [PGBJ20, LLSY21, CGS22b, PGB24], which has led to successful and efficient algorithms. The idea there is also to use a

187 pessimistic constraint so that we can always find the ‘true’ optimal policy in it. Similarly, here having  
 188 a pessimistic choice of  $\mathbf{A}$  results in a bigger alternative set for  $\boldsymbol{\mu}$  than the alternative set we get using  
 189  $\mathbf{A}$ , and a bigger feasible policy set that always includes the true optimal policy.

### 190 3.2 Lagrangian relaxation of the lower bound with estimated constraints

191 When we have a bandit instance with mean vector  $\boldsymbol{\mu}$  and a constraint matrix  $\mathbf{A}$ , we try to find the most  
 192 confusing instance of  $\boldsymbol{\mu}$ , so that we can minimise the KL-divergence between these two instances  
 193 to make sure we have gathered enough statistical evidence to rule out all the confusing instances.  
 194 Extending the BAI lower bound of [GK16], [CJBD24] prove that if  $\mathbf{A}$  is known, expected stopping  
 195 time of any  $(1 - \delta)$ -correct and always-feasible algorithm satisfies

$$\mathbb{E}[\tau_\delta] \geq T_{\mathcal{F}}(\boldsymbol{\mu}) \ln \frac{1}{2.4\delta}. \quad (4)$$

196 Here, the reciprocal of the characteristic time is defined by an optimisation problem over the set of  
 197 alternative instances  $\Lambda_{\mathcal{F}}(\boldsymbol{\mu}) \triangleq \{\boldsymbol{\lambda} \in \mathcal{D} : \max_{\boldsymbol{\pi} \in \mathcal{F}} \boldsymbol{\lambda}^T \boldsymbol{\pi} > \boldsymbol{\lambda}^T \boldsymbol{\pi}_{\mathcal{F}}^*\}$ :

$$T_{\mathcal{F}}^{-1}(\boldsymbol{\mu}) \triangleq \sup_{\boldsymbol{\omega} \in \mathcal{F}} \inf_{\boldsymbol{\lambda} \in \Lambda_{\mathcal{F}}(\boldsymbol{\mu})} \sum_{a=1}^K \omega_a d(\boldsymbol{\mu}_a, \boldsymbol{\lambda}_a) \triangleq \boldsymbol{\omega}^\top d(\boldsymbol{\mu}, \boldsymbol{\lambda}). \quad (5)$$

198  $\Lambda_{\mathcal{F}}(\boldsymbol{\mu})$ , aka the Alt-set, is the set of all bandit instances that have mean vectors in a bounded subset  
 199  $\mathcal{D} \in \mathbb{R}^K$  but a different optimal policy than that of  $\boldsymbol{\mu} \in \mathcal{D}$ .

200 Now, given an estimate  $\hat{\mathcal{F}}_t$  of the feasible policies at any  $t > 0$  (Eq. (3)), the corresponding Alt-set is

$$\Lambda_{\hat{\mathcal{F}}_t}(\boldsymbol{\mu}) \triangleq \{\boldsymbol{\lambda} \in \mathcal{D} : \max_{\boldsymbol{\pi} \in \hat{\mathcal{F}}_t} \boldsymbol{\lambda}^T \boldsymbol{\pi} > \boldsymbol{\lambda}^T \boldsymbol{\pi}_{\hat{\mathcal{F}}_t}^*\} \quad (6)$$

201 Since the estimated feasible policy set  $\hat{\mathcal{F}}_t$  is superset of the original feasible policy set  $\mathcal{F}$ , we observe  
 202 that  $\Lambda_{\mathcal{F}}(\boldsymbol{\mu}) \subseteq \Lambda_{\hat{\mathcal{F}}_t}(\boldsymbol{\mu})$ . Enabled by these definitions, we are ready define to Lagrangian relaxation of  
 203 the lower bound. Specifically, we observe that

$$\begin{aligned} T_{\mathcal{F}}^{-1}(\boldsymbol{\mu}) &\triangleq \sup_{\boldsymbol{\omega} \in \mathcal{F}} \inf_{\boldsymbol{\lambda} \in \Lambda_{\mathcal{F}}(\boldsymbol{\mu})} \boldsymbol{\omega}^\top d(\boldsymbol{\mu}, \boldsymbol{\lambda}) \leq \inf_{\boldsymbol{l} \in \mathbb{R}_+^d} \sup_{\boldsymbol{\omega} \in \mathcal{F}} \inf_{\boldsymbol{\lambda} \in \Lambda_{\mathcal{F}}(\boldsymbol{\mu})} \boldsymbol{\omega}^\top d(\boldsymbol{\mu}, \boldsymbol{\lambda}) + \boldsymbol{l}^T \boldsymbol{\Gamma} \\ &\leq \inf_{\boldsymbol{l} \in \mathbb{R}_+^d} \min_{\mathbf{A}' \in \mathcal{C}_t} \sup_{\boldsymbol{\omega} \in \hat{\mathcal{F}}_t} \inf_{\boldsymbol{\lambda} \in \Lambda_{\hat{\mathcal{F}}_t}(\boldsymbol{\mu})} \boldsymbol{\omega}^\top d(\boldsymbol{\mu}, \boldsymbol{\lambda}) - \boldsymbol{l}^T \mathbf{A}' \boldsymbol{\omega} \end{aligned} \quad (7)$$

204 Equation (7) defines the Lagrangian relaxation of the characteristic time under unknown constraints,  
 205 i.e. denoted by  $T_{\hat{\mathcal{F}}_t}^{-1}(\boldsymbol{\mu})$ . For non-negative Lagrange multipliers  $\boldsymbol{l} \in \mathbb{R}_+^d$ , the first inequality is true due  
 206 to the existence of a slack for the true constraints  $\mathbf{A}$ . The second inequality is due to the pessimistic  
 207 choice of the estimated constraint. Equation (7) shows that the Lagrangian relaxation of the  $T_{\hat{\mathcal{F}}_t}(\boldsymbol{\mu})$   
 208 always serves as a lower bound of the characteristic time  $T_{\mathcal{F}}^{-1}(\boldsymbol{\mu})$  for known constraints [CJBD23],  
 209 and thus, leads to a valid lower bound to optimise for the expected stopping time  $\mathbb{E}[\tau_\delta]$ . For brevity,  
 210 we denote  $\mathbf{A}'$  achieving the minimum as  $\tilde{\mathbf{A}}$ , and omit  $t$  from  $\hat{\mathcal{F}}_t$  if it is true for any  $t \in \mathbb{N}$ .

211 The Lagrangian relaxation leads to a natural question:

212 *Does the dual of the optimization problem for  $T_{\hat{\mathcal{F}}_t}^{-1}(\boldsymbol{\mu})$  yield the same solution as the primal one?*

213 **Theorem 1** (Strong Duality and Range of Lagrange Multipliers). *For a bounded sequence of  $\{l_t\}_{t \in \mathbb{N}}$ ,  
 214 strong-duality holds for the optimisation problem stated in Equation (7), i.e.*

$$\inf_{\boldsymbol{l} \in \mathbb{R}_+^d} \min_{\mathbf{A}' \in \mathcal{C}} \sup_{\boldsymbol{\omega} \in \hat{\mathcal{F}}} \inf_{\boldsymbol{\lambda} \in \Lambda_{\hat{\mathcal{F}}}(\boldsymbol{\mu})} \boldsymbol{\omega}^\top d(\boldsymbol{\mu}, \boldsymbol{\lambda}) - \boldsymbol{l}^T \tilde{\mathbf{A}} \boldsymbol{\omega} = \sup_{\boldsymbol{\omega} \in \hat{\mathcal{F}}} \min_{\boldsymbol{l} \in \mathcal{L}} \inf_{\boldsymbol{\lambda} \in \Lambda_{\mathcal{F}}(\boldsymbol{\mu})} \boldsymbol{\omega}^\top d(\boldsymbol{\mu}, \boldsymbol{\lambda}) - \boldsymbol{l}^T \tilde{\mathbf{A}} \boldsymbol{\omega}. \quad (8)$$

215 Here,  $\mathcal{L} \triangleq \{\boldsymbol{l} \in \mathbb{R}^d \mid 0 \leq \|\boldsymbol{l}\|_1 \leq \frac{1}{\gamma} T_{\hat{\mathcal{F}}}^{-1}(\hat{\boldsymbol{\mu}})\}$ , where  $\gamma \triangleq \min_{i \in [1, d]} \{-\tilde{\mathbf{A}}^i \boldsymbol{\omega}^*\}$ , i.e. the minimum  
 216 slack for pessimistic constraints w.r.t. the optimal allocation.

217 Hereafter, we use the RHS of Eq. (10) as  $T_{\hat{\mathcal{F}}}^{-1}(\boldsymbol{\mu})$ . Theorem 1 provides us a hypercube to search the  
 218 minimising Lagrangian multipliers, and thus, turns into a linear programming problem.

219 *Connections with Lagrangian-based methods in bandits.* In regret minimisation literature Lagrangian-  
 220 based optimistic-pessimistic methods [TPRL20, SSF23] are usually used to not only devise a no-regret

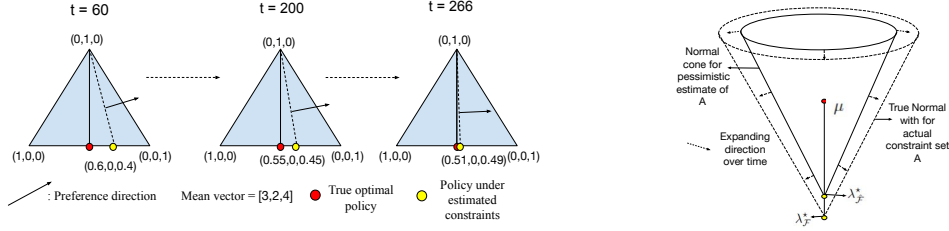


Figure 1: Convergence of the feasible set and optimal policy. Figure 2: Evolution of the normal cone.

221 learner but to control and get a sub-linear constraint violation guarantees [LLSY21, BCC24]. Our  
 222 proposed algorithm LAGEX is a prime example where the "self-boundedness" of the dual variables  
 223 results in better constraint violation guarantees (Figure 7 and 6). It would be interesting to see how  
 224 LAGEX performs in regret minimisation setting.

225 **Inner optimisation problem.** Now, we peel the layers of the optimisation problem in Eq. (10). For  
 226 known constraints, [CBJD23] has leveraged results from convex analysis [BV04] to show that the  
 227 most confusing instance for  $\mu$  lie in the boundary of the normal cone  $\Lambda_{\mathcal{F}}(\mu)^C$  (solid cone in Figure 2)  
 228 spanned by the active constraints  $\mathbf{A}_{\pi_{\mathcal{F}}^*}$  for  $\pi_{\mathcal{F}}^*$ .  $\mathbf{A}_{\pi_{\mathcal{F}}^*}$  is a sub-matrix of  $\mathbf{A}$  consisting at least  $K$   
 229 linearly independent rows. Specifically, they show that  $\mathcal{D}(\omega, \mu, \mathcal{F}) \triangleq \inf_{\lambda \in \Lambda_{\mathcal{F}}(\mu)} \omega^\top d(\mu, \lambda) =$   
 230  $\min_{\pi' \in \nu_{\mathcal{F}}(\pi_{\mathcal{F}}^*)} \min_{\lambda: \lambda^\top (\pi_{\mathcal{F}}^* - \pi') = 0} \omega^\top d(\mu, \lambda).$

231 In our setting, we are estimating both the mean vectors and the constraints. Hence, the normal cone  
 232 is also estimated at every step. Let  $\tilde{\mathbf{A}}_{\pi_{\mathcal{F}}^*}$  be the sub-matrix spanned by the active constraints for  
 233  $\pi_{\mathcal{F}}^*$ . Since we are working with the pessimistic estimate of  $\mathbf{A}$ , the vector space spanned by the  
 234 linearly independent rows of  $\mathbf{A}_{\pi_{\mathcal{F}}^*}$  is a subset of the vector space spanned by those of  $\tilde{\mathbf{A}}_{\pi_{\mathcal{F}}^*}$ . Since  
 235 the estimated Alt-set  $\Lambda_{\hat{\mathcal{F}}}(\mu)$  is always a superset of the true Alt-set  $\Lambda_{\mathcal{F}}(\mu)$ , the normal cone around  
 236  $\pi_{\mathcal{F}}^*$  is always a subset of the true normal cone  $\pi_{\mathcal{F}}^*$ . We illustrate them in Figure 2. We now extend  
 237 the projection lemma for known constraints to the Lagrangian formulation with unknown constraints.

238 **Proposition 1** (Projection Lemma for Unknown Constraints). *For any  $\omega \in \hat{\mathcal{F}}$  and  $\mu \in \mathcal{D}$ , the*  
 239 *following projection lemma holds for the Lagrangian relaxation,*

$$\mathcal{D}(\omega, \mu, \hat{\mathcal{F}}) = \min_{l \in \mathcal{L}} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \min_{\lambda: \lambda^\top (\pi_{\hat{\mathcal{F}}}^* - \pi') = 0} \omega^\top d(\mu, \lambda) - l^\top \tilde{\mathbf{A}} \omega. \quad (9)$$

240 This reduces the inner minimisation problem to a less intensive discrete optimisation, where we only  
 241 have to search over the neighbouring vertices of the optimal policy in  $\hat{\mathcal{F}}$  for a solution. Now, a natural  
 242 question arises around this formulation:

243 *Can we continuously track the lower bound defined by projection lemma for unknown constraints?*

244 **Theorem 2.** *For a sequence  $\{\hat{\mathcal{F}}_t\}_{t \in \mathbb{N}}$  and  $\{\hat{\lambda}_t\}_{t \in \mathbb{N}}$ , we show that (a)  $\lim_{t \rightarrow \infty} \hat{\mathcal{F}}_t \rightarrow \mathcal{F}$ , (b)  $\lambda^*$  is*  
 245 *unique, and (c)  $\lim_{t \rightarrow \infty} \hat{\lambda}_t \rightarrow \lambda^*$ . Thus, for any  $\omega \in \mathcal{F}$  and  $\mu$ ,  $\lim_{t \rightarrow \infty} \mathcal{D}(\omega, \mu, \hat{\mathcal{F}}) \rightarrow \mathcal{D}(\omega, \mu, \mathcal{F})$*   
 246 *where  $\lambda^*$  is such that for any  $\lambda^* \in \arg \min_{\lambda \in \Lambda_{\mathcal{F}}(\mu)} \omega^\top d(\mu, \lambda)$ .*

247 **Outer optimisation problem.** As the convergence of  $\hat{\mathcal{F}}_t$ , we are left with the outer optimisation  
 248 problem in Equation (10). Since it is a linear problem in  $\omega$ , we can use a linear programming method  
 249 which would lead to one of the vertices of  $\hat{\mathcal{F}}$ . But to be sure of an existence of a solution at each  
 250  $t \in \mathbb{N}$ , we try and derive some well-behavedness properties of the optimal allocation  $\omega^*(\mu)$ . First,  
 251 we observe that our estimates of the mean vector converge to  $\mu$  as  $t \rightarrow \infty$ . Hence, we also get  
 252  $\lim_{t \rightarrow \infty} \mathcal{D}(\omega, \hat{\mu}_t, \hat{\mathcal{F}}) \rightarrow \mathcal{D}(\omega, \mu, \mathcal{F})$ . Now, we ensure well-behavedness and existence of an optimal  
 253 allocation for all  $t > 0$ .

254 **Theorem 3.** (Existence of unique optimal allocation) *For any  $\mu \in \mathcal{D}$ , the optimization problem*  
 255  *$\max_{\pi \in \hat{\mathcal{F}}} \mu^\top \pi$  has a unique solution if  $\omega^*(\mu)$  satisfies the conditions: 1. Both the sets  $\hat{\mathcal{F}}$  and*  
 256  *$\omega^*(\mu)$  are closed and convex. 2.  $\forall \mu \in \mathcal{D}$  and  $\omega \in \hat{\mathcal{F}}$ , the function  $\lim_{t \rightarrow \infty} \mathcal{D}(\omega, \hat{\mu}_t, \hat{\mathcal{F}})$  is*  
 257 *continuous. 3. Reciprocal of the characteristic time function  $\lim_{t \rightarrow \infty} T_{\hat{\mathcal{F}}}^{-1}(\mu)$  is continuous  $\forall \mu \in \mathcal{D}$ .*  
 258 *4.  $\mu \in \mathcal{D} : \mu \rightarrow \omega^*(\mu)$  is upper hemi-continuous.*

259 **Characterising the lower bound for Gaussians.** Since we can derive explicit form of the optimisa-  
 260 tion problem for Gaussian reward distributions, we characterise it further to relate our lower bound  
 261 with the lower bound for known constraints.

**Theorem 4.** Let  $\{P_a\}_{a \in [K]}$  be Gaussian distributions with equal variance  $\sigma^2 > 0$ , and  $\text{Diag}(1/\omega_a)$  be a  $K$ -dimensional diagonal matrix with  $a$ -th diagonal entry  $1/\omega_a$ . Then, we get

$$T_{\hat{\mathcal{F}}}^{-1}(\boldsymbol{\mu}) = \max_{\boldsymbol{\omega} \in \hat{\mathcal{F}}} \min_{l \in \mathcal{L}} \min_{\boldsymbol{\pi}' \in \nu_{\hat{\mathcal{F}}}(\boldsymbol{\pi}_{\hat{\mathcal{F}}}^*)} \left\{ \frac{1}{2\sigma^2} \frac{\|\boldsymbol{\pi}_{\hat{\mathcal{F}}}^* - \boldsymbol{\pi}'\|_{\boldsymbol{\mu}\boldsymbol{\mu}^T}^2}{\|\boldsymbol{\pi}_{\hat{\mathcal{F}}}^* - \boldsymbol{\pi}'\|_{\text{Diag}(1/\omega_a)}^2} - \boldsymbol{l}^T \tilde{\mathbf{A}} \boldsymbol{\omega} \right\}.$$

Corollary 1 gives explicit upper and lower bound on characteristic bound for Gaussian rewards. It also shows explicit connection to the lower bounds under known constraints. Finally, it tells that the hardness of the pure exploration under constraints depends inversely on the condition number of the estimated and true constraints, which quantify their invertibility.

**Corollary 1.** Let  $d_{\boldsymbol{\pi}}^2 \triangleq \frac{\|\boldsymbol{\pi}_{\hat{\mathcal{F}}}^* - \boldsymbol{\pi}\|_{\boldsymbol{\mu}\boldsymbol{\mu}^T}^2}{\|\boldsymbol{\pi}_{\hat{\mathcal{F}}}^* - \boldsymbol{\pi}\|_2^2}$  be the norm of the projection of  $\boldsymbol{\mu}$  on the policy gap  $(\boldsymbol{\pi}_{\hat{\mathcal{F}}}^* - \boldsymbol{\pi})$ . The characteristic time  $T_{\hat{\mathcal{F}}}(\boldsymbol{\mu})$  satisfies two bounds. (a)  $\frac{2\sigma^2}{C_{\text{known}} + 2C_{\text{unknown}}} \leq T_{\hat{\mathcal{F}}}(\boldsymbol{\mu}) \leq \frac{2\sigma^2 K}{C_{\text{known}}}$ , where  $C_{\text{unknown}} \triangleq \min_{\boldsymbol{\pi}' \in \nu_{\hat{\mathcal{F}}}(\boldsymbol{\pi}_{\hat{\mathcal{F}}}^*)} d_{\boldsymbol{\pi}'}^2$  and  $C_{\text{known}} = \min_{\boldsymbol{\pi}'' \in \nu_{\mathcal{F}}(\boldsymbol{\pi}_{\mathcal{F}}^*)} d_{\boldsymbol{\pi}''}^2$ . (b)  $T_{\hat{\mathcal{F}}}(\boldsymbol{\mu}) \geq \frac{H}{\kappa_{\text{known}}^2 + 2\kappa_{\text{unknown}}^2}$ .  $H$  is the sum of squares of gaps.  $\kappa_{\text{known}}$  and  $\kappa_{\text{unknown}}$  are condition numbers of  $\mathbf{A}$  and  $\tilde{\mathbf{A}}$ .

**Connection to existing lower bounds.** (a) *Pure exploration under known constraints.* The upper and lower bounds on characteristic time coincides with the existing lower bound under known constraints, i.e. when  $C_{\text{unknown}} = 0$ . (b) *BAI without constraints.* In BAI, we consider only deterministic policies playing an arm. Thus,  $d_{\boldsymbol{\pi}_a} = \frac{\boldsymbol{\mu}^T(\boldsymbol{\pi}_{\hat{\mathcal{F}}}^* - \boldsymbol{\pi}_a)}{\|\boldsymbol{\pi}_{\hat{\mathcal{F}}}^* - \boldsymbol{\pi}_a\|_2} = \boldsymbol{\mu}^* - \boldsymbol{\mu}_a$ , i.e. the sub-optimality gap for arm  $a$ . Here,  $\boldsymbol{\mu}^*$  is the mean of the best arm. Then, the ‘known’ term in the denominator is the minimum squared sub-optimality gap and matches the BAI bound [CBJD24]. But the ‘unknown’ term in the lower bound is due to the unknown constraints. The term is the squared sub-optimality gap under the estimated feasible set  $\hat{\mathcal{F}}$ . Thus, our upper bound on characteristic time matches that of BAI [KCG16], while the lower bound has an added term due to the extra cost of exploring under unknown constraints. (c) *Safe linear BAI.* [WWJ21] stated a lower bound in for BAI under safety constraints but when the reward generation has a linear structure, and while the rewards and constraints are jointly generated for an arm. This is a bit different setting. Similar to our lower bound, they also show an added cost for constraints which is not directly comparable.

## 4 LATS and LAGEX: Algorithm design and analysis

Now, we propose two algorithms to conduct pure exploration with the Lagrangian relaxation of the lower bound, and derive upper bounds on their sample complexities.

**Assumption 2** (Distributional assumptions on rewards and constraints). *We require two distributional assumptions on rewards and constraints. (i) Reward distributions  $\{P_a\}_{a=1}^K$  are sub-Gaussian one parameter exponential family with mean vector  $\boldsymbol{\mu} \in \mathcal{D}$ . (ii) Each constraint follows a sub-Gaussian  $K$ -parameter exponential family parameterised by  $\mathbf{A}^i$  for  $i \in [d]$ .*

These assumptions are standard in bandits under constraints [CBJD23, DK19, PGBJ20, PGB24].

**Algorithm design.** Any algorithm in pure exploration setting comprises of three main components.

**Component 1 : Stopping rule.** Stopping rule of an exploring algorithm decides when to stop sampling. While exploring, once we gather enough statistical information about the parameters in the system, the test statistic crosses the stopping threshold with the chosen confidence level  $\delta$  and exploration stops to recommend the best policy.

**Theorem 5.** *The Chernoff stopping rule to ensure  $(1 - \delta)$ -correctness and  $(1 - \delta)$ -feasibility is*

$$\inf_{\boldsymbol{\lambda} \in \Lambda_{\hat{\mathcal{F}}_t}(\hat{\boldsymbol{\mu}}_t)} \sum_{a=1}^K N_{a,t} d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\lambda}_a) > \beta(t, \delta) \wedge \|(\tilde{\mathbf{A}}_t - \mathbf{A})\boldsymbol{\omega}_t\|_{\infty} > \rho(t, \delta),$$

where  $\beta(t, \delta) \triangleq 3S_0 \log(1 + \log N_{a,t}) + S_0 \mathcal{T} \left( \frac{(K \wedge d) + \log \frac{1}{\delta}}{S_0} \right)$ , and  $\rho(t, \delta)$  is in Lemma 3.

First, we need to check the Chernoff condition under the estimated alternative set as we do not know the true one. Secondly, we also need to ensure that the constraint matrix has also concentrated around the true matrix. Thus, the exploration stops when both the events occur together ensuring concentration to both the correct and feasible policy set.

---

**Algorithm 1** LATS - LAgrangian Track and Stop

---

```

1: Input : Time Horizon  $T$ , Confidence level  $\delta > 0$ 
2: Initialization :  $\hat{\mathbf{A}}_0 = \mathbf{0}_{d \times K}$ ,  $\hat{\boldsymbol{\mu}}_0 = \mathbf{0}_K$ ,  $\Sigma_0 = \lambda \mathbf{1}_K$ ,  $l_0$ 
3: Play each arm once to set  $\mu_1$  and  $\hat{\mathbf{A}}_1$ .
4: while  $\beta(t, \delta) > \mathcal{D}(\omega_t, \hat{\boldsymbol{\mu}}_t, \hat{\mathcal{F}}_t, l_t^*) \wedge \|(\tilde{\mathbf{A}}_t - \mathbf{A})\omega_t\|_\infty > \rho(t, \delta)$  do  $\rightsquigarrow$  Proposition 1)
5:   Optimal Policy:  $\omega_t^* \in \arg \max_{\omega \in \hat{\mathcal{F}}_t} \mathcal{D}(\omega, \hat{\boldsymbol{\mu}}_t, \hat{\mathcal{F}}_t)$ 
6:   Optimize Lagrange Multiplier :  $l_t^* \in \arg \min_{l \in \mathcal{L}_t} \mathcal{D}(\omega_t^*, \hat{\boldsymbol{\mu}}_t, \hat{\mathcal{F}}_t)$ 
7:   C-Tracking: Play  $a_t \in \arg \min_{a \in [1, K]} N_{a,t} - \sum_{s=1}^t \omega_{a,s}^*$ 
8:   Feedback : Observe reward  $r_t$  and cost  $\mathbf{A}_{a_t}$ , and update  $\Sigma_{t+1}$ ,  $\hat{\boldsymbol{\mu}}_{t+1}$  and  $\hat{\mathbf{A}}_{t+1}$ 
9: end while
10: Recommended policy:  $\pi_{\hat{\mathcal{F}}_t}^* = \arg \max_{\pi \in \hat{\mathcal{F}}_t} \hat{\boldsymbol{\mu}}_t^T \pi$ 

```

---



---

**Algorithm 2** LAGEX - LAgrangian Game EXplorer

---

```

1: Input and Initialisation as same as LATS
2: Play each arm once to set  $\mu_1$  and  $\hat{\mathbf{A}}_1$ .
3: while  $g(t, \delta) > \mathcal{D}(\omega_t, \hat{\boldsymbol{\mu}}_t, \hat{\mathcal{F}}_t, l_t^*) \wedge \|(\tilde{\mathbf{A}}_t - \mathbf{A})\omega_t\|_\infty > \rho(t, \delta)$  do  $\rightsquigarrow$  Proposition 1)
4:   Optimal allocation  $\omega_t \rightsquigarrow$  Using AdaGrad via Theorem 4
5:   Optimize Lagrange Multiplier:  $l_t^* \in \arg \min_{l \in \mathcal{L}_t} \mathcal{D}(\omega_t, \hat{\boldsymbol{\mu}}_t, \hat{\mathcal{F}}_t)$ 
6:   Compute confusing instance :  $\lambda_t \rightsquigarrow$  Via Proposition 1 plugging in  $\omega_t, l_t^*$ 
7:   Confidence intervals: for all  $a \in [K]$   $[\alpha_{t,a}, \beta_{t,a}] : \{\zeta : N_{a,t} d(\mu_{t,a}, \zeta) \leq g(t)\}$ 
    $U_t^a = \max \left\{ \frac{g(t)}{N_{a,t}}, d(\alpha_{t,a}, \lambda_{t,a}), d(\beta_{t,a}, \lambda_{t,a}) \right\}$ 
8:   Update loss for regret minimizer:  $l(\omega_t) = \langle \omega_t, U_t \rangle - l_t^{*T} \tilde{\mathbf{A}}_t \omega_t$ 
9:   C-Tracking : Play  $a_t \in \arg \min_{a \in [1, K]} N_{a,t} - \sum_{s=1}^t \omega_{a,s}^*$ 
10:  Feedback : Observe reward  $r_t$  and cost  $\mathbf{A}_{a_t}$ , and update  $\Sigma_{t+1}$ ,  $\hat{\boldsymbol{\mu}}_{t+1}$  and  $\hat{\mathbf{A}}_{t+1}$ 
11: end while
12: Recommended:  $\pi_{\hat{\mathcal{F}}_t}^* = \arg \max_{\pi \in \hat{\mathcal{F}}_t} \hat{\boldsymbol{\mu}}_t^T \pi$ 

```

---

**Component 2: Recommendation rule.** Once the stopping rule is fired, the agent recommends a policy based on the current estimate of  $\hat{\boldsymbol{\mu}}_t$  according the rule  $\pi_{\hat{\mathcal{F}}_{\tau_\delta}}^* = \arg \max_{\pi \in \hat{\mathcal{F}}_{\tau_\delta}} \hat{\boldsymbol{\mu}}_t^T \pi$ .

**Component 3: Sampling strategy.** We present two novel sampling algorithms: LATS and LAGEX.

**a. LATS.** The algorithm LATS (Algorithm 1) uses a Track and Stop strategy adapted to the unknown constraint setting. We use markers in the algorithm which are novel approaches used to handle the challenge of estimating the feasible space per step. The algorithm first warms up the parameter estimates by playing each arm once. Then until the test statistic jumps the threshold, first in line 5, it calculates the optimal policy under the estimated feasible space at the current step solving the Lagrangian relaxed optimization problem in Proposition 1 by plugging in the best choice of Lagrangian multiplier optimized in previous step. Using this current optimal policy we optimize the Lagrangian multiplier maintaining the bounds of it's 1-norm stated in Theorem 1. It uses C-tracking [GK16] to track the action taken per step. At line 9 and 10, we basically observe the instantaneous reward and cost feedback and update the parameter estimates based on them.

**Theorem 6.** Let  $\mathfrak{s}$  be the shadow price  $\mathfrak{s} \triangleq \frac{\Gamma_{\max}}{\Gamma_{\min}}$  of the slack  $\Gamma$ . For any  $\alpha > 1$ , the expected stopping time of LATS satisfies  $\lim_{\delta \rightarrow 0} \frac{\mathbb{E}[\tau_\delta]}{\log(1/\delta)} \leq \alpha T_{\mathcal{F}}(\boldsymbol{\mu})(1 + \mathfrak{s})$ .

**b. LAGEX.** Historically, algorithms base on track and stop mechanism seems to fail in case of larger problems where efficient optimization becomes a challenge due to the use of a max-min oracle per step (Line 5, Algorithm 2). To improve on this we land on the two-player zero sum game approach introduced in [DKM19b]. The second algorithm 2 we introduce in this work also starts by playing each arm once to initially start the estimation of the parameter in the system. Then it uses a **allocation player** (We have used AdaGrad) to optimize the allocation  $\omega_t$  in line 5 against the most confusing instance w.r.t current estimate of  $\hat{\boldsymbol{\mu}}_t$  optimized by a **instance player** which minimizes  $\sum_{a=1}^K \omega_{a,t} d(\hat{\boldsymbol{\mu}}_{a,t}, \lambda_a) - l_t^{*T} \tilde{\mathbf{A}}_t \omega_t$  w.r.t  $\lambda \in \lambda_{\hat{\mathcal{F}}_t}(\boldsymbol{\mu})$  by plugging in the chosen allocation, estimated feasible set  $\tilde{\mathbf{A}}_t$  and the optimized Lagrangian multiplier  $l_t^*$ . Since our search space in closed and convex, the allocation player enjoys sub-linear regret of order  $\mathcal{O}(\sqrt{t \log t})$ , whereas the instance player computes the best confusing instance using the Lagrangian formulation of the weighted projection lemma stated in Proposition 1. Then, in line 11, Adagrad loss function is updated with a loss by introducing optimism as  $U_t$  defined in line 10 with an extra term that is novel in the literature to track the unknown constraint set. LAGEX also uses C-tracking similar to LATS to track the actions taken per step. Then it goes on to observe the instantaneous reward and cost to update the estimates for the next optimization step.



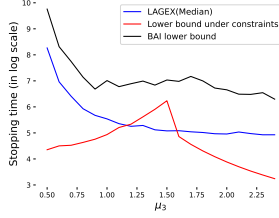


Figure 3: Lower bounds with and without constraints, and 500 runs of LAGEX for  $\mu_3 \in [0.5, 2.5]$ .

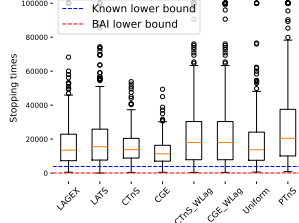


Figure 4: Sample complexity (median $\pm$ std.) of algorithms over 500 runs for **hard env.**

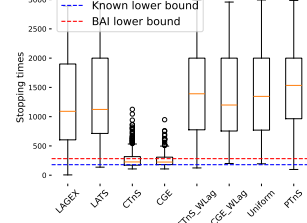


Figure 5: Sample complexity (median $\pm$ std.) of algorithms over 500 runs for **easy env.**

**Theorem 7.** *The expected sample complexity of LAGEX satisfies  $\lim_{\delta \rightarrow 0} \frac{\mathbb{E}(\tau_\delta)}{\ln(1/\delta)} \leq T_{\mathcal{F}}(\mu) + 2\epsilon$ .*

*Implications.* Theorem 6 and 7 show both LATS and LAGEX are asymptotically stable. Also, for LATS, the effect of the unknown constraints in the form of **shadow price** arises in a multiplicative way, whereas for LAGEX, it is additive. Thus, LAGEX should show a lower sample complexity than LATS, which is later validated in the experiments.

## 5 Experimental analysis

Now, we experimentally evaluate performances of LAGEX and LAGTS across environments and against baseline algorithms. Code available in this Link.

**Setup.** We run the algorithms on a 64-bit 13th Gen Intel® Core™ i7-1370P  $\times$  20 processor machine with 32GB ram. We evaluate using a set of environments with mean vectors  $[1.5, 1.0, \mu_3, 0.4, 0.3, 0.2, 0.1]$ . We impose two linear constraints  $\pi_1 + \pi_2 + \pi_3 \leq 0.5$  and  $\pi_4 + \pi_5 \leq 0.5$ . We conduct two experiments to validate universality and efficiency of LAGEX and LATS. We set  $\delta = 0.01$  for all the experiments.

**Experiment 1: Universality.** We vary  $\mu_3$  from 0.5 to 2.5. For each environment, we plot the corresponding BAI lower bounds (in red) and lower bounds under constraints (in blue) in Figure 3. We observe that the constraint problem gets easier with increasing  $\mu_3$ . In contrast, the BAI problem changes non-monotonically. BAI problem gets harder when  $\mu_3$  is around 1.5 as the suboptimality gap gets very small. But the constraint problem stays easier than BAI. In Figure 3, we also plot the median sample complexity of LAGEX across these environments over 500 runs. We observe that LAGEX grows parallel to the lower bound under constraints and can track it across environments.

**Experiment 2: Efficiency against existing algorithms.** We compare LAGEX and LATS with the two algorithms under known constraints, i.e. CTnS and CGE [CBJD23]. Also, to understand the utility of Lagrangian relaxation, we implement versions of CTnS and CGE with estimated constraints. In these variants, we solve the constrained optimisation problems without Lagrangian relaxation but with estimated constraints. We also compare with PTnS (Projected Track and Stop), a variant of TnS, where the algorithm solves a standard BAI problem and projects the allocation to the estimated feasible space. We run all these algorithms in two environment: (i) **hard** with  $\mu_3 = 0.5$  and (ii) **easy** with  $\mu_3 = 1.3$ . We call the first environment hard as it is harder than BAI and similarly, the second environment easy. We plot three quantiles of the sample complexities of these algorithms over 500 runs in Figure 4 and 5. We observe that (i) among the algorithms with unknown constraints incur the least sample complexity, and (ii) we pay a minimal cost than the known constraint Lagrangian algorithms in **hard env** whereas the price of estimating constraints is prominent in **easy env**.

## 6 Conclusion and Future Works

We study the problem of pure exploration under unknown linear constraints. This problem requires tracking both mean vectors and constraints to recommend a correct and feasible policy. We encompass this effect with a Lagrangian relaxation of the lower bound for known constraints. We further design an pessimistic estimate of constraints to ensure identification of the optimal feasible policy. These tools allows us to propose two algorithms LATS and LAGEX. We prove their sample complexity upper bounds, and conduct numerical experiments to find that LAGEX is the most efficient among baselines. In reality, constraints can be non-linear. Our concentration bounds are tailored for linearity. It would be interesting and challenging to extend our Lagrangian-based technique to nonlinear constraints.

## References

- [AAT19] Sanae Amani, Mahnoosh Alizadeh, and Christos Thrampoulidis. Linear stochastic bandits under safety constraints, 2019.
- [ACBF02] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multi-armed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- [AD14] Shipra Agrawal and Nikhil R. Devanur. Bandits with concave rewards and convex knapsacks. EC '14, page 989–1006, New York, NY, USA, 2014. Association for Computing Machinery.
- [AD16] Shipra Agrawal and Nikhil Devanur. Linear contextual bandits with knapsacks. 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.
- [ADL16] Shipra Agrawal, Nikhil R. Devanur, and Lihong Li. An efficient algorithm for contextual bandits with knapsacks, and an extension to concave objectives. In Vitaly Feldman, Alexander Rakhlin, and Ohad Shamir, editors, *29th Annual Conference on Learning Theory*, volume 49 of *Proceedings of Machine Learning Research*, pages 4–18, Columbia University, New York, New York, USA, 23–26 Jun 2016. PMLR.
- [AKR21a] Maryam Aziz, Emilie Kaufmann, and Marie-Karelle Riviere. On multi-armed bandit designs for dose-finding clinical trials. *The Journal of Machine Learning Research*, 22(1):686–723, 2021.
- [AKR21b] Maryam Aziz, Emilie Kaufmann, and Marie-Karelle Riviere. On multi-armed bandit designs for dose-finding trials. *Journal of Machine Learning Research*, 22(14):1–38, 2021.
- [AyPS11] Yasin Abbasi-yadkori, Dávid Pál, and Csaba 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.
- [BB62] Robert E. Bechhofer and Saul Blumenthal. A sequential multiple-decision procedure for selecting the best one of several normal populations with a common unknown variance, ii: Monte carlo sampling results and new computing formulae. *Biometrics*, 18(1):52–67, 1962.
- [BCC24] Martino Bernasconi, Matteo Castiglioni, and Andrea Celli. No-regret is not enough! bandits with general constraints through adaptive regret minimization, 2024.
- [Ber63] C. Berge. *Topological Spaces: Including a Treatment of Multi-valued Functions, Vector Spaces and Convexity*. Macmillan, 1963.
- [BKS18] Ashwinkumar Badanidiyuru, Robert Kleinberg, and Aleksandrs Slivkins. Bandits with knapsacks. *J. ACM*, 65(3), mar 2018.
- [BMS09] Sébastien Bubeck, Rémi Munos, and Gilles Stoltz. Pure exploration in multi-armed bandits problems. In *Algorithmic Learning Theory: 20th International Conference, ALT 2009, Porto, Portugal, October 3-5, 2009. Proceedings 20*, pages 23–37. Springer, 2009.
- [BMS10] Sébastien Bubeck, Rémi Munos, and Gilles Stoltz. Pure exploration for multi-armed bandit problems, 2010.
- [BV04] Stephen Boyd and Lieven Vandenbergh. *Convex Optimization*. Cambridge University Press, 2004.
- [CBJD23] Emil Carlsson, Debabrota Basu, Fredrik D. Johansson, and Devdatt Dubhashi. Pure Exploration in Bandits with Linear Constraints. In *EWRL 2023 – European Workshop on Reinforcement Learning*, Brussels, Belgium, September 2023.

- [CJBD24] Emil Carlsson, Debabrota Basu, Fredrik Johansson, and Devdatt Dubhashi. Pure exploration in bandits with linear constraints. In *International Conference on Artificial Intelligence and Statistics*, pages 334–342. PMLR, 2024.
- [CGS22a] Tianrui Chen, Aditya Gangrade, and Venkatesh Saligrama. Doubly-optimistic play for safe linear bandits. *arXiv preprint arXiv:2209.13694*, 2022.
- [CGS22b] Tianrui Chen, Aditya Gangrade, and Venkatesh Saligrama. Strategies for safe multi-armed bandits with logarithmic regret and risk. In *International Conference on Machine Learning*, pages 3123–3148. PMLR, 2022.
- [CLK<sup>+</sup>14] Shouyuan Chen, Tian Lin, Irwin King, Michael R Lyu, and Wei Chen. Combinatorial pure exploration of multi-armed bandits. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
- [CMC21] James Cheshire, Pierre Menard, and Alexandra Carpentier. The influence of shape constraints on the thresholding bandit problem, 2021.
- [CWM<sup>+</sup>22a] Romain Camilleri, Andrew Wagenmaker, Jamie H Morgenstern, Lalit Jain, and Kevin G Jamieson. Active learning with safety constraints. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 33201–33214. Curran Associates, Inc., 2022.
- [CWM<sup>+</sup>22b] Romain Camilleri, Andrew Wagenmaker, Jamie H Morgenstern, Lalit Jain, and Kevin G Jamieson. Active learning with safety constraints. *Advances in Neural Information Processing Systems*, 35:33201–33214, 2022.
- [DK19] Rémy Degenne and Wouter M. Koolen. Pure exploration with multiple correct answers. In *Neural Information Processing Systems*, 2019.
- [DKM19a] Rémy Degenne, Wouter M Koolen, and Pierre Ménard. Non-asymptotic pure exploration by solving games. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [DKM19b] Rémy Degenne, Wouter M. Koolen, and Pierre Ménard. Non-asymptotic pure exploration by solving games, 2019.
- [EDMM02a] Eyal Even-Dar, Shie Mannor, and Y. Mansour. Pac bounds for multi-armed bandit and markov decision processes. In *Annual Conference Computational Learning Theory*, 2002.
- [EDMM02b] Eyal Even-Dar, Shie Mannor, and Yishay Mansour. Pac bounds for multi-armed bandit and markov decision processes. In *Computational Learning Theory: 15th Annual Conference on Computational Learning Theory, COLT 2002 Sydney, Australia, July 8–10, 2002 Proceedings 15*, pages 255–270. Springer, 2002.
- [FJR19] Tanner Fiez, Lalit Jain, Kevin G Jamieson, and Lillian Ratliff. Sequential experimental design for transductive linear bandits. *Advances in neural information processing systems*, 32, 2019.
- [FN22] Fathima Zarin Faizal and Jayakrishnan Nair. Constrained pure exploration multi-armed bandits with a fixed budget. *arXiv preprint arXiv:2211.14768*, 2022.
- [GK16] Aurélien Garivier and Emilie Kaufmann. Optimal best arm identification with fixed confidence. In Vitaly Feldman, Alexander Rakhlin, and Ohad Shamir, editors, *29th Annual Conference on Learning Theory*, volume 49 of *Proceedings of Machine Learning Research*, pages 998–1027, Columbia University, New York, New York, USA, 23–26 Jun 2016. PMLR.
- [GK21a] Aurélien Garivier and Tomáš Kocák. Epsilon Best Arm Identification in Spectral Bandits. In *Thirtieth International Joint Conference on Artificial Intelligence (IJCAI-21)*, pages 2636–2642, Montreal, Canada, August 2021. International Joint Conferences on Artificial Intelligence Organization.

[GK21b] Aurélien Garivier and Emilie Kaufmann. Non-asymptotic sequential tests for overlapping hypotheses and application to near optimal arm identification in bandit models, 2021.

[GMRM18] Aurélien Garivier, Pierre Ménard, Laurent Rossi, and Pierre Menard. Thresholding bandit for dose-ranging: The impact of monotonicity, 2018.

[HTA24] Spencer Hutchinson, Berkay Turan, and Mahnoosh Alizadeh. Directional optimism for safe linear bandits, 2024.

[ISSS22] Nicole Immorlica, Karthik Sankararaman, Robert Schapire, and Aleksandrs Slivkins. Adversarial bandits with knapsacks. *J. ACM*, 69(6), nov 2022.

[JMKK21] Marc Jourdan, Mojmír Mutný, Johannes Kirschner, and Andreas Krause. Efficient pure exploration for combinatorial bandits with semi-bandit feedback. In *Proceedings of the 32nd International Conference on Algorithmic Learning Theory*, volume 132, 2021.

[JN14] Kevin Jamieson and Robert Nowak. Best-arm identification algorithms for multi-armed bandits in the fixed confidence setting. In *2014 48th Annual Conference on Information Sciences and Systems (CISS)*, pages 1–6. IEEE, 2014.

[KCG16] Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best arm identification in multi-armed bandit models, 2016.

[KCJ23] Newton Mwai Kinyanjui, Emil Carlsson, and Fredrik D. Johansson. Fast treatment personalization with latent bandits in fixed-confidence pure exploration. *Trans. Mach. Learn. Res.*, 2023, 2023.

[KGAYR17] Abbas Kazerouni, Mohammad Ghavamzadeh, Yasin Abbasi-Yadkori, and Benjamin Van Roy. Conservative contextual linear bandits, 2017.

[KK21] Emilie Kaufmann and Wouter M. Koolen. Mixture martingales revisited with applications to sequential tests and confidence intervals. *Journal of Machine Learning Research*, 22(246):1–44, 2021.

[KSS18] Julian Katz-Samuels and Clay Scott. Feasible arm identification. In *International Conference on Machine Learning*, pages 2535–2543. PMLR, 2018.

[LJD<sup>+</sup>17] Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband: A novel bandit-based approach to hyperparameter optimization. *The Journal of Machine Learning Research*, 18(1):6765–6816, 2017.

[LLSY21] Xin Liu, Bin Li, Pengyi Shi, and Lei Ying. An efficient pessimistic-optimistic algorithm for stochastic linear bandits with general constraints, 2021.

[LPJ22] Simon Lindstål, Alexandre Proutiere, and Andreas Johnsson. Measurement-based admission control in sliced networks: A best arm identification approach. In *GLOBE-COM 2022-2022 IEEE Global Communications Conference*, pages 1484–1490. IEEE, 2022.

[LS20] Tor Lattimore and Csaba Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.

[LSY21] Xiaocheng Li, Chunlin Sun, and Yinyu Ye. The symmetry between arms and knapsacks: A primal-dual approach for bandits with knapsacks. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 6483–6492. PMLR, 18–24 Jul 2021.

[LTHK22] David Lindner, Sebastian Tschieschek, Katja Hofmann, and Andreas Krause. Interactively learning preference constraints in linear bandits. In *International Conference on Machine Learning*, pages 13505–13527. PMLR, 2022.

- [LVR<sup>+</sup>19] Pieter JK Libin, Timothy Verstraeten, Diederik M Roijers, Jelena Grujic, Kristof Theys, Philippe Lemey, and Ann Nowé. Bayesian best-arm identification for selecting influenza mitigation strategies. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2018, Dublin, Ireland, September 10–14, 2018, Proceedings, Part III 18*, pages 456–471. Springer, 2019.
- [LZYL23] Shaoang Li, Lan Zhang, Yingqi Yu, and Xiangyang Li. Optimal arms identification with knapsacks. In *International Conference on Machine Learning*, pages 20529–20555. PMLR, 2023.
- [Ma] Will Ma. *Improvements and Generalizations of Stochastic Knapsack and Multi-Armed Bandit Approximation Algorithms: Extended Abstract*, pages 1154–1163.
- [MAAT20] Ahmadreza Moradipari, Sanae Amani, Mahnoosh Alizadeh, and Christos Thrampoulidis. Safe linear thompson sampling with side information, 2020.
- [MCP14] Stefan Magureanu, Richard Combes, and Alexandre Proutiere. Lipschitz bandits: Regret lower bounds and optimal algorithms, 2014.
- [MDP<sup>+</sup>11] Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In Dekang Lin, Yuji Matsumoto, and Rada Mihalcea, editors, *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics.
- [MJTN20] Blake Mason, Lalit Jain, Ardhendu Tripathy, and Robert Nowak. Finding all  $\epsilon$ -good arms in stochastic bandits. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 20707–20718. Curran Associates, Inc., 2020.
- [MPK21] Arnab Maiti, Vishakha Patil, and Arindam Khan. Multi-armed bandits with bounded arm-memory: Near-optimal guarantees for best-arm identification and regret minimization. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 19553–19565. Curran Associates, Inc., 2021.
- [Pau64] Edward Paulson. A sequential procedure for selecting the population with the largest mean from  $k$  normal populations. *The Annals of Mathematical Statistics*, 35(1):174–180, 1964.
- [PGB24] Aldo Pacchiano, Mohammad Ghavamzadeh, and Peter Bartlett. Contextual bandits with stage-wise constraints. *arXiv preprint arXiv:2401.08016*, 2024.
- [PGBJ20] Aldo Pacchiano, Mohammad Ghavamzadeh, Peter Bartlett, and Heinrich Jiang. Stochastic bandits with linear constraints, 2020.
- [SCBC23] Xuedong Shang, Igor Colin, Merwan Barlier, and Hamza Cherkaoui. Price of safety in linear best arm identification, 2023.
- [SJ19] Ashudeep Singh and Thorsten Joachims. Policy learning for fairness in ranking. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [Slu25] E. Slutsky. Über stochastische Asymptoten und Grenzwerte. *Metron* 5, Nr. 3, 3-89 (1925)., 1925.
- [SS18] Karthik Abinav Sankararaman and Aleksandrs Slivkins. Combinatorial semi-bandits with knapsacks. In Amos Storkey and Fernando Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 1760–1770. PMLR, 09–11 Apr 2018.

572 [SSF23] Aleksandrs Slivkins, Karthik Abinav Sankararaman, and Dylan J Foster. Contextual  
573 bandits with packing and covering constraints: A modular lagrangian approach via  
574 regression. In Gergely Neu and Lorenzo Rosasco, editors, *Proceedings of Thirty Sixth*  
575 *Conference on Learning Theory*, volume 195 of *Proceedings of Machine Learning*  
576 *Research*, pages 4633–4656. PMLR, 12–15 Jul 2023.

577 [TPRL20] Andrea Tirinzoni, Matteo Pirotta, Marcello Restelli, and Alessandro Lazaric. An  
578 asymptotically optimal primal-dual incremental algorithm for contextual linear bandits.  
579 In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances*  
580 *in Neural Information Processing Systems*, volume 33, pages 1417–1427. Curran  
581 Associates, Inc., 2020.

582 [TTCRJ12] Long Tran-Thanh, Archie Chapman, Alex Rogers, and Nicholas R. Jennings. Knapsack  
583 based optimal policies for budget-limited multi-armed bandits. In *Twenty-Sixth AAAI*  
584 *Conference on Artificial Intelligence (AAAI-12) (22/07/12 - 22/07/12)*, pages 1134–  
585 1140, April 2012.

586 [VBW15] Sofia Villar, Jack Bowden, and James Wason. Multi-armed bandit models for the  
587 optimal design of clinical trials: Benefits and challenges. *Statistical Science*, 30:199–  
588 215, 05 2015.

589 [WBSJ21] Lequn Wang, Yiwei Bai, Wen Sun, and Thorsten Joachims. Fairness of exposure  
590 in stochastic bandits. In Marina Meila and Tong Zhang, editors, *Proceedings of the*  
591 *38th International Conference on Machine Learning*, volume 139 of *Proceedings of*  
592 *Machine Learning Research*, pages 10686–10696. PMLR, 18–24 Jul 2021.

593 [WWJ21] Zhenlin Wang, Andrew Wagenmaker, and Kevin Jamieson. Best arm identification  
594 with safety constraints, 2021.

595 [WZZ23] Yuhang Wu, Zeyu Zheng, and Tingyu Zhu. Best arm identification with fairness  
596 constraints on subpopulations. In *2023 Winter Simulation Conference (WSC)*, pages  
597 540–551. IEEE, 2023.

598 [ZY24] Yafei Zhao and Long Yang. Constrained contextual bandit algorithm for limited-  
599 budget recommendation system. *Engineering Applications of Artificial Intelligence*,  
600 128:107558, 2024.

# Appendix

## Table of Contents

---

<b>A</b>	<b>Notations</b>	<b>16</b>
<b>B</b>	<b>Additional discussion on problem setting</b>	<b>17</b>
B.1	Extended related work . . . . .	17
B.2	Motivations: Reductions to and generalisations of existing settings . . . . .	18
<b>C</b>	<b>Strong duality and the Lagrangian multiplier: Proof of Theorem 1</b>	<b>19</b>
<b>D</b>	<b>Lagrangian Relaxation of Projection Lemma: Proof of Theorem 2</b>	<b>21</b>
<b>E</b>	<b>Characterization of the unique optimal policy: Proof of Theorem 3</b>	<b>24</b>
<b>F</b>	<b>Lagrangian Lower Bound for Gaussians: Proof of Theorem 4</b>	<b>25</b>
F.1	Bounds on Sample complexity: Proof of Corollary 1 Part (a) . . . . .	26
F.2	Impact of unknown linear constraints: Proof of Corollary 1 Part (b) . . . . .	28
<b>G</b>	<b>Sample Complexity upper bounds (Analysis of algorithms)</b>	<b>31</b>
G.1	Stopping Criterion . . . . .	31
G.2	Upper Bound of LATS . . . . .	33
G.3	Upper Bound for LAGEX . . . . .	35
G.4	Applications to existing problems . . . . .	39
<b>H</b>	<b>Constraint violations during exploration</b>	<b>41</b>
H.1	Upper Bound on Constraint Violation . . . . .	41
H.2	Experimental results . . . . .	42
H.3	Experiment on IMDB dataset . . . . .	42
<b>I</b>	<b><math>\epsilon</math>-good policies under unknown linear constraints</b>	<b>43</b>
<b>J</b>	<b>Technical results and known tools in BAI and pure exploration</b>	<b>44</b>
J.1	Concentration lemma for constraints . . . . .	44
J.2	Useful results from BAI and pure exploration literature . . . . .	45
J.3	Useful definitions and theorems from literature on continuity of convex functions	46

---

Notation	Definition
$\Delta_K$	K-simplex
$T$	Time Horizon
$K$	Number of Arms
$\mathbf{A}$	True constraint set
$\mathcal{F}$	True feasible set w.r.t $\mathbf{A}$ , $\mathcal{F} = \{\mathbf{A} \in \mathbb{R}^{N \times K} : \mathbf{A}\boldsymbol{\pi} \leq 0\}$
$\tilde{\mathbf{A}}_t$	Pessimistic estimate of constraint set at time t, $\tilde{\mathbf{A}}_t = \hat{\mathbf{A}} - f(t, \delta) \ \omega_t\ _{\Sigma_t^{-1}}$
$\tilde{\mathcal{F}}_t$	Estimated feasible set w.r.t pessimistic estimate $\tilde{\mathbf{A}}$ at time t, $\mathcal{F} = \{\tilde{\mathbf{A}}_t \in \mathbb{R}^{N \times K} : \tilde{\mathbf{A}}_t \boldsymbol{\pi} \leq 0\}$
$\mathbf{A}$	The action set of K possible choices
$\omega_t$	Policy chosen at time t
$a_t$	Action at time t among K possible actions
$N$	Number of Constraints
$\Gamma$	Slack
$\sigma^2$	Variance of the reward distribution (Gaussian) of arms
$T$	Time Horizon
$r_t, c_t$	Reward and cost observed at time t
$\delta$	Chosen confidence level
$l_t$	The Lagrangian multiplier at time t
$\Sigma_t$	The covariance matrix (Gram matrix) at round t
$\Lambda_{\mathcal{F}}(\boldsymbol{\mu})$	Set of alternative (confusing) instances for bandit instance $\boldsymbol{\mu}$
$\Lambda_{\tilde{\mathcal{F}}}(\boldsymbol{\mu})$	Estimated set of alternative (confusing) instances for bandit instance $\boldsymbol{\mu}$
$\nu_{\mathcal{F}}(\boldsymbol{\pi}_{\mathcal{F}}^*)$	Neighbourhood set around $\boldsymbol{\pi}_{\mathcal{F}}^*$
$\nu_{\tilde{\mathcal{F}}}(\boldsymbol{\pi}_{\tilde{\mathcal{F}}}^*)$	Neighbourhood set around $\boldsymbol{\pi}_{\tilde{\mathcal{F}}}^*$
$\tau_\delta$	Stopping time
$\boldsymbol{\pi}_{\mathcal{F}}^*$	True optimal policy w.r.t actual constraint set $\mathbf{A}$
$\boldsymbol{\pi}_{\tilde{\mathcal{F}}}^*$	Optimal Policy for the estimated feasible set



## B Additional discussion on problem setting

### B.1 Extended related work

**Historical pioneering works.** Literature on bandits has come a long way since the problem of optimal sequential sampling started with the works of [BB62] and [Pau64] with the assumption of the populations being normally distributed. To talk about pure exploration setting, [EDMM02b], [BMS10] should be mentioned as the first ones who worked in this specific setting for stochastic bandits.

**Existing work on adapting known constraints.** In Multi-armed bandit literature, people often introduce constraints as a notion of safety where they impose known constraints on the chosen arm or on the exploration process. [WWJ21] considers pure strategy (only one co-ordinate as chosen action) and imposes a safety threshold on the linear cost feedback of the chosen arm. On the other hand, the setting considered in [CBJD23] is closer as it tracks an optimal policy w.r.t to a known set of known constraints. On the other hand, [LLSY21] (Improvement over [PGBJ20] in MAB setting) generalized the known constraint regret minimization setting by assuming existence of a set of general constraints. Our work captures the hardness of not knowing the constraint set while tracking the lower bound and also in sample complexity upper bounds of Algorithm 1 and 2. Our work also introduce *shadow price* as a novel term in pure exploration literature which characterises the extra cost that arises due to tracking the unknown constraints.

**Learning unknown constraints.** [LTHK22] considers constrained linear best-arm identification arm are vectors with known rewards and a single unknown constraint (representing preferences) on the actions. Works on adapting to unknown constraints is discussed in the related work section of the main paper.

**Transductive Linear Bandit.** In this setting [FJR19] [CWM<sup>+</sup>22b] studies this setting with unknown linear constraints where we have to find the best safe arm in a finite set  $\mathcal{Z}$  different than actual arm set  $\mathcal{A}$ . Our setting generalises the setting in the sense that the finite feasible set  $\mathcal{Z}$  is not static, rather we track  $\mathcal{Z}_t$  per time step  $t \in \mathbb{N}$  and explore within that set to find the optimal allocation. At the end of exploration after hitting the stopping criterion at  $\tau_\delta$  the agent recommends the optimal policy inside the set  $\mathcal{Z}_{\tau_\delta}$ .

**Regret Minimization with Unknown constraints.** In bandit literature, constraints are often introduced in the setting to study regret minimization. [MAAT20] studies regret minimization using Linear Thompson Sampling (LTS) imposing known safety constraints on the chosen action. [AAT19] studies contextual bandits under unknown and unobserved linear constraint, whereas [KGAYR17] [PGBJ20] studies UCB based algorithms for regret minimization for linear bandits which assumes existence of a safe action space in case of unknown anytime linear constraint. In line with these, recent works [HTA24] [PGB24] [SCBC23] improved on regret guarantees and the first one relaxed the assumption of existence of a pessimistic safe space. [CGS22a] introduces doubly-optimistic setting to study safe linear bandit (SLB). [LLSY21] generalised the setting of [PGBJ20] not only relaxing the condition of having a safe action but also considered a set of general constraints and also captured the notion of both anytime and end-of-time constraints which we also see in [CBJD23]. [LLSY21] also shows the trade-off between maximising reward or minimizing regret and constraint violation using Lyapunov drift. In this work we do not focus on regret guarantees but finding the optimal policy with sample complexity as least as possible while tracking and satisfying a set of unknown linear constraints.

**BAI with Fairness Constraint.** Considering fairness constraint in our setting can be an interesting application to our setting. Recently [WZZ23] studied Best Arm Identification with fairness Constraints on Subpopulations (BAICS), where they have discussed the trade-off in the standard BAI complexity if there are finite number of subpopulations are given and the best chosen arm must perform well (not too bad) on all those subpopulations. Another important line of work [WBSJ21] [SJ19] explores regret analysis of BAI with positive merit-based exposure of fairness constraints where the chosen policy has to satisfy some fairness constraint across all its indices. Our setting comes as a direct application to these settings. Further discussion in Section B.2.

**BAI with Knapsack constraint.** While the existing literature on bandit with knapsack [BKS18] [AD16] [ISSS22], [AD14] [ADL16] [SS18] [Ma] focused on mainly regret minimization, our setting aligns more as a special case of the Optimal Arm identification with Knapsack setting in [LZYL23,

687 TTCRJ12, LSY21]. Though we aim to find the best policy rather than a specific arm in the constraint  
688 space. Our setting should be considered as a special case of these settings. Further discussion in  
689 Section B.2.

690 **Algorithms on Pure exploration.** Algorithm 1 is an extension of the Track and Stop(TnS) strategy  
691 from [GK16], while the motivation for Algorithm 2 comes from the Gamified Explorer strategy  
692 from [DKM19b] where the lower bound is treated as a zero-sum game between the allocation and  
693 the instance player. We refer to [GK21b],[GK16],[KCG16],[DK19], [BV04],[JMKG21] etc for  
694 important concentration inequalities, tracking lemmas.

695 **Dose-finding and Thresholding Bandits.** Another special case of our setting is Dose-finding or  
696 Thresholding bandits in structured MAB literature [CLK<sup>+</sup>14] generalized the problem, then a line  
697 of work [AKR21b] [GMRM18] [CMC21] aims to find the maximum safe dose for a specific drug  
698 in early stages of clinical trials. In some sense our setting generalizes this setting. If we have to  
699 administer more than drugs to a patient, our setting generalises to track the best possible proportion in  
700 which the drugs should be administered with maximum efficacy. Further discussions in Section B.2.

## 701 B.2 Motivations: Reductions to and generalisations of existing settings

702 Before delving into the details of the lower bounds and algorithms, we first clarify our motivation by  
703 showing how different setups studied in literature and their variations are special case of our setting.

704 **Thresholding Bandits.** Our setting encompasses the thresholding bandit problem [AKR21a]. Thresh-  
705 olding bandit is motivated from the safe dose finding problem in clinical trials, where one wants to  
706 identify the highest dose of a drug that is below a known safety level. This has also motivated the  
707 studies on safe arm identification [WWJ21]. Our setting generalises it further to detect the dose of the  
708 drug with highest efficacy while it is still below the safety level. We can formulate it as identifying  
709  $\pi^* = \arg \max \mu^T \pi$ , such that  $\mathbf{I} \pi \leq \mathbf{I} \theta$ . Rather, generalising the classical thresholding bandits,  
710 our formulation can further model the safe doses for the optimal cocktail of drugs, and  $\theta$  can have  
711 different values across drugs, i.e we can consider different thresholds for different drugs.

712 **Optimal policy under Knapsack.** Bandits under knapsack constraints have been studied both in best-  
713 arm identification [LZYL23, TTCRJ12, LSY21] and regret minimisation [BKS18, AD16, ISSS22,  
714 AD14, ADL16, SS18, Ma] literature. BAI under knapsacks is motivated by the fact that detecting an  
715 optimal arm might have additional resource constraints in addition to the number of required samples.  
716 This has led to study of BAI with knapsacks only under fixed-budget settings [LZYL23]. But as in  
717 regret-minimisation literature [SS18, Ma], one might want to recommend a policy that maximises  
718 utility while satisfying knapsack constraints. For example, we want to manage caches where the  
719 recommended memory allocation should satisfy a certain resource budget. Thus, the recommended  
720 policy has to satisfy  $\pi_\tau^* = \arg \max_{\pi \in C_A} \hat{\mu}_{\tau_\delta}^T \pi$ , where  $C_A \triangleq \{\mathbf{A} \pi_{\tau_\delta} \leq c\}$ . Naturally, this is a  
721 special case of our problem setting.

722 **Feasible arm selection.** We look at the pure exploration problem of feasible arm selection studied by  
723 [KSS18]. Here, we think of a problem of workers having a multi-dimension vector representation  
724 where each index denotes the accuracy of that worker being able to identify a specific class label in  
725 a classification task in hand. The problem turns to be a feasible arm selection from a simple BAI  
726 problem when we impose a feasibility constraint that for example, the chosen worker should show  
727 more than 90% accuracy across all labels. We can generalise this setting in the sense that we are  
728 now not looking for a specific worker, rather we want to make a team of workers that has the highest  
729 utility. The recommended policy at time  $t \in \mathbb{N}$ ,  $\max_{\pi \in \Delta_{K-1}} \mu_t^T \pi$  such that  $f^T \pi \geq \tau$  where  $\tau$  is the  
730 desired threshold level. The generalisation of the setting pitch in as thresholds of  $\tau$  can have different  
731 values corresponding to different workers.

732 **BAI with fairness across sub-populations.** The Best Arm Identification with fairness Constraints  
733 on Sub-population (BAICS) studied in [WZZ23] aims on selecting an arm that must be fair across  
734 all sub-populations rather than the whole population in standard BAI setting. Let, there are  $l$  sub-  
735 populations and  $\mu_a$  are the means corresponding to the  $a$ -th arm. Finding only the optimal arm  
736  $K_{\text{BAI}} = \arg \max_{k \in [K]} \mu_k$  may not be enough because it may not perform equally good for all  
737 the  $l$  sub-populations. Then the arm should belong to a set  $C := \{k \in [K] | \mu_{k,m} \geq 0, m \in [l]\}$   
738 where the observation for arm  $k$  and population  $m$  comes from  $\mathcal{N}(\mu_{k,m}, 1)$  It ensures that the  
739 chosen arm does not perform *too bad* for any sub-population. Let us think of a problem where  
740 there are  $l$  sub-groups of patients and we have  $K$  number of drugs to administer with reward means

741  $\mu_k, k \in [K]$ . We are looking for a combination of drugs rather than a single drug to administer as  
 742  $\pi^* = \arg \max_{\pi \in \Delta_K} \mu^T \pi$  such that  $\mathbb{1}_{\mu_m \geq 0}^T \pi = 1, \forall m \in [l]$ . Thus, BAICS is a special case of ours.

743 **Fairness of exposure in bandits.** [WBSJ21] introduced positive merit based exposure of fairness  
 744 constraints [SJ19] in stochastic bandits standing against the winner-takes-all allocation strategy that  
 745 are historically studied. The chosen allocation in this setting should satisfy the fairness constraint  
 746  $\frac{\pi_a^*}{f(\mu_a^*)} = \frac{\pi_{a'}^*}{f(\mu_{a'}^*)}, \forall a' \in [K]$  where  $f(\cdot)$  transform reward of an arm to a positive merit. Though  
 747 [WBSJ21] studied this setting in regret analysis, this setting in BAI setting is a direct application  
 748 of our setting as we are looking for an optimal policy  $\pi^* = \arg \max \mu^T \pi$  such that  $\pi^*$  satisfies  
 749  $\mathbf{A}_{f(\mu)}^T \pi = 0$  where  $\mathbf{A}_{f(\mu)}$  is of order  $\frac{K(K-1)}{2} \times K$  and  $\mathbf{A}_{f(\mu)}$  is expressed as,

$$(\mathbf{A}_{f(\mu)})_{ij} = \begin{cases} \frac{1}{f(\mu_a)} & \text{if } aK - \frac{1}{2}(a-1)(a-2) \leq j \leq aK - \frac{1}{2}a(a-1) \text{ and } i = a, \\ \frac{-1}{f(\mu_j)} & \text{if } aK - \frac{1}{2}(a-1)(a-2) \leq j \leq aK - \frac{1}{2}a(a-1) \text{ and } a < i \leq K, \\ 0 & \text{otherwise.} \end{cases}$$

750 For example, when  $K = 3$ ,  $\mathbf{A}_{f(\mu)} = [[\frac{1}{\mu_1}, -\frac{1}{\mu_2}, 0], [0, \frac{1}{\mu_2}, -\frac{1}{\mu_3}], [\frac{1}{\mu_1}, 0, -\frac{1}{\mu_3}]]$ .

## 751 C Strong duality and the Lagrangian multiplier: Proof of Theorem 1

752 **Theorem 1.** *For a bounded sequence of  $\{l_t\}_{t \in \mathbb{N}}$ , strong-duality holds for the optimisation problem*  
 753 *stated in Equation (7) i.e.*

$$\inf_{l \in \mathbb{R}^d} \min_{\mathbf{A}' \in \mathcal{C}} \sup_{\omega \in \hat{\mathcal{F}}} \inf_{\lambda \in \Lambda_{\hat{\mathcal{F}}}(\mu)} \omega^\top d(\mu, \lambda) - l^T \tilde{\mathbf{A}} \omega = \sup_{\omega \in \hat{\mathcal{F}}} \min_{l \in \mathcal{L}} \inf_{\lambda \in \Lambda_{\hat{\mathcal{F}}}(\mu)} \omega^\top d(\mu, \lambda) - l^T \tilde{\mathbf{A}} \omega. \quad (10)$$

754 Here,  $\mathcal{L} \triangleq \{l \in \mathbb{R}^d \mid 0 \leq \|l\|_1 \leq \frac{1}{\gamma} T^{-1}(\hat{\mu})\}$ , where  $\gamma \triangleq \min_{i \in [1, d]} \{-\tilde{\mathbf{A}}^i \omega^*\}$ , i.e. the minimum  
 755 slack for pessimistic constraints w.r.t. the optimal allocation.

756 *Proof.* This proof involves three steps. In the first step we prove convexity and other properties of  
 757 the sets involved in the main optimisation problem 10. Then in the next step we show that Slater's  
 758 sufficient conditions hold for  $\pi$  as a consequence of these properties. Once we prove the unique  
 759 optimality of  $\pi$  we state bounds on the L1-norm of the Lagrangian multiplier. We conclude by  
 760 establishing strong duality and proving the statement of the theorem.

761 **Step 1: Properties of perturbed feasible set and alt-set.** Let us first check the properties of  $\hat{\mathcal{F}}$ ,  
 762  $\Lambda_{\hat{\mathcal{F}}}(\mu)$  and  $\nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)$ . For that, let us remind the definitions of these sets. The estimated feasible set is  
 763 defined as  $\hat{\mathcal{F}} \triangleq \{\pi \in \Delta_K : \tilde{\mathbf{A}} \pi \leq 0\}$ . The set of alternative (confusing) instances for the optimal  
 764 policy  $\pi_{\hat{\mathcal{F}}}^*$  is  $\Lambda_{\hat{\mathcal{F}}}(\mu) \triangleq \{\lambda \in \mathbb{D} : \max_{\pi \in \hat{\mathcal{F}}} \lambda^T \pi > \lambda^T \pi_{\hat{\mathcal{F}}}^*\}$ . For  $\pi'$  being a neighbour of  $\pi_{\hat{\mathcal{F}}}^*$  or in  
 765 other words, an extreme point in  $\hat{\mathcal{F}}$ , we decompose the alternative set as the union of half-spaces as,

$$\Lambda_{\hat{\mathcal{F}}}(\mu) = \bigcup_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \left\{ \lambda : \lambda^T (\pi_{\hat{\mathcal{F}}}^* - \pi') < 0 \right\}$$

766 We should note,  $\pi'$  shares at least  $(K-1)$  active constraints with  $\pi_{\hat{\mathcal{F}}}^*$ . It is clear that  $\hat{\mathcal{F}}$  is bounded  
 767 and convex in  $\pi$ . Since, convex combination of any two extreme point  $\pi'_1, \pi'_2$  in the neighbourhood  
 768 of the optimal policy  $\nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)$  also shares  $K-1$  active constraints with  $\pi^*$ , so  $\nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)$  is convex in  
 769  $\pi'$ .

770 Let,  $\pi'_1$  and  $\pi'_2$  are two policies in the neighbourhood of  $\pi^*$ , such that for any alternative instance  $\lambda$ ,  
 771  $\lambda^T (\pi'_1 - \pi'_2) \geq 0$  Above equation implies that the policy  $\pi'_1$  is closer to the optimal policy in the  
 772 neighbourhood than the policy  $\pi'_2$

773 Therefore, any convex combination of these neighbourhood policy,  $\lambda^T (\pi^* - (a\pi'_1 + (1-a)\pi'_2)) =$   
 774  $\lambda^T (\pi^* - \pi'_2) - a\lambda^T (\pi'_1 - \pi'_2) \leq c$ . Therefore, the set  $\Lambda_{\hat{\mathcal{F}}}(\mu)$  is also bounded and convex in  $\pi$

775 Also, since we are working with pessimistic estimate of  $\mathbf{A}$ , the set  $\hat{\mathcal{F}}$  will always be non-empty,  
 776 because we will find at least one  $\tilde{\mathbf{A}}_0$  which is non-singular and it's inverse exists.

777 **Step 2: Slater's condition.** From step 1 of this proof we have the following properties

- 778 1.  $\hat{\mathcal{F}}$  is non-empty, bounded and convex in  $\pi$ .  
 779 2. The perturbed neighbourhood  $\nu_{\hat{\mathcal{F}}}(\pi^*)$  is convex for any  $\pi' \in \nu_{\hat{\mathcal{F}}}(\pi^*)$   
 780 3.  $\Lambda_{\hat{\mathcal{F}}}(\mu)$  is also bounded and convex in  $\pi$ .

781 Leveraging these three results we claim that there exists a  $\pi^*$  that uniquely solve the optimisation  
 782 problem in Equation (10) and satisfy the constraints with strict inequality. As a consequence of this  
 783 we claim Slater's sufficient conditions hold.

784 **Step 3: Bound on the Lagrangian multiplier.** Here, we try to bound the L-1 norm of the Lagrangian  
 785 multiplier. Since,  $\|l\|_1$  cannot be less than 0, then we already have a lower bound.

786 Now we refer to lemma 4 for the upper bound. An immediate implication of this result is that for  
 787 any dual optimal solution  $\mu^*$ , we have  $\|\mu^*\|_1 \leq \frac{1}{\gamma} (f(\bar{x}) - q^*)$ . Since Slater's conditions hold in our  
 788 case for  $\pi$ , we can write that the optimal solution of the Lagrangian dual,

$$0 \leq \|l_t^*\|_1 \leq \frac{1}{\gamma} \left( \mathcal{D}(\omega_t^*, \hat{\mu}_t, \hat{\mathcal{F}}_t) - \underbrace{\mathcal{D}(\pi^*, \hat{\mu}_t, \hat{\mathcal{F}}_t)}_{\text{not tractable}} \right) \leq \frac{1}{\gamma} \mathcal{D}(\omega_t^*, \hat{\mu}_t, \hat{\mathcal{F}}_t, l_{t-1}^*) = \frac{1}{\gamma} T_{\hat{\mathcal{F}}}^{-1}(\hat{\mu})$$

$$\text{where, } \gamma \triangleq \min_{i \in [1, d]} \{-\tilde{\mathbf{A}}^i \omega^*\}$$

789 Where,  $\pi^*$  is the pure-exploration solution. We can replace the dual function with primal  
 790  $\mathcal{D}(\omega_t^*, \hat{\mu}_t, \hat{\mathcal{F}}_t, l_{t-1}^*)$  because it is always upper bounded by the dual function, so we don't have  
 791 to explicitly calculate the dual function. Though it is not tractable anyway due to not knowing the  
 792 pure exploration solution.

793 **Step 4: Establishing strong duality.** Therefore the domain of the Lagrangian multiplier is also  
 794 bounded and convex. So again we say that  $l_t^*$  uniquely minimises Equation 10. We define  $\mathcal{L} \triangleq \{l \in$   
 795  $\mathbb{R}^d \mid 0 \leq \|l\|_1 \leq \frac{1}{\gamma} T_{\hat{\mathcal{F}}}^{-1}(\hat{\mu})\}$ , where  $\gamma \triangleq \min_{i \in [1, d]} \{-\tilde{\mathbf{A}}^i \omega^*\}$  Then according to **Heine-Borel's**  
 796 **theorem** (Theorem 10) we can say that these sets are compact as well. We can then conclude  
 797 that Strong duality holds which means that it perfectly make sense of solving the Lagrangian dual  
 798 formulation of the primal optimisation problem because there is no duality gap. We later on will  
 799 consider this formulation as two player zero sum game. Due to strong duality we claim that the agent  
 800 while playing this game, Nash equilibrium will be eventually established.

801 Now that everything is put into place we can conclude with the very statement of the theorem that  
 802 due to strong duality the following holds

$$\inf_{l \in \mathbb{R}_+^d} \min_{\mathbf{A}' \in \mathcal{C}} \sup_{\omega \in \hat{\mathcal{F}}} \inf_{\lambda \in \Lambda_{\hat{\mathcal{F}}}(\mu)} \omega^\top d(\mu, \lambda) - l^\top \tilde{\mathbf{A}} \omega = \sup_{\omega \in \hat{\mathcal{F}}} \min_{l \in \mathcal{L}} \inf_{\lambda \in \Lambda_{\hat{\mathcal{F}}}(\mu)} \omega^\top d(\mu, \lambda) - l^\top \tilde{\mathbf{A}} \omega.$$

803

□

## D Lagrangian Relaxation of Projection Lemma: Proof of Theorem 2

**Theorem 2.** For a sequence  $\{\hat{\mathcal{F}}_t\}_{t \in \mathbb{N}}$  and  $\{\hat{\lambda}_t\}_{t \in \mathbb{N}}$ , we show that (a)  $\lim_{t \rightarrow \infty} \hat{\mathcal{F}}_t \rightarrow \mathcal{F}$ , (b)  $\lambda^*$  is unique, and (c)  $\lim_{t \rightarrow \infty} \hat{\lambda}_t \rightarrow \lambda^*$ . Thus, for any  $\omega \in \mathcal{F}$  and  $\mu$ ,  $\lim_{t \rightarrow \infty} \mathcal{D}(\omega, \mu, \hat{\mathcal{F}}) \rightarrow \mathcal{D}(\omega, \mu, \mathcal{F})$  where  $\lambda^*$  is such that for any  $\lambda \in \mathcal{D} : \inf_{\lambda \in \Lambda_{\mathcal{F}}(\mu)} \omega^T d(\mu, \lambda) \geq \omega^T d(\mu, \lambda^*)$ .

*Proof.* Here, we prove the three parts of the theorem consecutively.

**Statement (a): Convergence of the limit**  $\lim_{t \rightarrow \infty} \hat{\mathcal{F}}$ . To begin with the proof of the first statement of Theorem 2 we leverage the results stated in Theorem 11. Let  $H(\tilde{\mathbf{A}}) \triangleq \{\pi \in \Delta_K : \tilde{\mathbf{A}}\pi \leq 0\}$  and the set function  $\tilde{\mathbf{A}} \rightarrow H(\tilde{\mathbf{A}}) \cap C$  where  $C = \hat{\mathcal{F}}$  is a non-empty compact (proven in Section C subset of  $\Delta_K$ ). Then the set  $H(\tilde{\mathbf{A}}) \cap C$  can be written as

$$H(\tilde{\mathbf{A}}) \cap C = \{\pi \in \hat{\mathcal{F}} : \tilde{\mathbf{A}}\pi \leq 0\}$$

To apply Theorem 11,  $\{\tilde{\mathbf{A}}^r, r \in \mathbb{N}\}$ , must be a convergent sequence of affine function. It is evident that  $\tilde{\mathbf{A}}^r$  for any  $r \in \mathbb{N}$  is an affine function since  $\mathbf{A}$  is linear in  $\mathbf{A}$  and the induced pessimism works as a translation. Then we can proceed to the next part of the proof of statement 1 where we prove that  $\{\tilde{\mathbf{A}}^r\}_{r \in \mathbb{N}}$  is a convergent sequence of functions. For ease of notation we will denote  $\tilde{\mathbf{A}}_t$  for the  $t$ -th element of the sequence  $\{\tilde{\mathbf{A}}^r\}_{r \in \mathbb{N}}$  for  $t \in \mathbb{N}$ .

The definition of the confidence radius for any constraint  $i \in [d]$  follows from the Definition 2 as  $f(\delta, t) \triangleq 1 + \sqrt{\frac{1}{2} \log \frac{K}{\delta} + \frac{1}{4} \log \det \Sigma_t}$ . It is evident from the definition that  $f(t, \delta)$  is a non-decreasing function w.r.t time and it grows with order of at least  $\mathcal{O}(\sqrt{\log t})$ .

We have from the definition of the confidence set, for all  $i \in [d]$

$$\begin{aligned} & \mathbb{P} \left( \hat{\mathbf{A}}_t^i - f(t, \delta) \|\omega_t\|_{\Sigma_t^{-1}} \leq \mathbf{A}^i \leq \hat{\mathbf{A}}_t^i + f(t, \delta) \|\omega_t\|_{\Sigma_t^{-1}} \right) \geq 1 - \delta \\ \implies & \mathbb{P} \left( -\frac{f(t, \delta) \|\omega_t\|_{\Sigma_t^{-1}}}{\sigma(\hat{\mathbf{A}}_t^i)} \leq \frac{\hat{\mathbf{A}}_t^i - \mathbf{A}^i}{\sigma(\hat{\mathbf{A}}_t^i)} \leq \frac{f(t, \delta) \|\omega_t\|_{\Sigma_t^{-1}}}{\sigma(\hat{\mathbf{A}}_t^i)} \right) \geq 1 - \delta \\ \implies & \mathbb{P} \left( -\frac{f(t, \delta) \|\omega_t\|_{\Sigma_t^{-1}}}{\sigma(\hat{\mathbf{A}}_t^i)} \leq Z \leq \frac{f(t, \delta) \|\omega_t\|_{\Sigma_t^{-1}}}{\sigma(\hat{\mathbf{A}}_t^i)} \right) \geq 1 - \delta \\ \implies & 2\Phi \left( \frac{f(t, \delta) \|\omega_t\|_{\Sigma_t^{-1}}}{\sigma(\hat{\mathbf{A}}_t^i)} \right) \geq 2 - \delta \\ \implies & \frac{f(t, \delta) \|\omega_t\|_{\Sigma_t^{-1}}}{\sigma(\hat{\mathbf{A}}_t^i)} \geq \Phi^{-1} \left( 1 - \frac{\delta}{2} \right) \\ \implies & f(t, \delta) \|\omega_t\|_{\Sigma_t^{-1}} \geq \sigma(\hat{\mathbf{A}}_t^i) \Phi^{-1} \left( 1 - \frac{\delta}{2} \right) \end{aligned}$$

where  $Z \triangleq \frac{\hat{\mathbf{A}}_t^i - \mathbf{A}^i}{\sigma(\hat{\mathbf{A}}_t^i)}$  and  $\Phi(\cdot)$  is the cumulative distribution function of  $\mathcal{N}(0, 1)$  distribution.

$\lim_{t \rightarrow \infty} f(t, \delta) \|\omega_t\|_{\Sigma_t^{-1}} \rightarrow 0$  since  $\sigma(\hat{\mathbf{A}}_t^i) = \mathcal{O}(\sqrt{\frac{\log t}{t}})$ . Leveraging CLT at this point we say

$$\hat{\mathbf{A}}_t^i \xrightarrow{d} \mathbf{A}^i, \forall i \in [d]$$

Then by Slutsky's theorem [Slu25], we conclude  $\hat{\mathbf{A}}_t^i - f(t, \delta) \|\omega_t\|_{\Sigma_t^{-1}} \xrightarrow{d} \mathbf{A}^i, \forall i \in [d]$ .

It implies that  $\{\tilde{\mathbf{A}}^r\}_{r \in \mathbb{N}}$  is a convergent sequence of function for  $\mathbf{A}$ . Now, we use Theorem 11 and get the following properties of the feasible set.

$$1. H(\tilde{\mathbf{A}}) \cap C \subset \lim_{r \rightarrow \infty} H(\tilde{\mathbf{A}}^r) \cap C$$

- 828 2.  $\lim_{r \rightarrow \infty} H(\tilde{\mathbf{A}}^r) \cap C$  is a closed convex superset of  $H(\tilde{\mathbf{A}}) \cap C$ .  
 829 3.  $H(\tilde{\mathbf{A}}) \cap C$  has non-empty interior because of the feasibility condition and no component in  
 830  $\tilde{\mathbf{A}}$  is identically 0.

$$\lim_{r \rightarrow \infty} H(\tilde{\mathbf{A}}^r) \cap C = H(\tilde{\mathbf{A}}) \cap C$$

- 831 4. Even if the set  $H(\tilde{\mathbf{A}}) \cap C$  has empty interior or some component if  $\tilde{\mathbf{A}}$  is identically zero, by  
 832 the last statement of the Theorem 11 we can say for any closed convex set  $Q$  of  $H(\tilde{\mathbf{A}}) \cap C$   
 833 we can design the function  $\{\tilde{\mathbf{A}}^r\}$  in such a way that  $\lim_{r \rightarrow \infty} H(\tilde{\mathbf{A}}^r) \cap C$  includes  $Q$ .

834 As the convergence of  $\tilde{\mathbf{A}}_t$  is guaranteed now asymptotically, we can guaranty convergence of the  
 835 following limit  $\lim_{t \rightarrow \infty} \hat{\mathcal{F}}_t \rightarrow \mathcal{F}$ .

836 **Statement (b) : Proof of Uniqueness of  $\lambda^*$**  Here, we try to prove if there exists a confusing instance  
 837  $\lambda^* \in \Lambda_{\hat{\mathcal{F}}}(\mu)$  which uniquely minimises the the function  $\mathcal{D}(\cdot)$  defined as

$$\mathcal{D}(\omega, \mu, \hat{\mathcal{F}}) \triangleq \inf_{l \in \mathcal{L}} \inf_{\lambda \in \Lambda_{\hat{\mathcal{F}}}(\mu)} \omega^T d(\mu, \lambda) - l^T \tilde{\mathbf{A}} \omega$$

838 We can observe that only the leading quantity on the R.H.S associated with the KL is dependent on  
 839  $\lambda$ . So, in this proof we will only show that  $\lambda^*$  minimizes the KL divergence uniquely and since  
 840 the KL is linearly dependent on the expression, proving this will be enough to ensure uniqueness of  $\lambda^*$ .  
 841

842 Now, from the properties of KL we know that  $d(\mu, \lambda)$  is convex on the pair  $(\mu, \lambda)$ . But it is also  
 843 strictly convex on  $\lambda$  if  $\text{supp}(\lambda) \subseteq \text{supp}(\mu)$  which is true in our case, since  $\mu, \lambda \in \mathcal{D} \subseteq \mathbb{R}^k$ .

844 Let us assume there are two local minima  $\lambda_1$  and  $\lambda_2$ , with the condition,

$$d(\mu, \lambda_1) \leq d(\mu, \lambda_2)$$

845 Then, we can write from the property of strict convexity, for some  $\{h : 0 < h < 1\}$ ,

$$d(\mu, h\lambda_1 + (1-h)\lambda_2) < hd(\mu, \lambda_1) + (1-h)d(\mu, \lambda_2)$$

846 Now, from the assumed condition on  $\lambda_1$  and  $\lambda_2$ , we can write —

$$\begin{aligned} d(\mu, \lambda_1) &\leq d(\mu, \lambda_2) \\ \implies hd(\mu, \lambda_1) &\leq hd(\mu, \lambda_2), \text{ since } h > 0 \\ \implies hd(\mu, \lambda_1) + (1-h)d(\mu, \lambda_2) &\leq hd(\mu, \lambda_2) + (1-h)d(\mu, \lambda_2) \\ \implies hd(\mu, \lambda_1) + (1-h)d(\mu, \lambda_2) &\leq d(\mu, \lambda_2) \end{aligned}$$

847 Putting this result in the strict convexity condition we get

$$d(\mu, h\lambda_1 + (1-h)\lambda_2) < d(\mu, \lambda_2)$$

848 which is a contradiction.

849 Thus, we can conclude that for a strictly convex function  $f(x)$  with  $\text{supp}(x)$  being convex as well, the  
 850 set of minimisers is either empty or singleton. Then, we can say  $\lambda^*$  uniquely minimizes the KL, or  
 851 say  $\mathcal{D}(\omega, \mu, \hat{\mathcal{F}})$ . Let us now once again remind the definition of perturbed alt-set  $\Lambda_{\hat{\mathcal{F}}}(\mu) \triangleq \{\lambda \in \mathbb{D} : \max_{\pi \in \hat{\mathcal{F}}} \lambda^T(\pi - \pi_{\hat{\mathcal{F}}}^*) > 0\}$ . Let us denote  $\nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)$  as the neighbourhood of  $\pi_{\hat{\mathcal{F}}}^*$ . Any  $\pi' \in \hat{\mathcal{F}}$  is  
 852 called a neighbour of  $\pi_{\hat{\mathcal{F}}}^*$ , if it is an extreme point of  $\hat{\mathcal{F}}$  and shares (K-1) active constraints with  $\pi_{\hat{\mathcal{F}}}^*$ .  
 853 Then, we can decompose the perturbed alt-set as  $\Lambda_{\hat{\mathcal{F}}}(\mu) = \bigcup_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \left\{ \lambda : \lambda^T(\pi_{\hat{\mathcal{F}}}^* - \pi') < 0 \right\}$ ,  
 854 which is a union of half-spaces for each neighbour. From this decomposition we can observe that  $\pi_{\hat{\mathcal{F}}}^*$   
 855 is not the optimal policy for  $\lambda$ , i.e,  $\{\exists \pi' \in \Lambda_{\hat{\mathcal{F}}}(\mu) : \lambda^T(\pi_{\hat{\mathcal{F}}}^* - \pi') < 0\}$ . Then, it follows similar  
 856 argument in [CBJD23] to argue that the most confusing instance w.r.t  $\mu$  lies in the boundary of the  
 857 normal cone, which lands us to Proposition 1.  
 858

859 For any  $\omega \in \hat{\mathcal{F}}$  and  $\mu \in \mathcal{D}$ , the following projection lemma holds for the Lagrangian relaxation,

$$\mathcal{D}(\omega, \mu, \hat{\mathcal{F}}) = \min_{l \in \mathcal{L}} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \min_{\lambda : \lambda^T(\pi_{\hat{\mathcal{F}}}^* - \pi') = 0} \omega^T d(\mu, \lambda) - l^T \tilde{\mathbf{A}} \omega. \quad (11)$$

860 **Statement (c): Convergence of the sequence**  $\{\hat{\lambda}_n\}_{n \in \mathbb{N}}$ . In known constraint setting the agent has  
 861 access to  $\mathcal{F}$ . That means there is the actual sequence  $\{\lambda_n\}_{n \in \mathbb{N}}$  for which  $\lambda_n \rightarrow \lambda^*$  as  $n \rightarrow \infty$  since  
 862  $\mathcal{D}(\omega, \mu, \hat{\mathcal{F}})$  is convex and continuous on  $\Lambda_{\hat{\mathcal{F}}}(\mu)$ . But in this setting we try to estimate  $\mathcal{F}$  as  $\hat{\mathcal{F}}_n$  at  
 863 each time step  $n \in \mathbb{N}$ . So there exists the  $\{\hat{\lambda}_n\}_{n \in \mathbb{N}}$  such that  $\hat{\lambda}_n \in \Lambda_{\hat{\mathcal{F}}_n}(\mu)$  and we have to ensure  
 864 it converges to the unique optimal  $\lambda^*$  i.e  $\lambda^* \in \Lambda_{\mathcal{F}}(\mu) \subseteq \lim_{n \rightarrow \infty} \Lambda_{\hat{\mathcal{F}}_n}(\mu)$  implies  $\{\hat{\lambda}_n\} \rightarrow \lambda^*$  as  
 865  $n \rightarrow \infty$

866 We use the fundamental theorem of limit to carry out this proof with the help of properties of the sets  
 867  $\hat{\mathcal{F}}$  and  $\Lambda$ . The properties we have already proven for these sets are

- 868 1.  $\hat{\mathcal{F}}_n$  for any  $n \in \mathbb{N}$  is a superset of  $\mathcal{F}$  due to the pessimistic choice of  $\mathbf{A}$ .
- 869 2.  $\hat{\mathcal{F}}_n$  is a non-empty compact subset of  $\Delta_K$  and  $\lim_{n \rightarrow \infty} \hat{\mathcal{F}}_n = \mathcal{F}$ .
- 870 3.  $\Lambda_{\hat{\mathcal{F}}_n}(\mu)$  is a closed convex set and it also is a superset of the real alt-set  $\Lambda_{\mathcal{F}}(\mu)$ .

871 Leveraging these properties we claim that for any  $\mu \in \mathcal{D}$ ,  $\lim_{n \rightarrow \infty} \Lambda_{\hat{\mathcal{F}}_n}(\mu) = \Lambda_{\mathcal{F}}(\mu)$ . Since we have  
 872 already proven uniqueness of  $\lambda$  in statement 2, we say  $\hat{\lambda}_n$  uniquely minimises  $\Lambda_{\hat{\mathcal{F}}_n}(\mu)$ . Now from  
 873 the  $(\epsilon, \delta)$ -definition of limits we say if  $\Lambda_{\hat{\mathcal{F}}_n}(\mu)$  is an  $\epsilon$ -cover of  $\Lambda_{\mathcal{F}}(\mu)$  for  $\epsilon > 0$ , then  $|\hat{\lambda}_n - \lambda^*| \leq \delta$   
 874 for  $\delta > 0$  sufficiently small. It implies for a sequence of  $\{\hat{\lambda}_n\}_{n \in \mathbb{N}}$  we claim  $\lim_{n \rightarrow \infty} \hat{\lambda}_n = \lambda^*$  i.e  
 875 the sequence convergence. Therefore we conclude by the statement itself

$$\{\hat{\lambda}_n\}_{n \in \mathbb{N}} \rightarrow \lambda^*$$

876 Hence, proved. □

## E Characterization of the unique optimal policy: Proof of Theorem 3

**Theorem 3.** For any  $\mu \in \mathcal{D}$ , the optimization problem  $\max_{\pi \in \hat{\mathcal{F}}} \mu^T \pi$  has a unique solution if  $\omega^*(\mu)$  satisfies the following conditions:

1. Both the sets  $\hat{\mathcal{F}}$  and  $\omega^*(\mu)$  are closed and convex.
2.  $\forall \mu \in \mathcal{D}$  and  $\omega \in \hat{\mathcal{F}}$ , the function  $\lim_{t \rightarrow \infty} \mathcal{D}(\omega, \hat{\mu}_t, \hat{\mathcal{F}})$  is continuous.
3. Reciprocal of the characteristic time function  $\lim_{t \rightarrow \infty} T_{\hat{\mathcal{F}}}^{-1}(\mu)$  is continuous  $\forall \mu \in \mathcal{D}$ .
4.  $\mu \in \mathcal{D} : \mu \rightarrow \omega^*(\mu)$  is upper hemi-continuous.

*Proof.* The theorem has four statements as the sufficient condition for the existence of unique optimal policy. So naturally we will dictate the proof structure in four steps and prove the statements one by one.

**Statement 1: Convexity of feasible space and optimal set function.** Let us first analyse the properties of  $\hat{\mathcal{F}}$ . For any two member of  $\omega_1, \omega_2 \in \hat{\mathcal{F}}$  satisfying  $\tilde{\mathbf{A}}\omega_1 \leq 0$  and  $\tilde{\mathbf{A}}\omega_2 \leq 0$ , their convex combination for any  $\alpha \in [0, 1]$ ,

$$\tilde{\mathbf{A}}(\alpha\omega_1 + (1 - \alpha)\omega_2) = \alpha\tilde{\mathbf{A}}\omega_1 + (1 - \alpha)\tilde{\mathbf{A}}\omega_2 \leq 0$$

Therefore we can say  $\hat{\mathcal{F}}$  is convex because it is closed under convex operation. We claim  $\hat{\mathcal{F}}$  is also closed since

1. The complement of  $\hat{\mathcal{F}}$ ,  $\hat{\mathcal{F}}^c \triangleq \{\pi \in \Delta_K : \tilde{\mathbf{A}}\pi > 0\}$  is an open set.
2. we have already proven the limit of  $\hat{\mathcal{F}}$  to be  $\mathcal{F}$  which is always contained by  $\hat{\mathcal{F}}$ .

The elements in the domain of optimal allocation set function must be included in  $\hat{\mathcal{F}}$ . So compactness of  $\omega^*(\mu)$  is a direct consequence of compactness of  $\hat{\mathcal{F}}$ .

**Statement 2: Continuity of limit.** We have already proven in Section D that  $\lim_{t \rightarrow \infty} \hat{\mathcal{F}} \rightarrow \mathcal{F}$ . Also by convexity of KL and CLT we claim  $\hat{\mu}_t \rightarrow \mu$  as  $t \rightarrow \infty$  and since  $\omega$  is linear in  $\mathcal{D}(\omega, \mu, \hat{\mathcal{F}})$  it will converge to  $\omega^*(\mu)$  as  $t \rightarrow \infty$ , also due to convergence of  $\hat{\mu}_t$ . Then we can say that the limiting value is same as the value if we plug in the limits in  $\mathcal{D}$  i.e  $\lim_{t \rightarrow \infty} \mathcal{D}(\omega, \hat{\mu}_t, \hat{\mathcal{F}}_t) = \mathcal{D}(\omega^*(\mu), \mu, \mathcal{F})$ . So we ensure the continuity of  $\lim_{t \rightarrow \infty} \mathcal{D}(\omega, \hat{\mu}_t, \hat{\mathcal{F}})$ .

**Statement 3: Continuity of limit of inverse sampling complexity.** This statement directly follows from the statement 2. Due to convexity of KL-divergence and convergence of  $\hat{\mathcal{F}}$ , the limiting value exists and it is equal to the inverse of characteristic time with the limiting value.

**Statement 4: Upper hemi-continuity of optimal allocation function.** We refer to [MCP14] (see Lemma 11) for this proof. We denote  $Q(\tilde{\mathbf{A}}') \triangleq \lim_{\hat{\mathcal{F}} \rightarrow \mathcal{F}} \max_{\omega \in \Delta_K} \left\{ \sum_{a=1}^K \omega_a d(\mu, \lambda) - l^T \tilde{\mathbf{A}}'' \omega \mid \tilde{\mathbf{A}}'' \omega \leq 0, \omega_a \geq 0 \forall i \in [K] \right\} = \omega(\mu)$  where  $\tilde{\mathbf{A}}' \in \mathbb{R}^{K \times K}$  is the rank-1 update of  $\tilde{\mathbf{A}}$  which is a sub-matrix of  $\tilde{\mathbf{A}}$  with  $K$  number of active constraints. We define limiting set as

$$Q^*(\tilde{\mathbf{A}}'') = \left\{ \omega : \lim_{\hat{\mathcal{F}} \rightarrow \mathcal{F}} \sum_{a=1}^K \omega_a d(\mu, \lambda) = Q(\tilde{\mathbf{A}}'') \mid \tilde{\mathbf{A}}'' \omega \leq 0, \omega_a \geq 0 \forall i \in [K] \right\} = \omega^*(\mu)$$

As a direct consequence of Lemma 11 we get the following results

1. The function  $\omega^*(\mu)$  is continuous in  $(\mathbb{R}^{K \times K}) \times \mathbb{R}^K$
2.  $\omega^*(\mu)$  is upper-hemicontinuous on  $(\mathbb{R}^{K \times K}) \times \mathbb{R}^K$

Leveraging these four sufficient statements ensure that there exist unique solution for the optimization problem  $\max_{\pi \in \hat{\mathcal{F}}} \mu^T \pi, \forall \mu \in \mathcal{D}$  i.e the image set of the set-valued function  $\omega^*(\cdot)$  is singleton.  $\square$



## 913 F Lagrangian Lower Bound for Gaussians: Proof of Theorem 4

914 **Theorem 4.** Let  $\{P_a\}_{a \in [K]}$  be Gaussian distributions with equal variance  $\sigma^2 > 0$ , and  $\text{Diag}(1/\omega_a)$   
 915 be a  $K$ -dimensional diagonal matrix with  $a$ -th diagonal entry  $1/\omega_a$ . Then, we get

$$T_{\hat{\mathcal{F}}}^{-1}(\mu) = \max_{\omega \in \hat{\mathcal{F}}} \min_{l \in \mathcal{L}} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \left\{ \frac{1}{2\sigma^2} \frac{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2}{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2} - l^T \tilde{\mathbf{A}}\omega \right\}.$$

916 *Proof.* We start the proof by the definition of  $\mathcal{D}(\omega, \mu, \hat{\mathcal{F}})$  as per Equation (9)

$$\begin{aligned} \mathcal{D}(\omega, \mu, \hat{\mathcal{F}}) &= \min_{l \in \mathcal{L}} \min_{\lambda \in \Lambda_{\hat{\mathcal{F}}}(\mu)} \left\{ \sum_{a=1}^k \omega_a d(\mu_a, \lambda_a) - l^T \tilde{\mathbf{A}}\omega \right\} \\ &= \min_{l \in \mathcal{L}} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \min_{\lambda: \lambda^T(\pi_{\hat{\mathcal{F}}}^* - \pi')=0} \left\{ \sum_{a=1}^k \omega_a d(\mu_a, \lambda_a) - l^T \tilde{\mathbf{A}}\omega \right\} \rightsquigarrow \text{via Proposition 1} \end{aligned} \quad (12)$$

917 The Lagrangian formulation of  $\mathcal{D}(\omega, \mu, \hat{\mathcal{F}})$  is written as

$$\mathcal{L}(\omega, \mu, \hat{\mathcal{F}}) = \min_{l \in \mathcal{L}} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \left\{ \sum_{a=1}^K \omega_a d(\mu_a, \lambda_a) - l^T \tilde{\mathbf{A}}\omega - \gamma \sum_{a=1}^K \lambda_a v_a \right\}$$

918 where  $v_a \triangleq (\pi_{\hat{\mathcal{F}}}^* - \pi')_a$ .

919 We assume both the instances  $\mu$  and  $\lambda$  follow Gaussian distribution with same variance  $\sigma^2$ .

920 Then, we can rewrite the Lagrangian putting the value of the KL as —

$$\mathcal{L}(\omega, \mu, \hat{\mathcal{F}}) = \min_{\gamma \in \mathbb{R}_+} \min_{l \in \mathcal{L}} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \left\{ \sum_{a=1}^K \omega_a \frac{(\mu_a - \lambda_a)^2}{2\sigma^2} - l^T \tilde{\mathbf{A}}\omega - \gamma \sum_{a=1}^K \lambda_a v_a \right\} \quad (13)$$

921 Differentiating the Lagrangian w.r.t  $\lambda_a$  and equating it to 0, we get

$$\begin{aligned} \nabla_{\lambda_a} \mathcal{L}(\lambda, \omega, \mu) &= 0 \\ \text{or, } -\frac{\omega_a(\mu_a - \lambda_a)}{\sigma^2} - \gamma v_a &= 0 \\ \text{or, } \lambda_a &= \mu_a + \frac{\gamma v_a \sigma^2}{\omega_a} \end{aligned}$$

922 Then putting back the value of  $\lambda_a$  in Equation 13 we get

$$\mathcal{L}(\omega, \mu, \hat{\mathcal{F}}) = \min_{\gamma \in \mathbb{R}_+} \min_{l \in \mathcal{L}} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \left\{ \sum_{a=1}^K \gamma^2 \frac{v_a^2 \sigma^2}{2\omega_a} - l^T \tilde{\mathbf{A}}\omega - \gamma \sum_{a=1}^K \mu_a v_a \right\} \quad (14)$$

923 Again differentiating the Lagrangian w.r.t  $\gamma$  and equating it to 0, we get

$$\begin{aligned} \nabla_{\gamma} \mathcal{L}(\omega, \mu, \hat{\mathcal{F}}) = 0 &\implies -\sum_{a=1}^K \mu_a v_a - \gamma \sum_{a=1}^K \frac{\sigma^2}{\omega_a} v_a = 0 \\ &\implies \gamma = -\frac{\mu^T v}{\sum_{a=1}^K \frac{\sigma^2}{\omega_a} v_a^2} \end{aligned}$$

924 Putting the value of  $\gamma$  in Equation 14, we get —

$$\mathcal{L}(\omega, \mu, \hat{\mathcal{F}}) = \min_{l \in \mathcal{L}} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \left\{ \frac{(\mu^T v)^2}{2\sigma^2 \sum_{a=1}^K \frac{v_a^2}{\omega_a}} - l^T \tilde{\mathbf{A}}\omega \right\}$$

$$= \min_{l \in \mathcal{L}} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \left\{ \frac{1}{2\sigma^2} \frac{[\mu^T(\pi_{\hat{\mathcal{F}}}^* - \pi')]^2}{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2} - l^T \tilde{\mathbf{A}} \omega \right\}$$

Therefore inverse characteristic time for Lagrangian relaxation with unknown constraints satisfies,

$$T_{\hat{\mathcal{F}}}^{-1}(\mu) = \max_{\omega \in \hat{\mathcal{F}}} \min_{l \in \mathcal{L}} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \left\{ \frac{1}{2\sigma^2} \frac{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2}{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2} - l^T \tilde{\mathbf{A}} \omega \right\}$$

□

## F.1 Bounds on Sample complexity: Proof of Corollary 1 Part (a)

**Corollary 1. Part (a)** Let  $d_{\pi}^2 \triangleq \frac{\|\pi_{\hat{\mathcal{F}}}^* - \pi\|_{\mu\mu^T}^2}{\|\pi_{\hat{\mathcal{F}}}^* - \pi\|_2^2}$  be the norm of the projection of  $\mu$  on the policy gap  $(\pi_{\hat{\mathcal{F}}}^* - \pi)$ . Then, the characteristic time  $T_{\hat{\mathcal{F}}}(\mu)$  satisfies  $\frac{2\sigma^2}{C_{\text{known}} + 2C_{\text{unknown}}} \leq T_{\hat{\mathcal{F}}}(\mu) \leq \frac{2\sigma^2 K}{C_{\text{known}}}$ , where  $C_{\text{unknown}} \triangleq \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} d_{\pi'}^2$  and  $C_{\text{known}} = \min_{\pi'' \in \nu_{\mathcal{F}}(\pi_{\hat{\mathcal{F}}}^*)} d_{\pi''}^2$ .

*Proof.* Here, we derive explicit expression for gaussian characterisation of the lower and upper bound on the characteristic time. We start the proof with the difference in sample complexity between unknown and known constraint setting

$$\begin{aligned} \mathcal{D}(\omega, \mu, \hat{\mathcal{F}}) - \mathcal{D}(\omega, \mu, \mathcal{F}) &= \min_{l \in \mathcal{L}} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \left\{ \frac{1}{2\sigma^2} \frac{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2}{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2} - l^T \tilde{\mathbf{A}} \omega \right\} \\ &\quad - \min_{\pi'' \in \nu_{\mathcal{F}}(\pi_{\hat{\mathcal{F}}}^*)} \frac{1}{2\sigma^2} \frac{\|\pi_{\hat{\mathcal{F}}}^* - \pi''\|_{\mu\mu^T}^2}{\|\pi_{\hat{\mathcal{F}}}^* - \pi''\|_{\text{Diag}(1/\omega_a)}^2} \end{aligned} \quad (15)$$

Let us remind, due to pessimistic choice of  $\tilde{\mathbf{A}}$ ,  $\mathcal{F} \subseteq \hat{\mathcal{F}}$ .  $\pi'$  is a neighbour of  $\pi_{\hat{\mathcal{F}}}^*$  if it is an extreme point in the polytope  $\hat{\mathcal{F}}$  and shares (K-1) active constraints with  $\pi_{\hat{\mathcal{F}}}^*$ . Then  $\pi_{\hat{\mathcal{F}}}^*$  and  $\pi''$  lies in the interior of  $\hat{\mathcal{F}}$  i.e, they can be expressed as a convex combination of  $\pi_{\hat{\mathcal{F}}}^*$  and  $\pi'$ . Let,  $\exists 0 \leq t_1 \leq 1 : \pi_{\hat{\mathcal{F}}}^* = t_1 \pi_{\hat{\mathcal{F}}}^* + (1 - t_1) \pi'$  and  $\exists 0 \leq t_2 \leq 1 : \pi'' = t_2 \pi_{\hat{\mathcal{F}}}^* + (1 - t_2) \pi'$ . Then,  $(\pi_{\hat{\mathcal{F}}}^* - \pi') = t_1(\pi_{\hat{\mathcal{F}}}^* - \pi')$  and  $(\pi'' - \pi') = t_2(\pi_{\hat{\mathcal{F}}}^* - \pi')$ . Then,

$$\begin{aligned} \|\pi_{\hat{\mathcal{F}}}^* - \pi''\|_{\text{Diag}(1/\omega_a)}^2 &= \|(\pi_{\hat{\mathcal{F}}}^* - \pi') - (\pi'' - \pi')\|_{\text{Diag}(1/\omega_a)}^2 \\ &\leq 2 \left\{ \|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2 + \|\pi'' - \pi'\|_{\text{Diag}(1/\omega_a)}^2 \right\} \\ &= 2 \{t_1^2 + t_2^2\} \|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2 \end{aligned}$$

Putting the above inequality in Equation 15, we get —

$$\begin{aligned} &\mathcal{D}(\omega, \mu, \hat{\mathcal{F}}) - \mathcal{D}(\omega, \mu, \mathcal{F}) \\ &\leq \min_{l \in \mathcal{L}} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \left\{ \frac{1}{2\sigma^2} \frac{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2}{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2} - l^T \tilde{\mathbf{A}} \omega \right\} \\ &\leq \frac{1}{2\sigma^2} \min_{l \in \mathcal{L}} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \min_{\pi'' \in \nu_{\mathcal{F}}(\pi_{\hat{\mathcal{F}}}^*)} \left[ \frac{2 \{t_1^2 + t_2^2\} \|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2 - \|\pi_{\hat{\mathcal{F}}}^* - \pi''\|_{\mu\mu^T}^2}{2 \{t_1^2 + t_2^2\} \|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2} \right. \\ &\quad \left. - l^T \tilde{\mathbf{A}} \omega \right] \end{aligned} \quad (16)$$

Now, we have already proven in Lemma 5 that  $(-l^T \tilde{\mathbf{A}} \omega) \leq \mathcal{D}(\omega, \mu, \hat{\mathcal{F}}) \psi$  where,  $\psi \triangleq \frac{\|(\tilde{\mathbf{A}} - \mathbf{A})\omega\|_{\infty} + \max_{i \in [1, N]} \Gamma_i}{\gamma}$ . Also,

$$\|\pi_{\hat{\mathcal{F}}}^* - \pi''\|_{\mu\mu^T}^2 = \|(\pi_{\hat{\mathcal{F}}}^* - \pi') - (\pi'' - \pi')\|_{\mu\mu^T}^2$$

$$\begin{aligned}
&\geq 2 \left\{ \|\pi_{\mathcal{F}}^* - \pi'\|_{\mu\mu^T}^2 - \|(\pi'' - \pi')\|_{\mu\mu^T}^2 \right\} \\
&= 2 \left\{ t_1^2 \|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2 - t_2^2 \|(\pi_{\hat{\mathcal{F}}}^* - \pi')\|_{\mu\mu^T}^2 \right\} \\
&= 2(t_1^2 - t_2^2) \|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2
\end{aligned}$$

942 Therefore Equation 16 gives us

$$\begin{aligned}
&(1 + \psi) \mathcal{D}(\omega, \mu, \hat{\mathcal{F}}) \\
&\leq \mathcal{D}(\omega, \mu, \mathcal{F}) + \frac{1}{2\sigma^2} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \min_{\pi'' \in \nu_{\mathcal{F}}(\pi_{\mathcal{F}}^*)} \frac{2\{t_1^2 + t_2^2\} - 2(t_1^2 - t_2^2)}{2\{t_1^2 + t_2^2\}} \frac{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2}{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2} \\
&= \mathcal{D}(\omega, \mu, \mathcal{F}) + \frac{1}{2\sigma^2} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \min_{\pi'' \in \nu_{\mathcal{F}}(\pi_{\mathcal{F}}^*)} \frac{2t_2^2}{\{t_1^2 + t_2^2\}} \frac{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2}{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2} \quad (17)
\end{aligned}$$

943 Now to get the lower bound on characteristic time in unknown constraint setting, we take maximum  
944 over  $\omega$  in Equation 17,

$$\begin{aligned}
T_{\hat{\mathcal{F}}}^{-1}(\mu) &\leq \frac{1}{(1 + \psi)} \left\{ T_{\mathcal{F}}^{-1}(\mu) + \frac{1}{2\sigma^2} \max_{\omega} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \min_{\pi'' \in \nu_{\mathcal{F}}(\pi_{\mathcal{F}}^*)} \frac{2t_2^2}{\{t_1^2 + t_2^2\}} \frac{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2}{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2} \right\} \\
&< T_{\mathcal{F}}^{-1}(\mu) + \frac{1}{\sigma^2} \max_{\omega} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \frac{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2}{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2} \quad (18)
\end{aligned}$$

945 To get the lower bound on characteristic time we minimize  $\|\pi_{\mathcal{F}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2$  with the simplex  
946 constraint  $\sum_{a=1}^K \omega_a = 1$  for each and every neighbour in  $\nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)$  and  $\nu_{\mathcal{F}}(\pi_{\mathcal{F}}^*)$ . the lagrangian  
947 formulation gives us the solution,  $\omega_a = \frac{|\pi_{\hat{\mathcal{F}}}^* - \pi'|_a}{\sum_{a=1}^K |\pi_{\hat{\mathcal{F}}}^* - \pi'|_a}$ . Therefore,  $\|\pi_{\mathcal{F}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2 \geq$   
948  $\|\pi_{\mathcal{F}}^* - \pi'\|_1^2 \geq \|\pi_{\mathcal{F}}^* - \pi'\|_2^2$ , since  $\forall a \in [1, K] : \|\pi_{\mathcal{F}}^* - \pi'\|_a \leq 1$ . Then, putting the value of  $\omega_a$  in  
949 Equation 18,

$$T_{\hat{\mathcal{F}}}^{-1}(\mu) \leq T_{\mathcal{F}}^{-1}(\mu) + \frac{1}{\sigma^2} \max_{\omega} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \frac{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2}{\|\pi_{\mathcal{F}}^* - \pi'\|_2^2}$$

950 For ease of comparison, we follow the same notation used in [CBJD23] and define for any  $\pi \in \Delta_{K-1}$ ,  
951  $d_{\pi}^2 = \frac{\|\pi_{\mathcal{F}}^* - \pi\|_{\mu\mu^T}^2}{\|\pi_{\mathcal{F}}^* - \pi\|_2^2}$ , which is the squared distance between  $\mu$  and the hyper-plane  $(\pi_{\mathcal{F}}^* - \pi) = \mathbf{0}$ .  
952 Therefore, using [CBJD23, Corollary 1], we get

$$T_{\hat{\mathcal{F}}}^{-1}(\mu) \leq \frac{1}{2\sigma^2} \left( \min_{\pi'' \in \nu_{\mathcal{F}}(\pi_{\mathcal{F}}^*)} d_{\pi''}^2 + 2 \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} d_{\pi'}^2 \right)$$

953 Therefore lower bound on the characteristic time is given by,

$$T_{\hat{\mathcal{F}}}(\mu) \geq \left( \frac{2\sigma^2}{\min_{\pi'' \in \nu_{\mathcal{F}}(\pi_{\mathcal{F}}^*)} d_{\pi''}^2 + 2 \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} d_{\pi'}^2} \right)$$

954 Now, let us find an upper bound on the characteristic time. We have already shown that,  
955  $\|\pi_{\mathcal{F}}^* - \pi''\|_{\mu\mu^T}^2 \geq 2(t_1^2 - t_2^2) \|\pi_{\hat{\mathcal{F}}}^* - \pi''\|_{\mu\mu^T}^2$ , which also implies  $\|\pi_{\mathcal{F}}^* - \pi''\|_{\text{Diag}(1/\omega_a)}^2 \geq$   
956  $2(t_1^2 - t_2^2) \|\pi_{\hat{\mathcal{F}}}^* - \pi''\|_{\text{Diag}(1/\omega_a)}^2$ .  
957

958 Using this result in Equation 15,

$$\begin{aligned}
&\mathcal{D}(\omega, \mu, \hat{\mathcal{F}}) - \mathcal{D}(\omega, \mu, \mathcal{F}) \\
&\geq \frac{1}{2\sigma^2} \left\{ \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \frac{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2}{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2} - \min_{\pi'' \in \nu_{\mathcal{F}}(\pi_{\mathcal{F}}^*)} \frac{\|\pi_{\mathcal{F}}^* - \pi''\|_{\mu\mu^T}^2}{2(t_1^2 - t_2^2) \|\pi_{\hat{\mathcal{F}}}^* - \pi''\|_{\text{Diag}(1/\omega_a)}^2} \right\}
\end{aligned}$$

$$= \frac{1}{2\sigma^2} \left\{ \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \min_{\pi'' \in \nu_{\mathcal{F}}(\pi_{\mathcal{F}}^*)} \frac{2(t_1^2 - t_2^2) \|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2 - \|\pi_{\mathcal{F}}^* - \pi''\|_{\mu\mu^T}^2}{2(t_1^2 - t_2^2) \|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2} \right\} \quad (19)$$

959 We have already shown,  $\|\pi_{\mathcal{F}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2 \leq 2(t_1^2 + t_2^2) \|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2$ . Therefore  
 960  $\|\pi_{\mathcal{F}}^* - \pi''\|_{\mu\mu^T}^2 \leq 2(t_1^2 + t_2^2) \|\pi_{\hat{\mathcal{F}}}^* - \pi''\|_{\mu\mu^T}^2$ . Consequently Equation 19 gives,

$$\begin{aligned} & \mathcal{D}(\omega, \mu, \hat{\mathcal{F}}) - \mathcal{D}(\omega, \mu, \mathcal{F}) \\ & \geq \frac{1}{2\sigma^2} \left\{ \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \min_{\pi'' \in \nu_{\mathcal{F}}(\pi_{\mathcal{F}}^*)} \frac{2(t_1^2 - t_2^2) \|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2 - 2(t_1^2 + t_2^2) \|\pi_{\hat{\mathcal{F}}}^* - \pi''\|_{\mu\mu^T}^2}{2(t_1^2 - t_2^2) \|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2} \right\} \\ & = \frac{1}{2\sigma^2} \left\{ \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} \min_{\pi'' \in \nu_{\mathcal{F}}(\pi_{\mathcal{F}}^*)} \left( \frac{2t_2^2}{t_2^2 - t_1^2} \right) \frac{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\mu\mu^T}^2}{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_{\text{Diag}(1/\omega_a)}^2} \right\} \geq 0 \\ \implies & T_{\hat{\mathcal{F}}}^{-1}(\mu) \geq T_{\mathcal{F}}^{-1}(\mu) \geq \frac{1}{2\sigma^2 K} \left\{ \min_{\pi'' \in \nu_{\mathcal{F}}(\pi_{\mathcal{F}}^*)} d_{\pi''}^2 \right\} \end{aligned}$$

961 Therefore, upper bound on the characteristic time is given by,

$$T_{\hat{\mathcal{F}}}(\mu) \leq \frac{2\sigma^2 K}{\min_{\pi'' \in \nu_{\mathcal{F}}(\pi_{\mathcal{F}}^*)} d_{\pi''}^2}$$

962 Let,  $C_{\text{unknown}} \triangleq \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} d_{\pi'}^2$ , and  $C_{\text{known}} \triangleq \min_{\pi'' \in \nu_{\mathcal{F}}(\pi_{\mathcal{F}}^*)} d_{\pi''}^2$ , and the result follows.  $\square$

## 963 F.2 Impact of unknown linear constraints: Proof of Corollary 1 Part (b)

964 **Corollary 1. Part (b)** Let  $d_{\pi}^2 \triangleq \frac{\|\pi_{\mathcal{F}}^* - \pi\|_{\mu\mu^T}^2}{\|\pi_{\mathcal{F}}^* - \pi\|_2^2}$  be the norm of the projection of  $\mu$  on the policy gap  
 965  $(\pi_{\mathcal{F}}^* - \pi)$ . Then, the characteristic time  $T_{\hat{\mathcal{F}}}(\mu)$  satisfies  $T_{\hat{\mathcal{F}}}(\mu) \geq \frac{H}{\kappa_{\text{known}}^2 + 2\kappa_{\text{unknown}}^2}$ .  $H$  is the sum  
 966 of squares of gaps.  $\kappa_{\text{known}}$  and  $\kappa_{\text{unknown}}$  are condition numbers of  $\mathbf{A}$  and  $\tilde{\mathbf{A}}$ .

967 *Proof.* We dictate this proof in total four steps. The first step mainly deals with a sub-matrix  $\tilde{\mathbf{A}}'$  of  
 968  $\tilde{\mathbf{A}}$  with  $K$  number of active constraints where it will be shown that we can get any neighbouring  
 969 policy of  $\pi_{\hat{\mathcal{F}}}^*$  just by a rank-1 update of  $\tilde{\mathbf{A}}'$  which means we just one of the constraints inactive for  
 970  $\pi_{\hat{\mathcal{F}}}^*$ . Using this result we will find an upper bound in the next step on the deviation between  $\pi_{\hat{\mathcal{F}}}^*$  and  
 971 any of its neighbouring policy sharing  $K - 1$  number of active constraints. Once we find the upper  
 972 bound on this deviation, in the third step we give a new expression for the characteristic time and we  
 973 conclude the fourth step by reducing a new lower bound characterised the condition number of  $\tilde{\mathbf{A}}'$   
 974 and  $\tilde{\mathbf{A}}$ , where  $\tilde{\mathbf{A}}$  is the sub-matrix of the actual constraint matrix  $\mathbf{A}$  with at least  $K$  number of active  
 975 constraints.

976 **Step 1 : Optimal policy to neighbouring policy via rank-1 update.** We have the pessimistic  
 977 estimate of  $\mathbf{A}$  as  $\tilde{\mathbf{A}}$ , which gives us the perturbed feasible space  $\hat{\mathcal{F}}$ . Let,  $\tilde{\mathbf{A}}'$  be a sub-matrix of  $\tilde{\mathbf{A}}$  such  
 978 that it consists of  $K$  linearly independent rows of  $\tilde{\mathbf{A}}$  active at  $\pi_{\hat{\mathcal{F}}}^*$ . We can then say  $\pi_{\hat{\mathcal{F}}}^* \in \text{Null}(\tilde{\mathbf{A}}')$ ,  
 979 where,  $\text{Null}(\tilde{\mathbf{A}}')$  is the null space of  $\tilde{\mathbf{A}}'$ . Now for some neighbour of  $\pi_{\hat{\mathcal{F}}}^*$ ,  $\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)$  belongs to  
 980 the null space of some  $\tilde{\mathbf{A}}''$ , where this  $\tilde{\mathbf{A}}''$  can be expressed as a rank-1 update of  $\tilde{\mathbf{A}}'$ . Specifically,  
 981  $\tilde{\mathbf{A}}'' = \tilde{\mathbf{A}}' + e_r(a''_r - a'_r)^T$ . Here,  $a'_r$  is the column corresponding to the  $r$ -th constraint of  $\tilde{\mathbf{A}}'$ . We  
 982 just want to replace this column with a new column  $a''_r$ , so we set  $e_r$  as a vector which has 1 at the  
 983  $r$ -th position and 0 everywhere else.

984 **Step 2: Bounding distance between neighbor and optimal policy.** From Equation 18, we have —

$$T_{\hat{\mathcal{F}}}^{-1}(\mu) < T_{\mathcal{F}}^{-1}(\mu) + \frac{1}{2\sigma^2} \min_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_{\hat{\mathcal{F}}}^*)} 2 \frac{(\mu^T(\pi_{\hat{\mathcal{F}}}^* - \pi'))^2}{\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|_2^2} \quad (20)$$

985 So, it becomes evident that we need to get an upper bound on the quantity  $\|\pi_{\hat{\mathcal{F}}}^* - \pi'\|$ . Let us start  
 986 with,

$$\tilde{\mathbf{A}}'(\pi' - \pi_{\hat{\mathcal{F}}}^*) = \tilde{\mathbf{A}}'\pi' \text{ since, } \pi_{\hat{\mathcal{F}}}^* \in \text{Null}(\tilde{\mathbf{A}}')$$

$$\begin{aligned}
&= \left\{ \tilde{\mathbf{A}}'' + e_r(a'_r - a''_r)^T \right\} \boldsymbol{\pi}' \\
&= e_r(a'_r - a''_r)^T \boldsymbol{\pi}' \\
&= e_r a'^T_r \boldsymbol{\pi}', \text{ since } a''_r \in \text{column space of } \tilde{\mathbf{A}}'' \\
\Rightarrow (\boldsymbol{\pi}' - \boldsymbol{\pi}_{\tilde{\mathcal{F}}}^*) &= \tilde{\mathbf{A}}'^{-1}(e_r a'^T_r) \boldsymbol{\pi}'
\end{aligned}$$

987 We denote  $\xi \triangleq a'^T_r \boldsymbol{\pi}'$  as it is the slack for the new  $r$ -th row/constraint. We also define  $\sigma_{\min}(\tilde{\mathbf{A}}')$  and  
988  $\sigma_{\max}(\tilde{\mathbf{A}}')$  be respectively the minimum and maximum singular value of  $\tilde{\mathbf{A}}'$ . Also, let  $\kappa_{\text{unknown}} \triangleq$   
989  $\frac{\sigma_{\max}(\tilde{\mathbf{A}}')}{\sigma_{\min}(\tilde{\mathbf{A}}')}$  be the minimum condition number for  $\tilde{\mathbf{A}}'$ . Then by property of singular value of a matrix,  
990 it follows —

$$\begin{aligned}
\frac{1}{\sigma_{\min}(\tilde{\mathbf{A}}')} &= \sigma_{\min}(\tilde{\mathbf{A}}'^{-1}) \leq \min_{v: \|v\|_2=1} \|\tilde{\mathbf{A}}'^{-1}v\| \leq \|\tilde{\mathbf{A}}'^{-1}e_r\| \leq \max_{v: \|v\|_2=1} \|\tilde{\mathbf{A}}'^{-1}v\| \\
&\leq \sigma_{\max}(\tilde{\mathbf{A}}'^{-1}) = \frac{1}{\sigma_{\min}(\tilde{\mathbf{A}}')} \\
\Rightarrow \frac{|\xi|}{\sigma_{\min}(\tilde{\mathbf{A}}')} &\leq \|\boldsymbol{\pi}_{\tilde{\mathcal{F}}}^* - \boldsymbol{\pi}'\|_2 \leq \frac{\xi}{\sigma_{\max}(\tilde{\mathbf{A}}')}
\end{aligned}$$

991 **Step 3 : A new expression for Characteristic time.** Plugging in the above obtained bound on  
992  $\|\boldsymbol{\pi}_{\tilde{\mathcal{F}}}^* - \boldsymbol{\pi}'\|_2$  in the expression of inverse of characteristic time in Theorem 4, we get a new expression  
993 for the characteristic time.

$$\begin{aligned}
\frac{1}{2\sigma^2} \frac{\|\boldsymbol{\pi}_{\tilde{\mathcal{F}}}^* - \boldsymbol{\pi}'\|_{\boldsymbol{\mu}\boldsymbol{\mu}^T}^2}{\|\boldsymbol{\pi}_{\tilde{\mathcal{F}}}^* - \boldsymbol{\pi}'\|_{\text{Diag}(1/\omega_a)}^2} &= \frac{1}{2\sigma^2} \frac{\|\tilde{\mathbf{A}}'^{-1}(e_r a'^T_r) \boldsymbol{\pi}'\|_{\boldsymbol{\mu}\boldsymbol{\mu}^T}^2}{\|\tilde{\mathbf{A}}'^{-1}(e_r a'^T_r) \boldsymbol{\pi}'\|_{\text{Diag}(1/\omega_a)}^2} \\
&= \frac{1}{2\sigma^2} \frac{\|\tilde{\mathbf{A}}'^{-1}e_r \xi\|_{\boldsymbol{\mu}\boldsymbol{\mu}^T}^2}{\|\tilde{\mathbf{A}}'^{-1}e_r \xi\|_{\text{Diag}(1/\omega_a)}^2} \\
&= \frac{1}{2\sigma^2} \frac{(\Delta^T \tilde{\mathbf{A}}'^{-1}e_r)^2}{\|\tilde{\mathbf{A}}'^{-1}e_r\|_{\text{Diag}(1/\omega_a)}^2}
\end{aligned}$$

994 where  $\Delta$  is the vector of sub-optimality gaps of arms, i.e  $\Delta_a \triangleq \boldsymbol{\mu}^* - \boldsymbol{\mu}_a$ . Then, the new expression  
995 for the inverse of characteristic time is as follows,

$$T_{\tilde{\mathcal{F}}}^{-1}(\boldsymbol{\mu}) = \max_{\boldsymbol{\omega} \in \tilde{\mathcal{F}}} \min_{l \in \mathcal{L}} \min_{\boldsymbol{\pi}' \in \nu_{\tilde{\mathcal{F}}}(\boldsymbol{\pi}_{\tilde{\mathcal{F}}}^*)} \left\{ \frac{1}{2\sigma^2} \frac{(\Delta^T \tilde{\mathbf{A}}'^{-1}e_r)^2}{\|\tilde{\mathbf{A}}'^{-1}e_r\|_{\text{Diag}(1/\omega_a)}^2} - l^T \tilde{\mathbf{A}} \boldsymbol{\omega} \right\} \quad (21)$$

996 **Step 4 : New expression for the lower bound.**

997 Combining Equation 18 and the new expression in Equation 21, we get —

$$\begin{aligned}
T_{\tilde{\mathcal{F}}}^{-1}(\boldsymbol{\mu}) &< T_{\mathcal{F}}^{-1}(\boldsymbol{\mu}) + \frac{1}{\sigma^2} \max_{\boldsymbol{\omega}} \min_{\boldsymbol{\pi}' \in \nu_{\tilde{\mathcal{F}}}(\boldsymbol{\pi}_{\tilde{\mathcal{F}}}^*)} \frac{\|\boldsymbol{\pi}_{\mathcal{F}}^* - \boldsymbol{\pi}'\|_{\boldsymbol{\mu}\boldsymbol{\mu}^T}^2}{\|\boldsymbol{\pi}_{\mathcal{F}}^* - \boldsymbol{\pi}'\|_{\text{Diag}(1/\omega_a)}^2} \\
&\leq T_{\mathcal{F}}^{-1}(\boldsymbol{\mu}) + \frac{1}{\sigma^2} \min_{\boldsymbol{\pi}' \in \nu_{\tilde{\mathcal{F}}}(\boldsymbol{\pi}_{\tilde{\mathcal{F}}}^*)} \frac{\|\boldsymbol{\pi}_{\mathcal{F}}^* - \boldsymbol{\pi}'\|_{\boldsymbol{\mu}\boldsymbol{\mu}^T}^2}{\|\boldsymbol{\pi}_{\mathcal{F}}^* - \boldsymbol{\pi}'\|_2^2} \\
&= T_{\mathcal{F}}^{-1}(\boldsymbol{\mu}) + \min_{\boldsymbol{\pi}' \in \nu_{\tilde{\mathcal{F}}}(\boldsymbol{\pi}_{\tilde{\mathcal{F}}}^*)} \frac{1}{2\sigma^2} \frac{(\Delta^T \tilde{\mathbf{A}}'^{-1}e_r)^2}{\|\tilde{\mathbf{A}}'^{-1}e_r\|_2^2} \\
&\leq T_{\mathcal{F}}^{-1}(\boldsymbol{\mu}) + \min_{\boldsymbol{\pi}' \in \nu_{\tilde{\mathcal{F}}}(\boldsymbol{\pi}_{\tilde{\mathcal{F}}}^*)} \frac{\Delta^2}{2\sigma^2} \frac{\sigma_{\max}(\tilde{\mathbf{A}}')}{\sigma_{\min}(\tilde{\mathbf{A}}')}
\end{aligned}$$

998 For ease of comparing, we denote  $H \triangleq \frac{2\sigma^2}{\Delta^2}$ . Therefore, the new expression of characteristic time  
 999 satisfies,

$$T_{\hat{\mathcal{F}}}^{-1}(\boldsymbol{\mu}) \leq \min_{\boldsymbol{\pi}'' \in \nu_{\mathcal{F}}(\boldsymbol{\pi}_{\mathcal{F}}^*)} \frac{\kappa_{\text{known}}^2}{H} + \min_{\boldsymbol{\pi}' \in \nu_{\tilde{\mathcal{F}}}(\boldsymbol{\pi}_{\tilde{\mathcal{F}}}^*)} \frac{\kappa_{\text{unknown}}^2}{H} = \frac{1}{H} (\kappa_{\text{known}}^2 + \kappa_{\text{unknown}}^2)$$

1000 where,  $\kappa_{\text{known}} \triangleq \frac{\sigma_{\max}(\hat{\mathbf{A}})}{\sigma_{\min}(\hat{\mathbf{A}})}$  and  $\kappa_{\text{unknown}} \triangleq \frac{\sigma_{\max}(\tilde{\mathbf{A}}')}{\sigma_{\min}(\tilde{\mathbf{A}}')}$ ,  $\hat{\mathbf{A}}$  and  $\tilde{\mathbf{A}}'$  being the sub-matrix of  $\mathbf{A}$  and  $\tilde{\mathbf{A}}$   
 1001 having Klinearly independent rows.

1002 Consequently expected stopping time is then lower bounded by

$$\mathbb{E}[\tau_{\delta}] \geq \frac{H}{\kappa_{\text{known}}^2 + \kappa_{\text{unknown}}^2} \text{kl}(\delta \| 1 - \delta)$$

1003

□

## 1004 G Sample Complexity upper bounds (Analysis of algorithms)

### 1005 G.1 Stopping Criterion

1006 **Theorem 5.** *The Chernoff stopping rule to ensure  $(1 - \delta)$ -correctness and  $(1 - \delta)$ -feasibility is*

$$\inf_{\lambda \in \Lambda_{\hat{\mu}_t}(\hat{\mu}_t)} \sum_{a=1}^K N_{a,t} d(\hat{\mu}_{a,t}, \lambda_a) > \beta(t, \delta) \wedge \|(\tilde{\mathbf{A}}_t - \mathbf{A})\omega_t\|_\infty > \rho(t, \delta),$$

1007 where  $\beta(t, \delta) \triangleq 3S_0 \log(1 + \log N_{a,t}) + S_0 \mathcal{T} \left( \frac{(K \wedge d) + \log \frac{1}{\delta}}{S_0} \right)$ , and  $\rho(t, \delta)$  is in Lemma 3.

1008 *Proof.* We dictate the proof in 2 steps. In the first step, we prove that the stopping time  $\tau_\delta$  is finite.  
 1009 Then, in next step, we give an explicit expression of the stopping threshold by upper bounding  
 1010 probability of the bad event for stopping time  $\tau_\delta$ .

1011 Let us first go through some notations.

$$1012 \pi_t \triangleq \arg \max_{\pi \in \mathcal{F}} \hat{\mu}_t^T \pi, \text{ where } \hat{\mu}_t \in \arg \min_{\lambda \in \Lambda_{\hat{\mu}_t}(\hat{\mu}_t)} \sum_{a=1}^K N_{a,t} d(\hat{\mu}_{a,t}, \lambda_a) - \mathbf{l}_t^T \tilde{\mathbf{A}}_t N_{a,t}.$$

1013 Algorithm 1 and 2 stops at a finite  $\tau_\delta \in \mathbb{N}$  if the events  $\inf_{\lambda \in \Lambda_{\hat{\mu}_t}(\hat{\mu}_t)} \{\exists t \in \mathbb{N} : \sum_{a=1}^K N_{a,t} d(\hat{\mu}_{a,t}, \lambda_a) > \beta(t, \delta)\}$  and  $\{\exists t \in \mathbb{N} : \|(\tilde{\mathbf{A}}_t - \mathbf{A})N_t\|_\infty \leq t\rho(t, \delta)\}$  jointly occurs.

1015 **Step 1: Finiteness of the stopping time.** A stopping time is finite if the parameters in the system  
 1016 converges to their true values in finite time, in our case the means of arms and the constraint  
 1017 matrix. Let us define  $\mathbf{A} \triangleq \{a \in [K] : \lim_{t \rightarrow \infty} N_{a,t} < \infty\}$  as a sampling rule i.e if an arm  
 1018 belongs to this set  $\mathbf{A}$ , it has been sampled finitely and otherwise the arm has been sampled enough  
 1019 number of times so that the mean of that arm has converged to its true value and the column in  
 1020 the constraint matrix corresponding to that arm as also converged. For arms  $a \in [K]$  and  $a \in$   
 1021  $\mathbf{A}^c$ ,  $\hat{\mu}_{a,t} \rightarrow \tilde{\mu}_a \neq \mu_a$  and  $(\tilde{\mathbf{A}})_{a,t} \rightarrow (\mathbf{A}')_a \neq (\mathbf{A})_a$ . When all parameters are concentrated  
 1022  $\mathbf{A} = \emptyset$ , we say  $\forall a \in [K] : \hat{\mu}_a \rightarrow \mu_a$  and  $\tilde{\mathbf{A}} \rightarrow \mathbf{A}$ . We also define the limit of this empirical  
 1023 sampling rule as  $\omega_\infty = \lim_{t \rightarrow \infty} \frac{N_{a,t}}{t} \forall a \in [K]$ . We then write the stopping condition in a new way  
 1024  $\inf_{\lambda \in \Lambda_{\hat{\mu}_t}(\hat{\mu}_t)} \{\exists t \in \mathbb{N} : \sum_{a=1}^K \frac{N_{a,t}}{t} d(\hat{\mu}_{a,t}, \lambda_a) > \frac{\beta(t, \delta)}{t}\}$  and  $\{\exists t \in \mathbb{N} : \|(\tilde{\mathbf{A}}_t - \mathbf{A}) \frac{N_t}{t}\|_\infty \leq \rho(t, \delta)\}$ .  
 1025 By continuity properties and knowing  $\beta(t, \cdot) \rightarrow \log \log t$  and  $\rho(t, \cdot) \rightarrow 0$  as  $t \rightarrow \infty$ , we claim  
 1026 by taking taking asymptotic limits both sides  $\inf_{\lambda \in \Lambda_{\hat{\mu}_t}(\hat{\mu}_t)} \sum_{a=1}^K \frac{\omega_{\infty,a}}{t} d(\hat{\mu}_{a,t}, \lambda_a) > 0$  and also  
 1027  $\|(\tilde{\mathbf{A}}_t - \mathbf{A})\omega_{\infty,a}\|_\infty < 0$ . We get strict inequality for the both the cases by the virtue of construction  
 1028 of the set  $\mathbf{A}$  such that for arms  $a \in \mathbf{A}$ ,  $\omega_\infty \neq 0$  and the KL-divergence is non-zero as  $\lambda_a \neq \mu$  since  
 1029  $\lambda \in \Lambda_{\hat{\mu}_t}(\hat{\mu}_t)$ . Also for the second strict inequality, since  $\tilde{\mathbf{A}}_t$  is the pessimistic estimate of  $\mathbf{A}$  at time  
 1030  $t$  the condition will hold.

1031 **Step 2: Probability of bad event to Stopping threshold.** Let  $\omega_t$  is the allocation associated to  $N_t$ .  
 1032 Then we define the bad event as

$$\begin{aligned} U_t &\triangleq \{\pi_{\tau_\delta} \neq \pi_{\mathcal{F}}^*\} \\ &= \bigcup_{\pi \neq \pi_{\mathcal{F}}^*} \left\{ \exists t \in \mathbb{N} : \pi_{t+1} = \pi \right. \\ &\quad \left. \wedge \left\{ \inf_{\lambda \in \Lambda_{\hat{\mu}_t}(\hat{\mu}_t)} \sum_{a=1}^K N_{a,t} d(\hat{\mu}_{a,t}, \lambda_a) > \beta(t, \delta) \vee \mathbf{A}\omega_t > 0 \right\} \right\} \\ (a) \quad &\subseteq \bigcup_{\pi \neq \pi_{\mathcal{F}}^*} \left\{ \exists t \in \mathbb{N} : \pi_{t+1} = \pi \right. \\ &\quad \left. \wedge \left\{ \inf_{\lambda \in \Lambda_{\hat{\mu}_t}(\hat{\mu}_t)} \sum_{a=1}^K N_{a,t} d(\hat{\mu}_{a,t}, \lambda_a) > \beta(t, \delta) \vee \|(\tilde{\mathbf{A}}_t - \mathbf{A})\omega_t\|_\infty > \rho(t, \delta) \right\} \right\} \end{aligned}$$

1033 The argument (a) holds because

$$\mathbb{P}\{0 > -\mathbf{A}\omega_t\} = \mathbb{P}\{\tilde{\mathbf{A}}_t\omega_t > (\tilde{\mathbf{A}}_t - \mathbf{A})\omega_t\} \leq \mathbb{P}\{\|\tilde{\mathbf{A}}_t\omega_t\|_1 > \|(\tilde{\mathbf{A}}_t - \mathbf{A})\omega_t\|_1 \geq \|(\tilde{\mathbf{A}}_t - \mathbf{A})\omega_t\|_\infty\}$$

$$\begin{aligned} &\leq \mathbb{P}\{\|\tilde{\mathbf{A}}_t \boldsymbol{\omega}_t\|_1 > \rho(t, \delta)\} \mathbb{P}\{\|(\tilde{\mathbf{A}}_t - \mathbf{A}) \boldsymbol{\omega}_t\|_\infty > \rho(t, \delta)\} \\ &= \mathbb{P}\{\|(\tilde{\mathbf{A}}_t - \mathbf{A}) \boldsymbol{\omega}_t\|_\infty > \rho(t, \delta)\} \end{aligned}$$

1034 since the event  $\{\|\tilde{\mathbf{A}}_t \boldsymbol{\omega}_t\|_1 > \rho(t, \delta)\}$  is a sure event.

1035 Therefore probability of this bad event

$$\begin{aligned} \mathbb{P}(U_t) &\leq \bigcup_{\boldsymbol{\pi} \neq \boldsymbol{\pi}_{\mathcal{F}}^*} \mathbb{P} \left\{ \exists t \in \mathbb{N} : \boldsymbol{\pi}_{t+1} = \boldsymbol{\pi} \wedge \inf_{\boldsymbol{\lambda} \in \Lambda_{\hat{\mathcal{F}}_t}(\hat{\boldsymbol{\mu}}_t)} \sum_{a=1}^K N_{a,t} d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\lambda}_a) \right\} > \beta(t, \delta) \Big\} \\ &\quad + \mathbb{P} \left\{ \exists t \in \mathbb{N} : \boldsymbol{\pi}_{t+1} = \boldsymbol{\pi} \wedge \|(\tilde{\mathbf{A}}_t - \mathbf{A}) \boldsymbol{\omega}_t\|_\infty > \rho(t, \delta) \right\} \\ &= \sum_{\boldsymbol{\pi} \neq \boldsymbol{\pi}_{\mathcal{F}}^*} \mathbb{P} \left\{ \exists t \in \mathbb{N} : \inf_{\boldsymbol{\lambda} \in \Lambda_{\hat{\mathcal{F}}_t}(\hat{\boldsymbol{\mu}}_t)} \sum_{a=1}^K N_{a,t} d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\lambda}_a) \right\} > \beta(t, \delta) \Big\} \\ &\quad + \sum_{\boldsymbol{\pi} \neq \boldsymbol{\pi}_{\mathcal{F}}^*} \mathbb{P} \left\{ \exists t \in \mathbb{N} : \|(\tilde{\mathbf{A}}_t - \mathbf{A}) \boldsymbol{\omega}_t\|_\infty > \rho(t, \delta) \right\} \end{aligned} \quad (22)$$

1036 The second cumulative probability can be bound using Lemma 3, i.e

$$\sum_{\boldsymbol{\pi} \neq \boldsymbol{\pi}_{\mathcal{F}}^*} \mathbb{P} \left\{ \exists t \in \mathbb{N} : \|(\tilde{\mathbf{A}}_t - \mathbf{A}) \boldsymbol{\omega}_t\|_\infty > \rho(t, \delta) \right\} \leq \frac{1}{t}$$

1037 for the choice of  $\rho(t, \delta)$  given in Lemma 3. We work with the first term in R.H.S of Equation (22).

$$\begin{aligned} &\sum_{\boldsymbol{\pi} \neq \boldsymbol{\pi}_{\mathcal{F}}^*} \mathbb{P} \left\{ \exists t \in \mathbb{N} : \inf_{\boldsymbol{\lambda} \in \Lambda_{\hat{\mathcal{F}}_t}(\hat{\boldsymbol{\mu}}_t)} \sum_{a=1}^K N_{a,t} d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\lambda}_a) \right\} > \beta(t, \delta) \Big\} \\ &= \sum_{\boldsymbol{\pi} \neq \boldsymbol{\pi}_{\mathcal{F}}^*} \mathbb{P} \left\{ \exists t \in \mathbb{N} : \inf_{\boldsymbol{\lambda} \in \Lambda_{\hat{\mathcal{F}}_t}(\hat{\boldsymbol{\mu}}_t)} \inf_{\boldsymbol{\lambda}' \in \Lambda_{\mathcal{F}}(\hat{\boldsymbol{\mu}}_t)} \sum_{a=1}^K N_{a,t} (d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\lambda}_a) - d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\lambda}'_a)) \right. \\ &\quad \left. + \inf_{\boldsymbol{\lambda}' \in \Lambda_{\mathcal{F}}(\hat{\boldsymbol{\mu}}_t)} \sum_{a=1}^K N_{a,t} d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\lambda}'_a) > \beta(t, \delta) \right\} \\ &\leq \sum_{\boldsymbol{\pi} \neq \boldsymbol{\pi}_{\mathcal{F}}^*} \mathbb{P} \left\{ \exists t \in \mathbb{N} : \inf_{\boldsymbol{\lambda}' \in \Lambda_{\mathcal{F}}(\hat{\boldsymbol{\mu}}_t)} \sum_{a=1}^K N_{a,t} d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\lambda}'_a) \leq \beta(t, \delta) \right\} \end{aligned}$$

1038 The last inequality holds because

$$\inf_{\boldsymbol{\lambda} \in \Lambda_{\hat{\mathcal{F}}_t}(\hat{\boldsymbol{\mu}}_t)} \inf_{\boldsymbol{\lambda}' \in \Lambda_{\mathcal{F}}(\hat{\boldsymbol{\mu}}_t)} \sum_{a=1}^K N_{a,t} (d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\lambda}_a) - d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\lambda}'_a)) \leq 0$$

1039 since under  $\hat{\mathcal{F}}_t$  and the bad event we are assuming that the estimated alt-set is still bigger than the  
 1040 actual alt-set. So any  $\boldsymbol{\lambda} \in \Lambda_{\hat{\mathcal{F}}_t}(\hat{\boldsymbol{\mu}}_t)$  will have a bigger distance with the estimate of  $\boldsymbol{\mu}$  than any  
 1041  $\boldsymbol{\lambda}' \in \Lambda_{\mathcal{F}}(\hat{\boldsymbol{\mu}}_t)$ . We define  $I_\pi \triangleq \text{Supp}(\boldsymbol{\pi}_{\mathcal{F}}^*) \Delta \text{Supp}(\boldsymbol{\pi})$  and also  $S_0 \triangleq \max_{\boldsymbol{\pi}} |I_\pi|$ . We note that  
 1042  $0 \leq S_0 \leq K$ .

1043 We get from Lemma 8 in [KK21] with the notation of  $\mathcal{T}(\cdot)$  follows from Lemma 8

$$\begin{aligned} &\sum_{\boldsymbol{\pi} \neq \boldsymbol{\pi}_{\mathcal{F}}^*} \mathbb{P} \left\{ \exists t \in \mathbb{N} : \sum_{a \in I_\pi} N_{a,t} d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\mu}_a) \geq \sum_{a \in I_\pi} 3 \log(1 + \log N_{a,t}) + |S_0| \mathcal{T} \left( \frac{\log \frac{|\boldsymbol{\pi}_{\hat{\mathcal{F}}_t}| - 1}{\delta}}{S_0} \right) \right\} \\ &\leq \delta \end{aligned}$$

1044 where  $\delta$  is chosen to be  $\frac{\delta}{|\boldsymbol{\pi}_{\hat{\mathcal{F}}_t}| - 1}$  such that  $\log \frac{|\boldsymbol{\pi}_{\hat{\mathcal{F}}_t}| - 1}{\delta} \leq \log(\frac{2^K}{\delta}) \leq (K \wedge d) + \log \frac{1}{\delta}$



Also  $\sum_{a \in I_\pi} 3 \log(1 + \log N_{a,t}) \leq 3S_0 \log(1 + \log N_{a,t})$ . Therefore the stopping threshold is given by

$$\beta(t, \delta) = 3S_0 \log(1 + \log N_{a,t}) + S_0 \mathcal{T} \left( \frac{(K \wedge d) + \log \frac{1}{\delta}}{S_0} \right)$$

In practice, we use  $S_0 = K$ . □

## G.2 Upper Bound of LATS

**Theorem 8.** *The sample complexity upper bound of LATS to yield a  $(1 - \delta)$ -correct optimal policy is*

$$\lim_{\delta \rightarrow 0} \frac{\mathbb{E}[\tau_\delta]}{\log(1/\delta)} \leq \alpha(1 + \mathfrak{s})T_{\mathcal{F}}(\mu),$$

where  $\mathfrak{s}$  is the shadow price of the true constraint  $\mathbf{A}$ ,  $T_{\mathcal{F}}(\mu)$  is the characteristic time under the true constraint (Equation (5)), and  $\delta \in (0, 1]$ .

*Proof.* We will prove this theorem in 5 steps. In the first step, we define what is considered to be the good event in our unknown constraint setting, then we go on bounding the probability of the complement of this good event in step 2. Once the parameter concentrations are taken care of, we show how we can lower bound the instantaneous complexity of the algorithm in step 3. In step 4, we finally prove the upper bound on stopping time for both good and bad events. We conclude with the asymptotic upper bound on stopping time i.e when  $\delta \rightarrow 0$  and  $\epsilon \rightarrow 0$  in step 5.

**Step 1: Defining the good event.** Given an  $\epsilon > 0$ , we define the good event  $G_T$  as,

$$G_T \triangleq \bigcap_{t=h(T)}^T \{ \|\hat{\mu}_t - \mu\|_\infty \leq \xi(\epsilon) \wedge \|(\tilde{\mathbf{A}}_t^i - \mathbf{A}^i)\omega\|_\infty \leq \phi(\epsilon) \forall \omega \in \hat{\mathcal{F}} \}$$

where,  $\xi(\epsilon) \leq \max_{\pi' \in \nu_{\hat{\mathcal{F}}}(\pi_\pi^*)} \frac{1}{5} \mu^T(\pi_\pi^* - \pi')$ , and  $\phi(\epsilon) \triangleq \dots$  for a given  $\epsilon > 0$ . The good event implies that the means and the constraints are well concentrated in an  $\epsilon$ -ball around their true values. Thus, we have to now bound the extra cost of their correctness and the number of samples required to reach the good events.

We also observe that  $\|\mu' - \mu\|_\infty \leq \xi(\epsilon)$  and  $\|(\tilde{\mathbf{A}}_t^i - \mathbf{A}^i)\omega\|_\infty \leq \phi(\epsilon)$  implies that  $\sup_{\omega' \in \omega^*(\mu')} \sup_{\omega \in \omega^*(\mu)} \|\omega' - \omega\| \leq \epsilon$  due to upper hemicontinuity of  $\omega^*(\mu)$  (Theorem 3).

**Step 2: Samples to Achieve the Good Event.** Now, let us bound the probability of complement of the good event,

$$\begin{aligned} \mathbb{P}(G_T^c) &= \sum_{t=h(T)}^T \left( \mathbb{P} \{ \|\hat{\mu}_t - \mu\|_\infty > \xi(\epsilon) \} + \mathbb{P} \{ \|(\tilde{\mathbf{A}}_t^i - \mathbf{A}^i)\omega\|_\infty > \phi(\epsilon) \} \right) \\ &\leq \sum_{t=h(T)}^T \mathbb{P} \{ \|\hat{\mu}_t - \mu\|_\infty > \xi(\epsilon) \} + \sum_{t=h(T)}^T \mathbb{P} \{ \|(\tilde{\mathbf{A}}_t^i - \mathbf{A}^i)\omega\|_\infty > \phi(\epsilon) \} \\ &\leq BT \exp \left( -CT^{\frac{1}{8}} \right) + K \sum_{t=h(T)}^T \frac{1}{t} \end{aligned}$$

The first inequality is due to the union bound. The second inequality is due to the Lemma 6 [GK16], which states that

$$\sum_{t=h(T)}^T \mathbb{P} \{ \|\hat{\mu}_t - \mu\|_\infty > \xi(\epsilon) \} \leq BT \exp \left( -CT^{\frac{1}{8}} \right),$$

and also due to Lemma 3 that proves concentration bound of the constraint matrix over time.

**Step 3: Tracking argument.** Now, we state how concentrating on means and constraints leads to good concentration on the allocations too. Since we use C-tracking, we can leverage the concentration

1072 in allocation by [DKM19b, Lemma 17]. We use this lemma than D-tracking or the tracking argument  
 1073 in [GK16, Lemma 7] because the optimal allocations might not be unique but the set  $\omega^*(\mu)$  is convex  
 1074 (Theorem 3).

1075 Hence, there exists a  $T_\epsilon$  such that under the good event and  $t \geq \max(T_\epsilon, h(T))$ , we have,

$$\begin{aligned} |(\mu - \hat{\mu}_t)^T \pi_{\mathcal{F}}^*| &\leq |(\mu - \hat{\mu}_t)^T \pi_{\hat{\mathcal{F}}}^*| + |(\mu - \hat{\mu}_t)^T (\pi_{\mathcal{F}}^* - \pi_{\hat{\mathcal{F}}}^*)| \leq 4\xi(\epsilon) \\ &\leq \max_{\pi' \in \nu_{\mathcal{F}}(\pi_{\mathcal{F}}^*)'} \mu^T (\pi_{\mathcal{F}}^* - \pi') \end{aligned}$$

1076 We have replaced the perturbed alt-set with the true alt-set because for  $t \geq \max(T_\epsilon, h(T))$ , we can  
 1077 ensure the convergence of  $\hat{\mathcal{F}}$  almost surely i.e  $\hat{\mathcal{F}} \xrightarrow{\text{a.s.}} \mathcal{F}$

1078 **Step 3: Complexity of identification under good event and constraint.** Now, we want to understand  
 1079 how hard it is to hit the stopping rule even under the good event. First, we define

$$C_{\epsilon, \hat{\mathcal{F}}} \triangleq \inf_{\substack{\mu': \|\mu' - \mu\|_\infty \leq \xi(\epsilon) \\ \omega': \|\omega' - \omega\|_\infty \leq 3\epsilon \\ \tilde{\mathbf{A}}': \|(\tilde{\mathbf{A}} - \mathbf{A})\omega\|_\infty \leq \phi(\epsilon)}} \mathcal{D}(\mu', \omega', \hat{\mathcal{F}}) - l^T \tilde{\mathbf{A}} \omega.$$

1080 Now leveraging Lemma 5, we obtain

$$(1 + \psi) \inf_{\substack{\mu': \|\mu' - \mu\|_\infty \leq \xi(\epsilon) \\ \omega': \|\omega' - \omega\|_\infty \leq 3\epsilon \\ \tilde{\mathbf{A}}': \|(\tilde{\mathbf{A}} - \mathbf{A})\omega\|_\infty \leq \phi(\epsilon)}} \mathcal{D}(\mu', \omega', \hat{\mathcal{F}}) \geq \inf_{\substack{\mu': \|\mu' - \mu\|_\infty \leq \xi(\epsilon) \\ \omega': \|\omega' - \omega\|_\infty \leq 3\epsilon \\ \tilde{\mathbf{A}}': \|(\tilde{\mathbf{A}} - \mathbf{A})\omega\|_\infty \leq \phi(\epsilon)}} \mathcal{D}(\mu', \omega', \hat{\mathcal{F}}) - l^T \tilde{\mathbf{A}} \omega,$$

1081 where definition of  $\psi$  follows from Lemma 5. It quantifies how the Lagrangian lower bound relates  
 1082 with the Likelihood Ratio Test-based quantity in the stopping time.

1083 Therefore by the C-tracking argument 9, we can state

$$\mathcal{D}(\hat{\mu}_t, N_t, \hat{\mathcal{F}}) \geq t \inf_{\substack{\mu': \|\mu' - \mu\|_\infty \leq \xi(\epsilon) \\ \omega': \|\omega' - \omega\|_\infty \leq 3\epsilon \\ \tilde{\mathbf{A}}': \|(\tilde{\mathbf{A}} - \mathbf{A})\omega\|_\infty \leq \phi(\epsilon)}} \mathcal{D}(\mu', \omega', \hat{\mathcal{F}}) \geq \frac{t C_{\epsilon, \hat{\mathcal{F}}}}{1 + \psi}. \quad (23)$$

1084 Here, LHS is the quantity that we use to stop and yield a  $(1 - \delta)$ -correct policy.

1085 **Step 4: Bounding the stopping time with good and bad events.** We denote  $\tau_\delta$  as the stopping time.  
 1086 So we can write upper bound on this stopping time for both good and bad events as

$$\min(\tau_\delta, T) \leq \max(\sqrt{T}, \beta(T, \delta)) + \sum_{t=T_\epsilon}^T \mathbb{1}_{\tau_\delta > T}$$

1087 By the correctness of the stopping time, the event  $\tau(\delta) > t$  happens if  $\mathcal{D}(\hat{\mu}_t, N_t, \hat{\mathcal{F}}) \leq \beta(t, \delta)$  for  
 1088 any  $t \leq T$ .

1089 Now using the lower bound on  $\mathcal{D}(\hat{\mu}_t, N_t, \hat{\mathcal{F}})$  (Equation 23), we get

$$\begin{aligned} T_\epsilon + \sum_{t=T_\epsilon}^T \mathbb{1}(\mathcal{D}(\hat{\mu}_t, N_t, \hat{\mathcal{F}}) \leq \beta(t, \delta)) &\leq \max(\sqrt{T}, \beta(T, \delta)) + \sum_{t=T_\epsilon}^T \mathbb{1}\left(t \frac{C_{\epsilon, \hat{\mathcal{F}}}}{1 + \psi} \leq \beta(T, \delta)\right) \\ &\leq \max(\sqrt{T}, \beta(T, \delta)) + \frac{\beta(T, \delta)(1 + \psi)}{C_{\epsilon, \hat{\mathcal{F}}}} \end{aligned}$$

1090 Let us define a  $T_\delta \triangleq \inf\{T \in \mathbb{N} : \max(\sqrt{T}, \beta(T, \delta)) + \frac{\beta(T, \delta)(1 + \psi)}{C_{\epsilon, \hat{\mathcal{F}}}} \leq T\}$ . To find a lower bound on  
 1091  $T_\delta$ , we refer to [GK16]. Specifically, let us define  $\beta(\eta) \triangleq \{\inf T : T - \max(\sqrt{T}, \beta(T, \delta)) \geq \frac{T}{1 + \eta}\}$   
 1092 for some  $\eta > 0$ . Therefore,

$$T_\delta \leq \beta(\eta) + \inf\{T \in \mathbb{N} : T \frac{C_{\epsilon, \hat{\mathcal{F}}}}{(1 + \psi_{\hat{\mathcal{F}}})(1 + \eta)} \geq \beta(T, \delta)\}$$

Thus, finally combining the results, we upper bound the stopping time as

$$\mathbb{E}[\tau_\delta] \leq T_\epsilon + T_\delta + T_{\text{bad}}.$$

Here,  $T_{\text{bad}} = \sum_{t=1}^{\infty} BT \exp(-CT^{1/8}) + K\zeta(1) < \infty$  is the sum of probability of the bad events over time.  $\zeta(\cdot)$  denotes the Euler-Riemann Zeta function.

**Step 5. Deriving the asymptotics.** Now, we leverage the continuity properties of the Lagrangian characteristic time under approximate constraint to show that we converge to traditional hardness measures as  $\epsilon$  and  $\delta$  tends to zero.

First, we observe that for some  $\alpha > 1$

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}[\tau_\delta]}{\log(1/\delta)} \leq \frac{\alpha(1 + \psi_{\hat{\mathcal{F}}})(1 + \eta)}{C_{\epsilon, \hat{\mathcal{F}}}}$$

Now, if also  $\epsilon \rightarrow 0$ , by the Equation 4, we get  $\tilde{\mathbf{A}} \rightarrow \mathbf{A}$ , and thus,  $\hat{\mathcal{F}} \rightarrow F$ .

Thus, by continuity properties in Theorem 3 and Theorem 2, we get that

$$\mathcal{D}(\hat{\mathcal{F}}, \cdot) \rightarrow \mathcal{D}(\mathcal{F}, \cdot) \text{ and } \psi \rightarrow \frac{\max_{i \in [1, N]} \Gamma}{\min_{i \in [1, N]} \Gamma} \triangleq \frac{\Gamma_{\max}}{\Gamma_{\min}} \triangleq \mathfrak{s}.$$

Here,  $\mathfrak{s}$  is the shadow price of the true constraint matrix, and quantifies the change in the constraint values due to one unit change in the policy vector.

Hence, we conclude that

$$\lim_{\delta \rightarrow 0} \frac{\mathbb{E}[\tau_\delta]}{\log(1/\delta)} \leq \alpha T_{\mathcal{F}}(\boldsymbol{\mu})(1 + \mathfrak{s}).$$

□

### G.3 Upper Bound for LAGEX

**Theorem 7.** *The expected sample complexity of LAGEX satisfies  $\lim_{\delta \rightarrow 0} \frac{\mathbb{E}(\tau_\delta)}{\ln(1/\delta)} \leq T_{\mathcal{F}}(\boldsymbol{\mu}) + 2\mathfrak{s}$ .*

*Proof.* We will do this proof in two parts. In part (a) we will assume that the current recommended policy is the correct policy and try to find an upper bound on the sample complexity of LAGEX. In the next part (b) we break that assumption and try to get an upper bound on the number of steps the recommended policy is not the correct policy.

**Part (a) : Current recommended policy is correct.** Proof structure of this part involves several steps. We start with defining the good event where we introduce a new event associated with the concentration event of the constraint set, then proceeding to prove concentration on that good event. Third step starts with the stopping criterion explained in G.1. In step 4 we define LAGEX as an approximate saddle point algorithm. The next step further transforms the stopping criterion with the help of allocation and instance player's regret that play the zero-sum game. We conclude with the asymptotic upper on the sample complexity characterised by the additive effect of the novel quantity shadow price  $\mathfrak{s}$ .

**Step 1: Defining the good event.** We start the proof first by defining the good event as

$$G_t = \{\forall t \leq T, \forall a \in [K] : N_{a,t} d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\mu}_a) \leq g(t) \wedge \|(\tilde{\mathbf{A}} - \mathbf{A})\boldsymbol{\omega}\|_\infty \leq \rho(t, \delta)\}$$

where,  $g(t) = 3 \log t + \log \log t$  and  $\rho(t, \delta)$  is defined in Lemma 2. The choice of  $g(t)$  is motivated from [DKM19a] which originates from the negative branch of the Lambert's W function. This eventually helps us upper bounding the cumulative probability of the bad event.

**Step 2: Concentrating to the good event** We denote  $G_t^c$  as the bad event where any one of the above events does not occur. Cumulative probability of this bad event

$$\sum_{s=1}^T \mathbb{P}(G_t^c) = \sum_{s=1}^T \mathbb{P}\left(\sum_{a=1}^K N_{a,s} d(\hat{\boldsymbol{\mu}}_{a,s}, \boldsymbol{\mu}_a) > g(s)\right)$$

$$+ \sum_{s=1}^T \mathbb{P} \left( \|(\tilde{\mathbf{A}}_s - \mathbf{A})\boldsymbol{\omega}\|_{\infty} > \rho(s, \delta) \forall \boldsymbol{\omega} \in \hat{\mathcal{F}}_T \right)$$

1126 We get the upper bound on  $\sum_{s=1}^T \mathbb{P} \left( \sum_{a=1}^K N_{a,s} d(\hat{\boldsymbol{\mu}}_{a,s}, \boldsymbol{\mu}_a) > g(s) \right) \leq \frac{\exp(2)}{t^3 \log t} (g(t) + g^2(t) \log t) \leq$   
 1127  $\infty$  as a direct consequence of [DKM19a, Lemma 6]. The second cumulative probability is bounded  
 1128 by  $\zeta(1)$  using Lemma 3, which is finite.

1129 In the next step, we work with the stopping criterion where we do not have access to  $\mathcal{F}$  rather a bigger  
 1130 feasible set  $\hat{\mathcal{F}}_t$ .

1131 **Step 3: Working with the stopping criterion.** The stopping criterion implies that

$$\beta(t, \delta) \geq \inf_{\boldsymbol{\lambda} \in \Lambda_{\hat{\mathcal{F}}_t}(\hat{\boldsymbol{\mu}}_t)} \sum_{a=1}^K N_{a,t} d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\lambda}_a),$$

1132 where the exact expression of  $\beta(t, \delta)$  is defined in Theorem 5.

1133 We use the C-tracking lemma (Lemma 7) to express the stopping in terms of allocations

$$\beta(t, \delta) \geq \inf_{\boldsymbol{\lambda} \in \Lambda_{\hat{\mathcal{F}}_t}(\hat{\boldsymbol{\mu}}_t)} \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\lambda}_a) - (1 + \sqrt{t})K \quad (24)$$

1134 L-Lipschitz property of KL divergence gives

$$|d(\boldsymbol{\mu}_a, \boldsymbol{\lambda}_a) - d(\hat{\boldsymbol{\mu}}_{a,s}, \boldsymbol{\lambda}_a)| \leq L \sqrt{2\sigma^2 \frac{g(s)}{N_{a,s}}} \quad (25)$$

1135 Using this result in Equation (24) we get

$$\begin{aligned} \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\lambda}_a) &\geq \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} d(\hat{\boldsymbol{\mu}}_a, \boldsymbol{\lambda}_a) - L \sqrt{2\sigma^2 K t g(t)} \\ &\geq \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} d(\hat{\boldsymbol{\mu}}_{a,s}, \boldsymbol{\lambda}_a) - L \sqrt{2\sigma^2 K t g(t)} \\ &\quad - 2L \sqrt{2\sigma^2 g(t)} (K^2 + 2\sqrt{2Kt}) \\ &\geq \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} d(\hat{\boldsymbol{\mu}}_{a,s}, \boldsymbol{\lambda}_a) - \mathcal{O}(\sqrt{t \log t}) \end{aligned}$$

1136 the penultimate inequality yields from using the Equation (25). Using this result in Equation (24)

$$\beta(t, \delta) \geq \inf_{\boldsymbol{\lambda} \in \Lambda_{\hat{\mathcal{F}}_t}(\hat{\boldsymbol{\mu}}_t)} \sum_{a=1}^K \omega_{a,s} d(\hat{\boldsymbol{\mu}}_{a,t}, \boldsymbol{\lambda}_a) - (1 + \sqrt{t})K - \mathcal{O}(\sqrt{t \log t}) \quad (26)$$

1137 **Step 4: LAGEX (Algorithm 2) as an optimistic saddle point algorithm** We follow the definition  
 1138 of approximate saddle point algorithm in [DKM19a]. LAGEX acts as an approximate saddle point  
 1139 algorithm if

$$\inf_{\boldsymbol{l}_t \in \mathbb{R}_+^d} \inf_{\boldsymbol{\lambda} \in \Lambda_{\hat{\mathcal{F}}_t}(\hat{\boldsymbol{\mu}}_t)} \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} d(\hat{\boldsymbol{\mu}}_{a,s}, \boldsymbol{\lambda}_a) - \boldsymbol{l}_t^T \tilde{\mathbf{A}}_t \boldsymbol{\omega}_t \geq \max_{\boldsymbol{\omega} \in \hat{\mathcal{F}}_t} \sum_{a=1}^K \sum_{s=1}^t \omega_a U_{a,s} - x_t \quad (27)$$

1140 where,  $U_{a,s} = \max \left\{ \frac{g(t)}{N_{a,s}}, \max_{\xi \in [\alpha_{a,s}, \beta_{a,s}]} d(\xi, \boldsymbol{\lambda}_{a,s}) \right\}$  and  $x_t = R_t^\omega + \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} C_{a,s}$ .

1141 The quantities  $R_t^\omega$  and  $C_{a,s}$  will be defined later on.

1142 **Step 5: Bounds cumulative regret of players** Algorithm 2 at each step solves a two player zero-sum  
 1143 game. First one is the allocation player who uses AdaGrad to maximize the inverse of characteristic

time function to find the optimal allocation. The regret of the AdaGrad player is defined at time  $t \in \mathbb{N}$  as

**Allocation player's regret.**

$$R_t^\omega = \max_{\omega \in \hat{\mathcal{F}}} \sum_{s=1}^t \sum_{a=1}^K \omega_a U_{a,s} - \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} U_{a,s}$$

We should note that AdaGrad enjoys regret of order  $R_t^\omega \leq \mathcal{O}(\sqrt{Qt})$  where  $Q$  is an upper bound on the losses such that  $Q \geq \max_{x,y \in [\mu_{\min}, \mu_{\max}]} d(x,y)$ .

**$\lambda$ -player's regret.**

$$R_\lambda^t = \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} d(\hat{\mu}_{a,s}, \lambda_{a,s}) - \inf_{\lambda \in \Lambda_{\hat{\mathcal{F}}}(\mu)} \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} d(\hat{\mu}_{a,s}, \lambda_a) \leq 0$$

The last inequality holds because we take infimum over  $\lambda$  in the perturbed alt-set. now let us prove that LAGEX is a optimistic saddle point algorithm. We define  $C_{a,s} \triangleq U_{a,s} - d(\hat{\mu}_{a,s}, \lambda_{a,s})$ . From the definition of regret of the  $\lambda$ -player we get

$$\begin{aligned} & \inf_{\lambda \in \Lambda_{\hat{\mathcal{F}}}(\mu)} \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} d(\hat{\mu}_{a,s}, \lambda_a) \\ & \geq \inf_{\mathbf{l} \in \mathbb{R}_+^d} \inf_{\lambda \in \Lambda_{\hat{\mathcal{F}}}(\mu)} \left\{ \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} U_{a,s} - \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} C_{a,s} + \mathbf{l}^T \tilde{\mathbf{A}}_t \boldsymbol{\omega}_t \right\} \end{aligned}$$

Then we have from Equation (26) and Equation (27)

$$\beta(t, \delta) \geq \inf_{\mathbf{l} \in \mathbb{R}_+^d} \inf_{\lambda \in \Lambda_{\hat{\mathcal{F}}}(\mu)} \left\{ \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} U_{a,s} - \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} C_{a,s} + \mathbf{l}^T \tilde{\mathbf{A}}_t \boldsymbol{\omega}_t \right\} \quad (28)$$

$$- (1 + \sqrt{t})K - \mathcal{O}(\sqrt{t \log t}) \quad (29)$$

$$\geq \inf_{\lambda \in \Lambda_{\hat{\mathcal{F}}}(\mu)} \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} U_{a,s} - \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} C_{a,s} \quad (30)$$

$$- (1 + \sqrt{t})K - \mathcal{O}(\sqrt{t \log t}) + \mathbf{l}_t^T \tilde{\mathbf{A}}_t \boldsymbol{\omega}_t \quad (31)$$

$$\geq \inf_{\lambda \in \Lambda_{\hat{\mathcal{F}}}(\mu)} \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} U_{a,s} - \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} C_{a,s} - (1 + \sqrt{t})K \quad (32)$$

$$- \mathcal{O}(\sqrt{t \log t}) - \mathcal{D}(\boldsymbol{\omega}_t, \hat{\mu}_t, \hat{\mathcal{F}}_t) \psi_t \quad (33)$$

where  $\psi_t$  is defined as Lemma 5. The penultimate inequality holds as we have replaced  $\inf_{\mathbf{l} \in \mathcal{L}} \mathbf{l}$  with  $\mathbf{l}_t$  i.e the optimised Lagrangian multiplier at the last step. Whereas the last inequality follows from Lemma 5. From the definition of the allocation player regret we have

$$\inf_{\lambda \in \Lambda_{\hat{\mathcal{F}}}(\mu)} \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} d(\hat{\mu}_{a,s}, \lambda_s) \geq \max_{\omega \in \hat{\mathcal{F}}} \sum_{s=1}^t \sum_{a=1}^K \omega_a U_{a,s} - R_t^\omega - \underbrace{\sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} C_{a,s}}_{T1}$$

which shows that LAGEX is a approximate saddle point algorithm with slack  $x_t = R_t^\omega + \sum_{s=1}^t \sum_{a=1}^K \omega_{a,s} C_{a,s}$ .

Now we have to ensure that  $T1$  in the slack is bounded.

$$\sum_{K+1}^t \sum_{a=1}^K \omega_{a,s} C_{a,s} \leq \sum_{K+1}^t \sum_{a=1}^K \omega_{a,s} \left\{ \frac{g(s)}{N_{a,s}} + 2L \sqrt{2\sigma^2 \frac{g(s)}{N_{a,s}}} \right\}$$

$$\begin{aligned}
&\leq g(t) \sum_{K+1}^t \sum_{a=1}^K \frac{\omega_{a,s}}{N_{a,s}} + 2L\sqrt{2\sigma^2 g(t)} \sum_{K+1}^t \sum_{a=1}^K \frac{\omega_{a,s}}{\sqrt{N_{a,s}}} \\
&\leq g(t) \left( K^2 + 2K \log \frac{t}{K} \right) + 2L\sqrt{2\sigma^2 g(t)} (K^2 + 2\sqrt{2Kt}) \\
&\leq \mathcal{O}(\sqrt{t \log t}).
\end{aligned} \tag{34}$$

1160 The inequalities hold due to the good event  $G_T$  and L-Lipschitz property of KL

$$\begin{aligned}
|d(\mu_a, \lambda_a) - d(\hat{\mu}_{a,s}, \lambda_a)| &\leq L \sqrt{2\sigma^2 \frac{g(s)}{N_{a,s}}} \\
\implies \sup_{\xi \in [\alpha_{a,s}, \beta_{a,s}]} U_{a,s} - d(\xi, \lambda_{a,s}) &\leq \max \left\{ 2L \sqrt{2\sigma^2 \frac{g(s)}{N_{a,s}}}, \frac{g(s)}{N_{a,s}} \right\}
\end{aligned}$$

1161 Now the stopping time expression changes to

$$\begin{aligned}
\beta(t, \delta) &\geq \max_{\omega \in \hat{\mathcal{F}}} \left( \sum_{a=1}^K \sum_{s=1}^t \omega_a d(\mu_a, \lambda_{a,s}) - \mathcal{D}(\omega_t, \hat{\mu}_t, \hat{\mathcal{F}}_t) \psi_t \right) - R_t^\omega - \mathcal{O}(\sqrt{t \log t}) - (1 + \sqrt{t})K \\
&\geq \max_{\omega \in \hat{\mathcal{F}}} \sum_{a=1}^K \sum_{s=1}^t \omega_a d(\mu_a, \lambda_{a,s}) - T_{\hat{\mathcal{F}}_t}^{-1}(\hat{\mu}_t) \psi_t - R_t^\omega - \mathcal{O}(\sqrt{t \log t}) - (1 + \sqrt{t})K
\end{aligned}$$

1162 **Step 6: Characteristic time** Accumulating Equation (28) and Equation (34)

$$\begin{aligned}
\max_{\omega \in \hat{\mathcal{F}}_t} \sum_{s=1}^t \sum_{a=1}^K \omega_a d(\mu_a, \lambda_{a,s}) &\geq t \inf_{\lambda \in \Lambda_{\hat{\mathcal{F}}_t}(\mu)} \max_{\omega \in \hat{\mathcal{F}}_t} \sum_{a=1}^K \omega_a d(\mu_a, \lambda_a) \\
&\geq t \max_{\omega \in \hat{\mathcal{F}}_t} \inf_{\lambda \in \Lambda_{\hat{\mathcal{F}}_t}(\mu)} \sum_{a=1}^K \omega_a d(\mu_a, \lambda_a) = t T_{\hat{\mathcal{F}}_t}^{-1}(\mu)
\end{aligned}$$

1163 Then the sample complexity is upper bounded by

$$t \leq T_{\hat{\mathcal{F}}_t}(\mu) \left( \beta(t, \delta) + R_t^\omega + \mathcal{O}(\sqrt{t \log t}) \right) + \psi_t$$

1164 Now asymptotically when  $\hat{\mathcal{F}}_t \rightarrow \mathcal{F}$  and  $\psi_t \rightarrow \mathfrak{s}$ , the expression for the characteristic time is given by

$$t \leq T_{\mathcal{F}}(\mu) \left( \beta(t, \delta) + R_t^\omega + \mathcal{O}(\sqrt{t \log t}) \right) + \mathfrak{s}$$

1165 **Part (b) : Current recommended policy is wrong.** To get on with the proof for this part we will use  
1166 similar argument as [CBJD23]. Though the argument was motivated by the work [DKM19b]. We  
1167 define the event

$$B_t \triangleq \left\{ \pi^* \neq \arg \max_{\pi \in \hat{\mathcal{F}}_t} \hat{\mu}_t^T \pi \right\}$$

1168 i.e the current recommendation policy is not correct which implies that the mean estimate or the  
1169 constraint estimate has not been concentrated yet. If we define Chernoff's information function as  
1170  $\text{ch}(u, v) \triangleq \inf_{z \in \mathcal{D}} (d(u, z) + d(v, z))$ . Therefore the current mean estimate will yield positive Cher-  
1171 noff's information since it has not been converged yet i.e  $\exists \epsilon > 0 : \text{ch}(\hat{\mu}_{a,t}, \mu_a) > \epsilon$ . Consequently  
1172 under the good event  $G_T$  defined earlier

$$\frac{g(t)}{N_{a,t}} \leq \epsilon$$

1173 since  $\text{ch}(\hat{\mu}_{a,t}, \mu_a) \leq d(\hat{\mu}_{a,t}, \mu_a)$ . Let at time  $s \in \mathbb{N}$ ,  $\pi^*$  be an extreme point in  $\hat{\mathcal{F}}_s$  that is not the  
1174 optimal policy. But since it is an extreme point in  $\hat{\mathcal{F}}_s$  that shares  $(K - 1)$  active constraints with  $\pi_{\hat{\mathcal{F}}_s}^*$ ,

1175 it has to be an optimal policy w.r.t  $\lambda \in \Lambda_{\hat{\mathcal{F}}_s}(\hat{\mu}_s) : \pi' = \arg \max_{\pi \in \hat{\mathcal{F}}_s} \lambda^T \pi \neq \pi_{\hat{\mathcal{F}}_s}^*$ . So we again  
 1176 define the event  $B_t$  as

$$B_t \triangleq \left\{ \lambda \in \Lambda_{\hat{\mathcal{F}}_t}(\hat{\mu}_t) : \pi' = \arg \max_{\pi \in \hat{\mathcal{F}}_t} \lambda^T \pi \neq \pi_{\hat{\mathcal{F}}_t}^* \right\}$$

1177 We again define  $n_{\pi'}(t)$  be the number of steps when  $\pi_s = \pi', s \in [t]$ . Therefore

$$\epsilon = \min_{l \in \mathcal{L}} \sum_{s=1, B_s} \sum_{a=1}^K \omega_{a,s} d(\hat{\mu}_{a,s}, \mu_a) \geq \min_{l \in \mathcal{L}} \sum_{\pi' \neq \pi_{\hat{\mathcal{F}}_s}^*} \inf_{\lambda \in B_s} \sum_{s=1, \pi_s=\pi'}^t \sum_{a=1}^K \omega_{a,s} d(\hat{\mu}_{a,s}, \mu_a) - l^T \tilde{\mathbf{A}}_t \omega_t \quad (35)$$

1178 Now to break the RHS of the above inequality we go back to step 5 of the proof of part (a) where  
 1179 we showed LAGEX is an approximate saddle point algorithm. In this case the slack will be  $x_t =$   
 1180  $R_{n_{\pi'}(t)}^\omega + \sum_{s=1, \pi_s=\pi'}^t \sum_{a=1}^K \omega_{a,s} C_{a,s}$ . Therefore we can write the RHS of Equation (35) as

$$\sum_{\pi' \neq \pi_{\hat{\mathcal{F}}_s}^*} \inf_{\lambda \in B_s} \sum_{s=1, \pi_s=\pi'}^t \sum_{a=1}^K \omega_{a,s} d(\hat{\mu}_{a,s}, \mu_a) \quad (36)$$

$$\geq \max_{\pi \in \hat{\mathcal{F}}_s} \min_{l \in \mathcal{L}} \underbrace{\sum_{s=1, \pi_s=\pi'}^t \sum_{a=1}^K \omega_{a,s} U_{a,s}}_{T1} - R_{n_{\pi'}(t)}^\omega - \underbrace{\sum_{s=1, \pi_s=\pi'}^t \sum_{a=1}^K \omega_{a,s} C_{a,s}}_{T2} + l^T \tilde{\mathbf{A}}_t \omega_t \quad (37)$$

1181 We apply the same logic as in [CBJD23] and [DKM19b] that  $\exists a' \in [K]$  for which  $U_{a',s} \geq \epsilon$ . That  
 1182 means the term  $T1$  grows at most linear with  $n_{\pi'}(t)$ . From the proof of part (a) it is clear that the  
 1183 term  $T2 = \mathcal{O}\left(\sqrt{n_{\pi'}(t) \log n_{\pi'}(t)}\right) \leq \mathcal{O}(\sqrt{t \log t})$  and the allocation player regret is bounded by  
 1184  $R_{n_{\pi'}(t)} = \mathcal{O}(\sqrt{Q n_{\pi'}(t)}) \leq \mathcal{O}(\sqrt{Q t})$ . That means the number of times the event  $B_t$  occurs is  
 1185 upper bounded by  $\mathcal{O}(\sqrt{t \log t})$ . Now the extra term in Equation (36) appearing with term  $T1$  and  $T2$   
 1186 induces same implication as part (a).

1187 Then incorporating part (a) and (b) to get the upper bound on the expected stopping time asymptoti-  
 1188 cally

$$\mathbb{E}[\tau_\delta] \leq T_{\hat{\mathcal{F}}_t}(\mu) \left( \beta(t, \delta) + R_t^\omega + \mathcal{O}(\sqrt{t \log t}) \right) + 2\psi_t \implies \lim_{\delta \rightarrow 0} \frac{\mathbb{E}[\tau_\delta]}{\log \frac{1}{\delta}} \leq T_{\mathcal{F}} + 2\mathfrak{s}.$$

1189 □

#### 1190 G.4 Applications to existing problems

1191 **End-of-time Knapsack.** We can model the BAI problem with end-of-time knapsack constraints as  
 1192 discussed in Section B.2. In such a setting the shadow price comes out to be  $\mathfrak{s} \leq c$  i.e the maximum  
 1193 consumable resource. So if we were to implement LATS for this the asymptotic sample complexity  
 1194 upper bound will translate to  $\alpha T_{\hat{\mathcal{F}}}(\mu)(1+c)$ , the multiplicative part being the effect of the end of  
 1195 time knapsack constraint. In case of LAGEX the unknown knapsack constraint will leave a additive  
 1196 effect quantified by  $2c$ . Recently people have deviated from only devising a no-regret learner in BwK  
 1197 rather people are interested to also give good sub-optimal guaranties on constraint violation as well.  
 1198 We think algorithms like LAGEX will perform well if we translate this model to our setting since it  
 1199 has shown not only good sample complexity but also better constraint violation guaranties as well.

1200

1201 **Fair BAI across subpopulations.** This problem is a direct consequence of our setting. The shadow  
 1202 price in this setting  $\mathfrak{s} = \frac{\max_{i \in K} \pi_i^*}{\min_{i \in K} \pi_i^*} > 1$  i.e the ratio between maximum and minimum non-zero  
 1203 weight in the recommended policy. Similar to the knapsack scenario, here also this ratio will appear  
 1204 as a extra cost of not knowing the fairness constraint in multiplicative way in case of LATS and gets

1205 added to the sample complexity upper bound of LAGEX.

1206

1207 **Pure exploration with Fairness of exposure.** We can think of a problem where we want to select a  
 1208 pool of employees from different sub-sections of a whole population for a task. As we want to max-  
 1209 imise the reward or utility of this selected group we also must also give fair exposure to all race or say  
 1210 gender. As discussed earlier in Section B.2, a direct application of our algorithms LATS and LAGEX  
 1211 to use them in the problem of pure exploration with unknown constraints on fairness of exposure.

1212 The shadow price in such a setting would be  $\mathfrak{s} = \frac{\max_{i,j \in [K]} \left( \frac{1}{\mu_i} - \frac{1}{\mu_j} \right)}{\min_{i,j \in [K]} \left( \frac{1}{\mu_i} - \frac{1}{\mu_j} \right)} = \frac{\max_{i,j \in [K]} (\mu_i - \mu_j)}{\min_{i,j \in [K]} (\mu_i - \mu_j)} \geq 1$ .

1213

1214 **Thresholding bandits.** The problem of Thresholding bandit is motivated from the safe dose  
 1215 finding problem in clinical trials, where one wants to identify the highest dose of a drug that is  
 1216 below a known safety level. From the translated optimisation problem in Section B.2 we easily  
 1217 find out the shadow price for this setting to be  $\mathfrak{s} = \frac{\max_{i \in [K]} (\pi - \theta)^i}{\min_{i \in [K]} (\pi - \theta)^i} \geq 1$ . This shadow price is  
 1218 similar to ours because the constraint structure is very similar. Our setting generalises thresholding  
 1219 bandit problem by giving the liberty of choosing different threshold levels for different support  
 1220 index of  $\pi$ . Similarly to other settings this shadow price will come as a price of handling dif-  
 1221 ferent unknown thresholds for every arm as addition in case of LAGEX and as multiplication for LATS.

1222

1223 **Feasible arm selection.** Feasible arm selection problem is motivated by the spirit of recommending  
 1224 an optimal arm which should satisfy a performance threshold. For example one might be interested to  
 1225 find a combination of food among a plethora of options which maximises the nutrient intake, rather  
 1226 the nutrient value of the food combination should exceed a threshold value. The structure of the  
 1227 optimisation problem for such a setting is discussed in detail in Section B.2. Then the shadow price  
 1228 comes out as  $\mathfrak{s} = \frac{\tau - f_{\min}}{\tau - f_{\max}} \geq 1$  where  $f \in \mathbb{R}^{\text{Supp}(\pi)}$  can be compared to a utility function. In our  
 1229 setting we will not have access to the true utility function rather we have to track it per step. This  
 1230 shadow price again get multiplied to the LATS sample complexity upper bound as a cost of not  
 1231 knowing the true utility of the arms, whereas we see a additive cost incurred in case of LAGEX.



## 1232 H Constraint violations during exploration

### 1233 H.1 Upper Bound on Constraint Violation

1234 In a linear programming problem we say constraint is violated if the chosen allocation fails to satisfy  
 1235 any of the true linear constraints. In other words when the event  $\mathbf{A}\omega_t \geq 0$ . We start with the  
 1236 optimization problem relaxed with slack if the constraints were known,

$$\begin{aligned} & \max_{\pi \in \mathcal{F}} \mu^T \pi \\ & \text{such that, } \mathbf{A}\pi + \Gamma \leq 0 \end{aligned}$$

1237 where,  $\Gamma$  is the slack. Cumulative violation of constraints can be expressed as,

$$\mathcal{V}_t = \sum_{s=1}^t \max_{i \in [K]} [\mathbf{A}^i \omega_t]_+$$

1238 where,  $[z]_+ = \max\{z, 0\}$ . Then, at any time step  $t \in [T]$ , instantaneous viola-  
 1239 tion is given by,  $v_t = \max_{i \in [K]} [\mathbf{A}^i \omega_t]_+$ . Since,  $\mathbf{A}$  is feasible, we define the game  
 1240 value as,  $\eta = \max_{\omega \in \mathcal{F}} \max_{i \in [K]} \mathbf{A}^i \omega \leq \Gamma$  and  $\eta = \max_{\omega \in \mathcal{F}} \max_{i \in [K]} \mathbf{A}^i \omega \geq$   
 1241  $\min_{\tilde{\mathbf{A}} \in \mathcal{C}_A^t} \max_{\omega \in \mathcal{F}_t} \max_{i \in [K]} \tilde{\mathbf{A}}_t^i \omega_t = \tilde{\mathbf{A}}_t^i \omega_t$ , holds because of pessimistic estimate of  $\mathbf{A}$ .

1242 Again, we define,  $i_{\min}(\omega) = \arg \min_{i \in [K]} \tilde{\mathbf{A}}^i \omega$  Then,

$$\begin{aligned} \max_{i \in [K]} [\mathbf{A}^i \omega_t] &= \max_{i \in [K]} (\mathbf{A}^i - \tilde{\mathbf{A}}_t^i) \omega_t + \max_{i \in [K]} \tilde{\mathbf{A}}_t^i \omega_t \\ &\leq \max_{i \in [K]} \|(\mathbf{A}^i - \tilde{\mathbf{A}}_t^i)\|_{\Sigma_t} \|\omega_t\|_{\Sigma_t^{-1}} + \max_{i \in [K]} \tilde{\mathbf{A}}_t^i \omega_t \\ &\leq f(t, \delta) \|\omega_t\|_{\Sigma_t^{-1}} + \max_{i \in [K]} \tilde{\mathbf{A}}_t^i \omega_t \end{aligned}$$

1243 Again,  $\max_{i \in [K]} \tilde{\mathbf{A}}_t^i \omega_t \leq \max_{i \in [K]} \mathbf{A}^i \omega_t \leq \Gamma$

1244 Then, the instantaneous violation becomes,

$$\begin{aligned} v_t &\leq [f(t, \delta) \|\omega_t\|_{\Sigma_t^{-1}} + \Gamma]_+ \\ &\leq \underbrace{f(t, \delta) \|\omega_t\|_{\Sigma_t^{-1}}}_{\rho(t, \delta)} + |\Gamma| \end{aligned}$$

1245 Let stopping time is denoted by  $\tau_\delta < T$  following the expression from the Stopping criterion section.  
 1246 Then, cumulative constraint violation is denoted by,

$$\begin{aligned} \mathcal{V}_{\tau_\delta} &= \sum_{t \leq \tau_\delta} s_t \\ &= \sum_{t \leq \tau_\delta} \rho(t, \delta) + |\Gamma| \\ &\leq 2 \sum_{t \leq \tau_\delta} |\Gamma| \mathbb{1}\{\rho(t, \delta) \leq |\Gamma|\} + 2 \sum_{t \leq \tau_\delta} \rho(t, \delta) \mathbb{1}\{|\Gamma| \leq \rho(t, \delta)\} \\ &\leq 2\tau_\delta |\Gamma| + 2 \sum_{t \leq \tau_\delta} \frac{(\rho(t, \delta))^2}{|\Gamma|}, \text{ where, } \mathbb{1}(u \leq v) \leq \frac{v}{u} \\ &\leq 2\tau_\delta |\Gamma| + \frac{6d^2 \log^2(1 + \frac{\tau_\delta + 1}{d}) + 12d \log(1 + \frac{\tau_\delta + 1}{d})(1 + \log \frac{K}{\delta})}{|\Gamma|} \end{aligned}$$

1247 The last inequality is a direct consequence of Lemma 2.

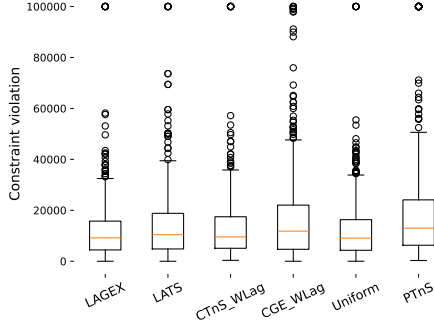


Figure 6: Constraint violation (median $\pm$ std.) algorithms over 500 runs for **hard environment**.

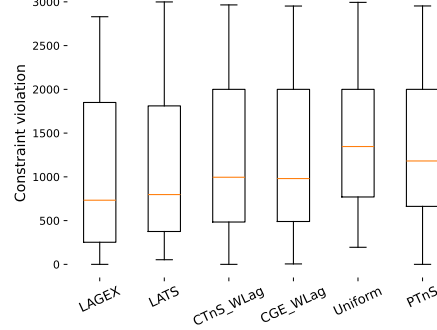


Figure 7: Constraint violation (median $\pm$ std.) of algorithms over 500 runs for **easy environment**.

## H.2 Experimental results

In all the experiments we have taken  $\delta = 0.01$ . Here, we will also explain what the abbreviations used in the plots. "CTnS-WLag" and "CGE-WLag" is basically the CTnS and CGE algorithm ([CBJD23]) under unknown constraints without Lagrangian relaxation. For the experiments in **easy environment** are clipped at 3000 for better visualisation. For CGE we have used  $g(t) = \log t$ . Each plot has been generated over 500 random seeds.

**Observation 2. LAGEX least violates constraints.** We compare the constraint violation i.e the number of times  $\mathbf{A}\omega_t > 0$ , where  $\mathbf{A}$  is the true constraint matrix.

**Observations in Environment 1.** From Fig. 6 we see that LAGEX shows least number of constraint violation across all the algorithms compared, followed by LATS and CTnS without Lagrangian relaxation shows the worst performance by having highest number of constraint violations. Though performance of Uniform explorer in terms of constraint violation is at par with LAGEX and LATS.

**Observations in Environment 2.** From Fig. 7 we can say our proposed algorithms LATS and LAGEX also better in terms of showing least number of constraint violations for the **easy environment**. An interesting application of these algorithm would be to explore BAI with end-of-time Knapsack constraint since LAGEX and LATS work "safer" than other algorithms in the unknown constraint setting with better constraint violation guaranties.

## H.3 Experiment on IMDB dataset

We evaluate our proposed algorithms Algorithm 1 and 2 with other algorithms using the publicly available and often used IMDB 50K dataset [MDP<sup>+</sup>11]. For ease of comparison we use the same bandit environment as [CBJD23] using 12 movies. We search for the optimal policy which allocates weight at most 0.3 to action movies and at least 0.3 to family and drama movies. The true optimal policy is  $[0.3, 0.3, 0, 0, 0.4, 0, 0, 0, 0, 0, 0, 0]$ . We assume  $\delta = 0.1$ . We compare the same set of algorithms as before.

**Observations 1. LAGEX shows better sample complexity** From figure 8 we can observe that the LAGEX (Algorithm 2) performs better any other algorithm in the unknown constraint setting. The algorithm LATS (Algorithm 1) performs also well on the IMDB environment but notably we cannot distinguish it's performance from the performance of the Uniform explorer as well. The diagram 8 also properly demonstrates the extra cost we had to pay due to not knowing the constraints in the allocation on the different genres of movies.

**Observation 2. LAGEX shows least constraint violation** Interestingly not only LAGEX performs as the most efficient algorithm among the all other algorithms in the unknown constraint setting but also it shows least constraint violation. It means LAGEX performs as the most safe algorithm.

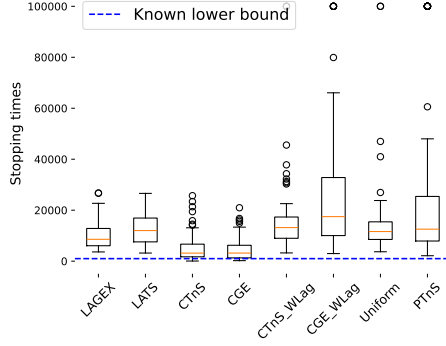


Figure 8: Sample complexity (median $\pm$ std.) of algorithms over 100 runs for **IMDB environment**.

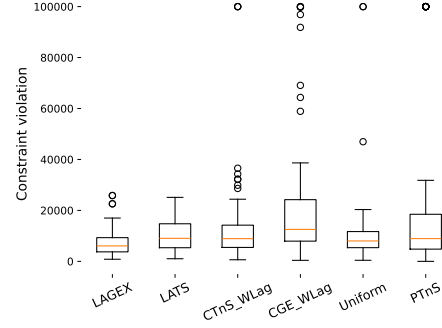


Figure 9: Constraint violation (median $\pm$ std.) of algorithms over 100 runs for **IMDB environment**.

## I $\epsilon$ -good policies under unknown linear constraints

Throughout this paper we have talked about converging to the optimal policy while tracking the estimates of the unknown linear constraints and unknown means of the reward distributions of  $K$  arm. One might not want to land on the exact optimal policy, rather may be interested in finding an  $\epsilon$ -good policy i.e the recommended policy is in a  $\epsilon$ -ball of the actual true optimal policy. This approach may get us to a much lower sample complexity lower bound [[GK21b],[MPK21], [GK21a], [MJTN20]]. We can extend our proposed Algorithms 1 and 2 can be extended to this setting by changing the definition of the set of alternative instances as  $\Lambda_{\hat{\mathcal{F}}}(\mu) \triangleq \left\{ \lambda \in \mathbb{D} : \lambda^T(\pi_{\hat{\mathcal{F}}}^* - \pi) > \epsilon \right\}$  for a pre-specified alpha as an input in the algorithms. So we will say the recommended policy in  $\epsilon$ -good for  $\pi_{\mathcal{F}}^*$  after the stopping rule fires, if  $\pi_{\hat{\mathcal{F}}}^*$  has converged to  $\pi_{\mathcal{F}}^*$  with concentration in means and constraint matrix. But there are two main problem that appears. Since  $\epsilon > 0$ , the most confusing instance for  $\mu$  will not lie on the boundary of the normal cone. So the projection lemma for lagrangian formulation (Proposition 1) is no more sufficient. Also, the algorithm may start oscillating when the allocation comes inside the epsilon ball of  $\pi_{\hat{\mathcal{F}}}^*$  among near-optimal policies since convexity of  $\omega^*(\mu)$  is no more ensured. To handle the first problem we can use the approximation error  $\epsilon$  as an *added pessimism* in the system. That means we are interested in the optimization problem

$$\begin{aligned} & \max_{\pi \in \mathcal{F}} \mu^T \pi \\ & \text{such that } \mathbf{A}\pi \leq \epsilon \mathbf{A}\mathbf{1}. \end{aligned}$$

Then to build a new superset for the new feasible set we use pessimistic estimate of  $\mathbf{A}$  i.e  $(\tilde{\mathbf{A}} - \epsilon)\pi \leq (\mathbf{A}\epsilon)\pi \leq \mathbf{A}\pi \leq 0$ . We would also want to track the sequence  $\{\epsilon_t\}_{t \in \mathbb{N}}$  and use it's concentration properties to find introduce new quantities in the lower bound that will capture the effect of the approximation. For the second hurdle we can add the notion of *sticky* approach from [DK19] that can help the agent stick to a specific  $\epsilon$ -good policy rather than oscillating.

## J Technical results and known tools in BAI and pure exploration

In this section, we will devise some technical lemma using the help of standard text on online linear regression to ensure the convergence of unknown constraints. We specifically give the expression of the radius of confidence ellipsoid mentioned in the main text in Equation 2. We then prove an upper bound on the *bad event* i.e when the constraint matrix is not concentrated around the true matrix. We also acknowledge some known theoretical results from BAI and pure exploration literature that are used in this work.

### J.1 Concentration lemma for constraints

Here, we want to get concentration on the deviation of the pessimistic estimate of the constraint matrix from the actual one quantified by  $\|(\tilde{\mathbf{A}} - \mathbf{A})\boldsymbol{\omega}\|_\infty$ , it becomes very crucial to prove upper bounds on sample complexity of our proposed algorithms. The following lemma ensures the concentration of the constraint matrix.

**Lemma 2.** *For the pessimistic estimate  $\tilde{\mathbf{A}}$  of  $\mathbf{A}$ , the following holds*

$$1. |(\tilde{\mathbf{A}}^i - \mathbf{A}^i)\boldsymbol{\omega}| \leq \rho(t, \delta) \text{ where } \rho(t, \delta) \triangleq f(t, \delta)\|\boldsymbol{\omega}\|_{\Sigma_t^{-1}}$$

$$2. \sum_{s=1}^t \|\boldsymbol{\omega}_s\|_{\Sigma_t^{-1}}^2 \leq 2d \log \left(1 + \frac{1+t}{d}\right)$$

$$3. \sum_{s=1}^t \rho(t, \delta) \leq \sqrt{2dt f^2(t, \delta) \log \left(1 + \frac{1+t}{d}\right)}$$

*Proof.* The first result gives control on the deviations  $\tilde{\mathbf{A}}\boldsymbol{\omega} - \mathbf{A}\boldsymbol{\omega}$  for  $\mathbf{A} \in \mathcal{C}_t^{\mathbf{A}}(\delta)$ . Then  $\forall i \in [d]$

$$|(\tilde{\mathbf{A}}^i - \mathbf{A}^i)\boldsymbol{\omega}| \leq |(\tilde{\mathbf{A}}_t^i - \mathbf{A}^i)\boldsymbol{\omega}| \leq 2 \sup_{\mathbf{A} \in \mathcal{C}_t^{\mathbf{A}}(\delta)} \|\tilde{\mathbf{A}}_t^i - \mathbf{A}^i\|_{\Sigma_t} \|\boldsymbol{\omega}\|_{\Sigma_t^{-1}} \leq f(t, \delta)\|\boldsymbol{\omega}\|_{\Sigma_t^{-1}} \leq \rho(t, \delta)$$

here, we define  $\rho(t, \delta) \triangleq f(t, \delta)\|\boldsymbol{\omega}\|_{\Sigma_t^{-1}}$ . The penultimate inequality follows from the definition of the confidence set defined in Equation 2. Now we want to derive an explicit expression of this upper bound. It is natural to use the concentration of the gram matrix  $\Sigma_t$  over time. We refer to [AyPS11] for the control over the behaviour of  $\Sigma_t$  and we directly get the second as

$$\sum_{s=1}^t \|\boldsymbol{\omega}_s\|_{\Sigma_t^{-1}}^2 \leq 2 \log \det \Sigma_{t+1} \leq 2d \log \left(1 + \frac{1+t}{d}\right)$$

Refer [AyPS11] for the context. Now we have to control the cumulative deviation because later on when we define the bad event based on this concentration we will need to know the cumulative behaviour of  $\rho(t, \delta)$ .

Then for an arbitrary sequence of actions  $\{\boldsymbol{\omega}_s\}_{s \in [T]}$

$$\sum_{s=1}^t \rho(t, \delta) \leq \sqrt{t \sum_{s=1}^t \rho^2(t, \delta)} \leq \sqrt{2dt f^2(t, \delta) \log \left(1 + \frac{1+t}{d}\right)}$$

where,  $\sum_{s=1}^t \rho^2(t, \delta) \leq 2dt f^2(t, \delta) \left(1 + \frac{1+t}{d}\right)$  using result 2 of this lemma. This holds because as we have already stated  $\{f(s, \delta)\}_{s \in [T]}$  is a non-decreasing sequence of function and  $f(t, \delta)$  is the maximum possible value in the set i.e  $\sum_{s=1}^t f^2(t, \delta) \leq t f^2(t, \delta)$

□

Now we proceed to state an upper bound on the cumulative probability of the *bad event* i.e the event  $\|(\tilde{\mathbf{A}}_t - \mathbf{A})\boldsymbol{\omega}\|_\infty > |(\tilde{\mathbf{A}}_t^i - \mathbf{A}^i)\boldsymbol{\omega}| > \rho(t, \delta)$ .

**Lemma 3.** *The cumulative probability of the bad event till time  $t < T$ ,*

$$\sum_{s=1}^t \mathbb{P} \left( \|(\tilde{\mathbf{A}}_t - \mathbf{A})\boldsymbol{\omega}\|_\infty > \rho(t, \delta) \right) \leq \zeta(1)$$

where  $\zeta(\cdot)$  is the Euler-Riemann zeta function.

1335 *Proof.* We already have stated in the main paper  $f(t, \delta) = 1 + \sqrt{\frac{1}{2} \log \frac{K}{\delta} + \frac{1}{4} \log \det \Sigma_t} \leq$   
1336  $1 + \sqrt{\frac{1}{2} \log \frac{K}{\delta} + \frac{d}{4} \log \left(1 + \frac{t}{d}\right)} \triangleq f'(t, \delta)$  by Lemma 2 and we define  
1337  $\rho'(t, \delta) \triangleq f'(t, \delta) \|\omega\|_{\Sigma_t^{-1}}$ . It implies  $\mathbb{P} \left( \exists t \in [T] \|(\tilde{\mathbf{A}}_t - \mathbf{A})\omega\|_{\infty} > \rho'(t, \delta) \right) \leq$   
1338  $\mathbb{P} \left( \exists t \in [T] \|(\tilde{\mathbf{A}}_t - \mathbf{A})\omega\|_{\infty} > \rho(t, \delta) \right) \leq \delta$ . Now if we replace  $\log \frac{1}{\delta}$  by  $u$ , we can write  
1339  $\mathbb{P} \left( \exists t \in [T], \forall i \in [d] \| \tilde{\mathbf{A}}_t^i - \mathbf{A}^i \|_{\Sigma_t} > 1 + \sqrt{\frac{1}{2} \log K + \frac{u}{2} + \frac{d}{4} \log \left(1 + \frac{t}{d}\right)} \right) \leq \exp(-u)$ . We can  
1340 directly assign  $\log t$  as the simplest and natural choice for  $u$ , since  $\sum_{s=1}^{\infty} \frac{1}{t} = \zeta(1)$ ,  $\zeta(\cdot)$  being the  
1341 Euler-Riemann zeta function. Though this integral is improper, it has a Cauchy principal value as  
1342 Euler-Mascheroni constant which means  $\sum_{s=1}^{\infty} \frac{1}{t} \approx \gamma = 0.577$ . So we assign  $u = \log t$

$$\sum_{s=1}^t \mathbb{P} \left( \|(\tilde{\mathbf{A}}_t - \mathbf{A})\omega\|_{\infty} > \rho(t, \delta) \right) \leq \sum_{s=1}^t \frac{1}{t} \leq \sum_{s=1}^{\infty} \frac{1}{t} \leq \zeta(1) \approx 0.577$$

1343 □

1344 **Lemma 4.** Let  $\bar{\mu} \geq 0$  be a vector, and consider the set  $Q_{\bar{\mu}} = \{\mu \geq 0 \mid q(\mu) \geq q(\bar{\mu})\}$ . Let  
1345 Slater condition hold. Then, the set  $Q_{\bar{\mu}}$  is bounded and, in particular, we have  $\|\mu\|_1 \leq \frac{1}{\gamma} (f(\bar{x}) -$   
1346  $q(\bar{\mu}))$ ,  $\forall \mu \in Q_{\bar{\mu}}$  where,  $\gamma = \min_{1 \leq j \leq m} \{-g_j(\bar{x})\}$  and  $\bar{x}$  is a Slater vector.  $f(\cdot)$  and  $q(\cdot)$  respectively  
1347 denotes the primal and the dual function of the optimization problem.

1348 Using the aforementioned lemmas we give a bound for the part in the inverse of the characteristic time  
1349 function that gets added for Lagrangian relaxation which eventually helps up landing on a unique  
1350 formulation of sample complexity upper bounds of our proposed algorithm.

1351 **Lemma 5.** For any  $l \in \mathcal{L}$  and  $\omega \in \Delta_K$

$$-l^T \tilde{\mathbf{A}} \omega \leq \mathcal{D}(\mu, \omega, \hat{\mathcal{F}}) \psi$$

1352 where,  $\psi = \frac{\|(\tilde{\mathbf{A}} - \mathbf{A})\omega\|_{\infty} + \max_{i \in [1, N]} (-\mathbf{A}^i \omega)}{\min_{i \in [1, N]} (-\tilde{\mathbf{A}}^i \omega)}$

1353 *Proof.* For any  $l \in \mathcal{L}$  and  $\omega \in \Delta_K$  we write

$$\begin{aligned} (-l^T \tilde{\mathbf{A}} \omega) &= l^T (-\tilde{\mathbf{A}} + \mathbf{A} - \mathbf{A}) \omega \\ &\leq \|l\|_1 \|(\mathbf{A} - \tilde{\mathbf{A}}) \omega\|_{\infty} + \|l\|_1 \max_{i \in [1, N]} (-\mathbf{A}^i \omega) \\ &\leq \|l\|_1 \left( \|(\mathbf{A} - \tilde{\mathbf{A}}) \omega\|_{\infty} + \max_{i \in [1, N]} \Gamma \right) \\ &\leq \mathcal{D}(\omega, \mu, \hat{\mathcal{F}}) \frac{\|(\tilde{\mathbf{A}} - \mathbf{A})\omega\|_{\infty} + \max_{i \in [1, N]} (-\mathbf{A}^i \omega)}{\min_{i \in [1, N]} (-\tilde{\mathbf{A}}^i \omega)} \end{aligned}$$

1354 Plugging in the definition of  $\psi$  mentioned in the statement of the lemma concludes the proof. □

## 1355 J.2 Useful results from BAI and pure exploration literature

1356 **Lemma 6.** (Lemma 19 in [GK16]) There exists two constants  $B$  and  $C$  (depends on  $\mu$  and  $\epsilon$ ), such  
1357 that—

$$\sum_{t=h(T)}^T \mathbb{P} \{ \|\hat{\mu}_t - \mu\|_{\infty} > \xi(\epsilon) \} \leq BT \exp(-CT^{\frac{1}{5}})$$

1358 **Lemma 7.** [GK16, Lemma 7] For all  $t \geq 1$  and  $\forall a \in [K]$ ,  $C$ -Tracking ensures  $N_{a,t} \geq \sqrt{t + K^2} - K$   
1359 and

$$\max_{a \in [K]} |N_{a,t} - \sum_{s=1}^t \omega_{a,s}| \leq K(1 + \sqrt{t})$$

**Lemma 8.** (Theorem 14 in [KK21]) Let  $\delta > 0, \nu$  be independent one-parameter exponential families with mean  $\mu$  and  $S \subset [d]$ . Then we have,

$$\mathbb{P}_\nu \left[ \exists t \in \mathbb{N} : \sum_{a \in S} \tilde{N}_{t,a} d_{KL}(\mu_{t,a}, \mu_a) \geq \sum_{a \in S} 3 \ln \left( 1 + \ln \left( \tilde{N}_{t,a} \right) \right) + |S| \mathcal{T} \left( \frac{\ln \left( \frac{1}{\delta} \right)}{|S|} \right) \right] \leq \delta.$$

1360 Here,  $\mathcal{T} : \mathbb{R}^+ \rightarrow \mathbb{R}^+$  is such that  $\mathcal{T}(x) = 2\tilde{h}_{3/2} \left( \frac{h^{-1}(1+x) + \ln \left( \frac{\pi^2}{3} \right)}{2} \right)$  with

$$\forall u \geq 1, \quad h(u) = u - \ln(u) \quad (38)$$

$$\forall z \in [1, e], \forall x \geq 0, \quad \tilde{h}_z(x) = \begin{cases} \exp \left( \frac{1}{h^{-1}(x)} \right) h^{-1}(x) & \text{if } x \geq h^{-1} \left( \frac{1}{\ln(z)} \right) \\ z(x - \ln(\ln(z))) & \text{else} \end{cases} \quad (39)$$

**Lemma 9.** (Lemma 17 in [DK19]) Under the good event  $G_T$ , there exists a  $T_\epsilon$  such that for  $T$  where  $h(T) \geq T_\epsilon$   $C$ -tracking will satisfy

$$\inf_{w \in w^*(\mu)} \left\| \frac{N_t}{t} - w \right\|_\infty \leq 3\epsilon, \forall t \geq 4 \frac{K^2}{\epsilon^2} + 3 \frac{h(T)}{\epsilon}$$

**Lemma 10.** (Theorem 2 in [DKM19b]) The sample complexity of GE is

$$\mathbb{E}[\tau] \leq T_0(\delta) + \frac{eK}{a}$$

where

$$T_0(\delta) = \max \left\{ t \in \mathbb{N} : t \leq T(\mu)c(t, \delta) + C_\mu \left( R_t^\lambda + R_t^w + O(\sqrt{t \log t}) \right) \right\}$$

1361 where  $R_t^\lambda$  is the regret of the instance player,  $R_t^w$  the regret of the allocation player and  $C_\mu$  an  
1362 instance-dependent constant.

### 1363 J.3 Useful definitions and theorems from literature on continuity of convex functions

1364 **Definition 2.** (Definition of Upper Hemicontinuity) We say that a set-valued function  $C : \Theta \rightarrow \omega$  is  
1365 upper hemicontinuous at the point  $\theta \in \Theta$  if for any open set  $S \subset \omega$  with  $C(\theta) \in S$  there exists a  
1366 neighborhood  $U$  around  $\theta$ , such that  $\forall x \in U, C(x)$  is a subset of  $S$ .

1367 **Theorem 9.** (Berge's maximum theorem, [Ber63]) Let  $X$  and  $\Theta$  be topological spaces. Let  $f : X \times \Theta \rightarrow \mathbb{R}$   
1368 be a continuous function and let  $C : \Theta \rightarrow \bar{X}$  be a compact-valued correspondence  
1369 such that  $C(\theta) \neq \emptyset \forall \theta \in \Theta$ . If  $C$  is continuous at  $\theta$  then  $f^*(\theta) = \sup_{x \in C(\theta)} f(x, \theta)$  is continuous  
1370 and  $C^* = \{x \in C(\theta) : f(x, \theta) = f^*(\theta)\}$  is upper hemicontinuous.

1371 **Theorem 10.** (Heine-Borel theorem, Eduard Heine and Émile Borel) For a subset  $S$  in  $\mathbb{R}^n$ , the  
1372 following two statements are equivalent

- 1373 1.  $S$  is closed and bounded.
- 1374 2.  $S$  is compact, means every open cover of  $S$  has a finite sub-cover.

1375 **Theorem 11.** Let  $C$  be a closed convex set with nonempty (topological) interior. Let  $f$  and  $\{f^r\}$  be  
1376 affine functions from  $E^n$  to  $E^m$  with  $f^r \rightarrow f$ . Then

$$1377 \text{ (II.1.2) } \overline{\lim}_{r \rightarrow \infty} (H(f^r) \cap C) \subset H(f) \cap C$$

$$1378 \text{ (II.1.3) } \underline{\lim}_{r \rightarrow \infty} (H(f^r) \cap C) \text{ is a closed convex subset of } H(f) \cap C,$$

1379 (II.1.4). If  $H(f) \cap C$  has nonempty interior and no component of  $f$  is identically zero, then  
1380  $\lim_{r \rightarrow \infty} (H(f^r) \cap C) = H(f) \cap C$

1381 **Lemma 11.** [MCP14, Lemma 13] Consider  $A \in (\mathbb{R}^+)^{k \times k}$ ,  $c \in (\mathbb{R}^+)^k$ , and  $\mathcal{T} \subset (\mathbb{R}^+)^{k \times k} \times (\mathbb{R}^+)^k$ .  
1382 Define  $t = (A, c)$ . Consider the function  $Q$  and the set-valued map  $Q^*$

$$Q(t) = \inf_{x \in \mathbb{R}^k} \{cx \mid Ax \geq 1, x \geq 0\}$$

$$Q^*(t) = \{x : cx \leq Q(t) \mid Ax \geq 1, x \geq 0\}.$$

1383 Assume that: For all  $t \in \mathcal{T}$ , all rows and columns of  $A$  are non-identically 0 and  $\min_{t \in \mathcal{T}} \min_k c_k > 0$ .  
1384 Then, 1.  $Q$  is continuous on  $\mathcal{T}$ , 2.  $Q^*$  is upper-hemicontinuous on  $\mathcal{T}$ .

## 1385 NeurIPS Paper Checklist

### 1386 1. Claims

1387 Question: Do the main claims made in the abstract and introduction accurately reflect the  
1388 paper's contributions and scope?

1389 Answer: [\[Yes\]](#)

1390 Justification: The abstract and introduction section clearly and accurately reflect the contri-  
1391 butions and scopes of this work.

1392 Guidelines:

- 1393 • The answer NA means that the abstract and introduction do not include the claims  
1394 made in the paper.
- 1395 • The abstract and/or introduction should clearly state the claims made, including the  
1396 contributions made in the paper and important assumptions and limitations. A No or  
1397 NA answer to this question will not be perceived well by the reviewers.
- 1398 • The claims made should match theoretical and experimental results, and reflect how  
1399 much the results can be expected to generalize to other settings.
- 1400 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
1401 are not attained by the paper.

### 1402 2. Limitations

1403 Question: Does the paper discuss the limitations of the work performed by the authors?

1404 Answer: [\[Yes\]](#)

1405 Justification: In Section 6 we have clearly stated the limitations of our work. Though our  
1406 focus in this work has been a set of unknown linear constraints, but the constraints also can  
1407 be non-linear in nature.

1408 Guidelines:

- 1409 • The answer NA means that the paper has no limitation while the answer No means that  
1410 the paper has limitations, but those are not discussed in the paper.
- 1411 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 1412 • The paper should point out any strong assumptions and how robust the results are to  
1413 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
1414 model well-specification, asymptotic approximations only holding locally). The authors  
1415 should reflect on how these assumptions might be violated in practice and what the  
1416 implications would be.
- 1417 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
1418 only tested on a few datasets or with a few runs. In general, empirical results often  
1419 depend on implicit assumptions, which should be articulated.
- 1420 • The authors should reflect on the factors that influence the performance of the approach.  
1421 For example, a facial recognition algorithm may perform poorly when image resolution  
1422 is low or images are taken in low lighting. Or a speech-to-text system might not be  
1423 used reliably to provide closed captions for online lectures because it fails to handle  
1424 technical jargon.
- 1425 • The authors should discuss the computational efficiency of the proposed algorithms  
1426 and how they scale with dataset size.
- 1427 • If applicable, the authors should discuss possible limitations of their approach to  
1428 address problems of privacy and fairness.
- 1429 • While the authors might fear that complete honesty about limitations might be used by  
1430 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
1431 limitations that aren't acknowledged in the paper. The authors should use their best  
1432 judgment and recognize that individual actions in favor of transparency play an impor-  
1433 tant role in developing norms that preserve the integrity of the community. Reviewers  
1434 will be specifically instructed to not penalize honesty concerning limitations.

### 1435 3. Theory Assumptions and Proofs

1436 Question: For each theoretical result, does the paper provide the full set of assumptions and  
1437 a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: For each theoretical result like theorems or lemmas, we have clearly stated the assumptions. The reader can find all the proofs in complete in various sections of the appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: We have extensively stated all the parameters and function values considered in the numerical simulations in Section 5 and we have also added additional information of the experiments in Section H.2 one need to reproduce all the result that conforms with all the claims made in this paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in



1492 some way (e.g., to registered users), but it should be possible for other researchers  
1493 to have some path to reproducing or verifying the results.

## 1494 5. Open access to data and code

1495 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
1496 tions to faithfully reproduce the main experimental results, as described in supplemental  
1497 material?

1498 Answer: [Yes]

1499 Justification: the experiments in this paper are simulations and uses synthetic data which is  
1500 clearly stated how to generate in Section 5. We have given the link to the python codes used  
1501 to do the experiments in Section 5

1502 Guidelines:

- 1503 • The answer NA means that paper does not include experiments requiring code.
- 1504 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/](https://nips.cc/public/guides/CodeSubmissionPolicy)  
1505 [public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1506 • While we encourage the release of code and data, we understand that this might not be  
1507 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
1508 including code, unless this is central to the contribution (e.g., for a new open-source  
1509 benchmark).
- 1510 • The instructions should contain the exact command and environment needed to run to  
1511 reproduce the results. See the NeurIPS code and data submission guidelines ([https://](https://nips.cc/public/guides/CodeSubmissionPolicy)  
1512 [nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1513 • The authors should provide instructions on data access and preparation, including how  
1514 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1515 • The authors should provide scripts to reproduce all experimental results for the new  
1516 proposed method and baselines. If only a subset of experiments are reproducible, they  
1517 should state which ones are omitted from the script and why.
- 1518 • At submission time, to preserve anonymity, the authors should release anonymized  
1519 versions (if applicable).
- 1520 • Providing as much information as possible in supplemental material (appended to the  
1521 paper) is recommended, but including URLs to data and code is permitted.

## 1522 6. Experimental Setting/Details

1523 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
1524 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
1525 results?

1526 Answer: [Yes]

1527 Justification: Section 5 and Section H.2 clearly describes how the synthetic data are generated  
1528 with all the details about parameters, optimizer etc.

1529 Guidelines:

- 1530 • The answer NA means that the paper does not include experiments.
- 1531 • The experimental setting should be presented in the core of the paper to a level of detail  
1532 that is necessary to appreciate the results and make sense of them.
- 1533 • The full details can be provided either with the code, in appendix, or as supplemental  
1534 material.

## 1535 7. Experiment Statistical Significance

1536 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
1537 information about the statistical significance of the experiments?

1538 Answer: [Yes]

1539 Justification: We have reported error bars, confidence intervals and statistical significance of  
1540 all the experiments suitably and defined correctly that conforms with the main claims of our  
1541 paper.

1542 Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We have stated the computer processor, memory and other specifications clearly in Section 5 which is used to carry out all the experiments given in this paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: The author of this paper are completely aware of the NeurIPS Code of Ethics and this paper totally conforms with it.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: This work can be considered generalisation or special case of some settings which can be used for positive social impacts like for example designing clinical trials for patients or designing a diet plan. We have stated some of the positive application of our work in the motivation of our work in Section B.2.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This work does not involve any data or models that have a high risk for misuse. It only uses synthetic data for simulation experiments.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Every assets that has been used in this paper are explicitly mentioned and properly respected.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets` has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [\[Yes\]](#)

Justification: All the new assets introduced in this paper for example, new algorithms and codes are properly documented.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [\[NA\]](#)

Justification: This work does not involve any crowdsourcing experiments or any experiments engaging human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

### 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

1700 Answer: [NA]  
1701 Justification: This work does not involve any crowdsourcing experiments or any experiments  
1702 engaging human subjects.  
1703 Guidelines:  
1704 • The answer NA means that the paper does not involve crowdsourcing nor research with  
1705 human subjects.  
1706 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
1707 may be required for any human subjects research. If you obtained IRB approval, you  
1708 should clearly state this in the paper.  
1709 • We recognize that the procedures for this may vary significantly between institutions  
1710 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
1711 guidelines for their institution.  
1712 • For initial submissions, do not include any information that would break anonymity (if  
1713 applicable), such as the institution conducting the review.