
Closing the Computational-Statistical Gap in Best Arm Identification for Combinatorial Semi-bandits

Ruo-Chun Tzeng

EECS
KTH, Stockholm, Sweden
rctzeng@kth.se

Po-An Wang

EECS
KTH, Stockholm, Sweden
wang9@kth.se

Alexandre Proutiere

EECS and Digital Futures
KTH, Stockholm, Sweden
alepro@kth.se

Chi-Jen Lu

Institute of Information Science
Academia Sinica, Taiwan
cjlu@iis.sinica.edu.tw

Abstract

We study the best arm identification problem in combinatorial semi-bandits in the fixed confidence setting. We present Perturbed Frank-Wolfe Sampling (P -FWS), an algorithm that (i) runs in polynomial time, (ii) achieves the instance-specific minimal sample complexity in the high confidence regime, and (iii) enjoys polynomial sample complexity guarantees in the moderate confidence regime. To the best of our knowledge, even for the vanilla bandit problems, no algorithm was able to achieve (ii) and (iii) simultaneously. With P -FWS, we close the computational-statistical gap in best arm identification in combinatorial semi-bandits. The design of P -FWS starts from the optimization problem that defines the information-theoretical and instance-specific sample complexity lower bound. P -FWS solves this problem in an online manner using, in each round, a single iteration of the Frank-Wolfe algorithm. Structural properties of the problem are leveraged to make the P -FWS successive updates computationally efficient. In turn, P -FWS only relies on a simple linear maximization oracle.

1 Introduction

An efficient method to design statistically optimal algorithms solving active learning tasks (e.g., regret minimization or pure exploration in bandits and reinforcement learning) consists in the following two-step procedure. The first step amounts to deriving, through change-of-measure arguments, tight information-theoretical fundamental limits satisfied by a wide class of learning algorithms. These limits are often expressed as the solution of an optimization problem, referred in this paper to as the *lower-bound problem*. Interestingly, this solution specifies the instance-specific optimal exploration process: it characterizes the limiting behavior of the adaptive sampling rule any statistically optimal algorithm should implement. In the second step, the learning algorithm is designed so that its exploration process approaches the solution of the lower-bound problem. This design yields statistically optimal algorithms, but typically requires to repeatedly solve the lower-bound problem. This method has worked remarkably well for simple learning tasks such as regret minimization or best-arm identification with fixed confidence in classical stochastic bandits [Lai87, GC11, GK16], but also in bandits whose arm-to-average reward function satisfies simple structural properties (e.g., Lipschitz, unimodal) [MCP14, WTP21].

The method also provides a natural way of studying the computational-statistical gap [KLLM22] for active learning tasks. Indeed, if solving the lower-bound problem in polynomial time is possible, one

may hope to devise learning algorithms that are both statistically optimal and computationally efficient. As of now, however, the computational complexity of the lower-bound problem remains largely unexplored, except for simple learning tasks. For example, in the case of best policy identification in tabular Markov Decision Processes, the lower-bound problem is non-convex [AMP21] and its complexity and approximability are unclear.

In this paper, we leverage the aforementioned two-step procedure to assess the computational-statistical gap for the best arm identification in combinatorial semi-bandits in the fixed confidence setting. We establish that, essentially, this gap does not exist (a result that was conjectured in [JMCK21]). Specifically, we present an algorithm that enjoys the following three properties: (i) it runs in polynomial time, (ii) its sample complexity matches the fundamental limits asymptotically in the high confidence regime, and (iii) its sample complexity is at most polynomial in the moderate confidence regime. Next, after formally introducing combinatorial semi-bandits, we describe our contributions and techniques in detail.

Best arm identification in combinatorial semi-bandits. In combinatorial semi-bandits [CBL12, CTMSP⁺15], the learner sequentially selects an action from a combinatorial set $\mathcal{X} \subset \{0, 1\}^K$. When in round t , the action $\mathbf{x}(t) = (x_1(t), \dots, x_K(t)) \in \mathcal{X}$ is chosen, the environment samples a K -dimensional vector $\mathbf{y}(t)$ whose distribution is assumed to be Gaussian $\mathcal{N}(\boldsymbol{\mu}, \mathbf{I})$. The learner then receives the detailed reward vector $\mathbf{x}(t) \odot \mathbf{y}(t)$ where \odot denotes the element-wise product (in other words, the learner observes the individual reward $y_k(t)$ of the arm k if and only if this arm is selected in round t , i.e., $x_k(t) = 1$). The parameter $\boldsymbol{\mu}$ characterizing the average rewards of the various arms is initially unknown. The goal of a learner is to identify the best action $\mathbf{i}^*(\boldsymbol{\mu}) = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x}, \boldsymbol{\mu} \rangle$ with a given level of confidence $1 - \delta$, for some $\delta > 0$ while minimizing the expected number of rounds needed. We assume that the best action is unique and denote by $\Lambda = \{\boldsymbol{\mu} \in \mathbb{R}^K : |\mathbf{i}^*(\boldsymbol{\mu})| = 1\}$ the set of parameters satisfying this assumption. The learner strategy is defined by three components: (i) a sampling rule dictating the sequence of the selected actions; (ii) a stopping time τ defining the last round where the learner interacts with the environment; (iii) a decision rule specifying the action $\hat{\mathbf{i}} \in \mathcal{X}$ believed to be optimal based on the data gathered until τ .

The sample complexity lower-bound problem. Consider the set of δ -PAC algorithms such that for any $\boldsymbol{\mu} \in \Lambda$, the best action is identified correctly with probability at least $1 - \delta$. We wish to find a δ -PAC algorithm with minimal expected sample complexity $\mathbb{E}_{\boldsymbol{\mu}}[\tau]$. To this aim, using classical change-of-measure arguments [GK16], we may derive a lower bound of the expected sample complexity satisfied by any δ -PAC algorithm. This lower bound is given by¹ $\mathbb{E}_{\boldsymbol{\mu}}[\tau] \geq T^*(\boldsymbol{\mu}) \operatorname{kl}(\delta, 1 - \delta)$. The characteristic time $T^*(\boldsymbol{\mu})$ is defined as the value of the following problem

$$T^*(\boldsymbol{\mu})^{-1} = \sup_{\boldsymbol{\omega} \in \Sigma} \inf_{\boldsymbol{\lambda} \in \operatorname{Alt}(\boldsymbol{\mu})} \left\langle \boldsymbol{\omega}, \frac{(\boldsymbol{\mu} - \boldsymbol{\lambda})^2}{2} \right\rangle, \quad (1)$$

where² $\Sigma = \{\sum_{\mathbf{x} \in \mathcal{X}} w_{\mathbf{x}} \mathbf{x} : \mathbf{w} \in \Sigma_{|\mathcal{X}|}\}$, $\operatorname{kl}(a, b)$ is the KL-divergence between two Bernoulli distributions with respective means a and b , and $\operatorname{Alt}(\boldsymbol{\mu}) = \{\boldsymbol{\lambda} \in \Lambda : \mathbf{i}^*(\boldsymbol{\lambda}) \neq \mathbf{i}^*(\boldsymbol{\mu})\}$ is the set of *confusing* parameters. As it turns out (see Lemma 1), $T^*(\boldsymbol{\mu})$ is at most quadratic in K , and hence the sample complexity lower bound is polynomial. (1) is a concave program over Σ [WTP21], and a point $\boldsymbol{\omega}^*$ in its solution set corresponds to an optimal allocation of action draws: an algorithm sampling actions according to $\boldsymbol{\omega}^*$ and equipped with an appropriate stopping rule would yield a sample complexity matching the lower bound. In this paper, we provide computationally efficient algorithms to solve (1) and show how these can be used to devise a δ -PAC best action identification algorithm with minimal sample complexity and running in polynomial time. We only assume that we have access to a computationally efficient Oracle, referred to as the LM (Linear Maximization) Oracle, identifying the best action should $\boldsymbol{\mu}$ be known (but for any possible $\boldsymbol{\mu}$). This assumption, made in all previous work on combinatorial semi-bandits (see e.g. [JMCK21, PBVP20]), is crucial as indeed, if there is no computationally efficient algorithm solving the offline problem $\operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x}, \boldsymbol{\mu} \rangle$ with known $\boldsymbol{\mu}$, there is no hope to solve its online version with unknown $\boldsymbol{\mu}$ in a computationally efficient manner. The assumption holds for a large array of combinatorial sets of actions [S⁺03], including m -sets, matchings, (source–destination)-paths, spanning trees, matroids (refer to [CCG21b] for a thorough discussion).

¹We present proof in Appendix K for completeness – see also [JMCK21].

² Σ_N denotes the $(N - 1)$ dimensional simplex.

The Most-Confusing-Parameter (MCP) algorithm. The difficulty of solving (1) lies in the inner optimization problem, i.e., in evaluating the objective function:

$$F_\mu(\omega) = \inf_{\lambda \in \text{Alt}(\mu)} \left\langle \omega, \frac{(\mu - \lambda)^2}{2} \right\rangle = \min_{x \neq i^*(\mu)} f_x(\omega, \mu) \quad (2)$$

where $f_x(\omega, \mu) = \inf_{\lambda \in \mathcal{C}_x} \langle \omega, \frac{(\mu - \lambda)^2}{2} \rangle$ and $\mathcal{C}_x = \{\lambda \in \mathbb{R}^K : \langle \lambda, i^*(\mu) - x \rangle < 0\}$. Evaluating $F_\mu(\omega)$ is required to solve (1), but also in the design of an efficient stopping rule. Our first contribution is MCP (Most-Confusing-Parameter), a polynomial time algorithm able to approximate $F_\mu(\omega)$ for any given μ and ω . The algorithm's name refers to the fact that by computing $F_\mu(\omega)$, we implicitly identify the *most confusing parameter* $\lambda^* \in \arg \inf_{\lambda \in \text{Alt}(\mu)} \langle \omega, \frac{(\mu - \lambda)^2}{2} \rangle$. The design of MCP leverages a Lagrangian relaxation of the optimization problem defining $f_x(\omega, \mu)$ and exploits the fact that the Lagrange dual function linearly depends on x . In turn, this linearity allows us to make use of the LM Oracle. From these observations, we show that computing $F_\mu(\omega)$ boils down to solving a two-player game, for which one of the players can simply update her strategy using the LM Oracle. We prove the following informally stated theorem quantifying the performance of the MCP algorithm (see Theorem 3 for a more precise statement).

Theorem 1. *For any (ω, μ) , the MCP algorithm with precision ϵ and certainty parameter θ returns \hat{F} and \hat{x} satisfying $\mathbb{P}_\mu[F_\mu(\omega) \leq \hat{F} \leq (1 + \epsilon)F_\mu(\omega)] \geq 1 - \theta$ and $\hat{F} = f_{\hat{x}}(\omega, \mu)$. The number of calls to the LM Oracle is, almost surely, at most polynomial in K , ϵ^{-1} , and $\ln \theta^{-1}$.*

The Perturbed Frank-Wolfe Sampling (P-FWS) algorithm. The MCP algorithm allows us to solve the lower-bound problem (1) for any given μ . The latter is initially unknown, but could be estimated. A Track-and-Stop algorithm [GK16] solving (1) with this plug-in estimator in each round would yield asymptotically minimal sample complexity. It would however be computationally expensive. To circumvent this difficulty, as in [WTP21], our algorithm, P-FWS, performs a single iteration of the Frank-Wolfe algorithm for the program (1) instantiated with an estimator of μ . To apply the Frank-Wolfe algorithm, P-FWS uses stochastic smoothing techniques to approximate the non-differentiable objective function F_μ by a smooth function. To estimate the gradient of the latter, P-FWS leverages both the LM Oracle and the MCP algorithm (more specifically its second output \hat{x}). Finally, P-FWS stopping rule takes the form of a classical Generalized Likelihood Ratio Test (GLRT) where the estimated objective function is compared to a time-dependent threshold. Hence the stopping rule also requires the MCP algorithm. We analyze the sample and computational complexities of P-FWS. Our main results are summarized in the following theorem (refer to Theorem 4 for details).

Theorem 2. *For any $\delta \in (0, 1)$, P-FWS is δ -PAC, and for any $(\epsilon, \tilde{\epsilon}) \in (0, 1)$ small enough, its sample complexity satisfies:*

$$\mathbb{E}_\mu[\tau] \leq \frac{(1 + \tilde{\epsilon})^2}{T^*(\mu)^{-1} - \epsilon} \times H\left(\frac{1}{\delta} \cdot \frac{c(1 + \tilde{\epsilon})^2}{T^*(\mu)^{-1} - \epsilon}\right) + \Psi(\epsilon, \tilde{\epsilon}),$$

where $H(x) = \ln(x) + \ln \ln(x)$, $c > 0$ is a universal constant, and $\Psi(\epsilon, \tilde{\epsilon})$ is polynomial in ϵ^{-1} , $\tilde{\epsilon}^{-1}$, K , $\|\mu\|_\infty$, and Δ_{\min}^{-1} , where $\Delta_{\min} = \min_{x \neq i^*(\mu)} \langle i^*(\mu) - x, \mu \rangle$. Under P-FWS, the number of LM Oracle calls per round is at most polynomial in $\ln \delta^{-1}$ and K . The total expected number of these calls is also polynomial.

To the best of our knowledge, P-FWS is the first polynomial time best action identification algorithm with minimal sample complexity in the high confidence regime (when δ tends to 0). Its sample complexity is also polynomial in K in the moderate confidence regime.

2 Preliminaries

We start by introducing some notation. We use bold lowercase letters (e.g., x) for vectors, and bold uppercase letter (e.g., \mathbf{A}) for matrices. \odot (resp. \oplus) denotes the element-wise product (resp. sum over \mathbb{Z}_2). For $x \in \mathbb{R}^K$, $i \in \mathbb{N}$, $x^i = (x_k^i)_{k \in [K]}$ is the i -th element-wise power of x . $D = \max_{x \in \mathcal{X}} \|x\|_1$ denotes the maximum number of arms part of an action. For any $\mu \in \Lambda$, we define the sub-optimality gap of $x \in \mathcal{X}$ as $\Delta_x(\mu) = \langle i^*(\mu) - x, \mu \rangle$, and the minimal gap as $\Delta_{\min}(\mu) = \min_{x \neq i^*(\mu)} \Delta_x(\mu)$. \mathbb{P}_μ (resp. \mathbb{E}_μ) denotes the probability measure (resp. expectation) when the arm rewards are parametrized by μ . Whenever it is clear from the context, we will drop μ for simplicity, e.g. $i^* = i^*(\mu)$, $\Delta_x = \Delta_x(\mu)$, and $\Delta_{\min} = \min_{x \neq i^*} \Delta_x$. Refer to Appendix A for an exhaustive table of notation.

2.1 The lower-bound problem

Classical change-of-measure arguments lead to the asymptotic sample complexity lower bound $\mathbb{E}_\mu[\tau] \geq T^*(\mu) \text{kl}(\delta, 1 - \delta)$ where the characteristic time is defined in (1). To have a chance to develop a computationally efficient best action identification algorithm, we need that the sample complexity lower bound grows at most polynomially in K . This is indeed the case as stated in the following lemma, whose proof is provided in Appendix K.

Lemma 1. *For any $\mu \in \Lambda$, $T^*(\mu) \leq 4KD\Delta_{\min}(\mu)^{-2}$.*

We will use first-order methods to solve the lower-bound problem, and to this aim, we will need to evaluate the gradient w.r.t. ω of $f_x(\omega, \mu)$. We can apply the envelop theorem [WTP21] to show that for $(\omega, \mu) \in \Sigma_+ \times \Lambda$,

$$\nabla_\omega f_x(\omega, \mu) = \frac{(\mu - \lambda_{\omega, \mu}^*(x))^2}{2},$$

where $\Sigma_+ = \Sigma \cap \mathbb{R}_{>0}^K$, $\lambda_{\omega, \mu}^*(x) = \operatorname{argmin}_{\lambda \in \text{cl}(\mathcal{C}_x)} \langle \omega, \frac{(\mu - \lambda)^2}{2} \rangle$ and $\text{cl}(\mathcal{C}_x)$ is the closure of \mathcal{C}_x (refer to Lemma 19 in Appendix G.2).

2.2 The Linear Maximization Oracle

As mentioned earlier, we assume that we have access to a computationally efficient Oracle, referred to as the LM (Linear Maximization) Oracle, identifying the best action if μ is known. More precisely, as in existing works in combinatorial semi-bandits [KWA⁺14, PPV19, PBVP20], we make the following assumption.

Assumption 1. (i) *There exists a polynomial-time algorithm identifying $i^*(v)$ for any $v \in \mathbb{R}^K$; (ii) \mathcal{X} is inclusion-wise maximal, i.e., there is no $x, x' \in \mathcal{X}$ s.t. $x < x'$; (iii) for each $k \in [K]$, there exists $x \in \mathcal{X}$ such that $x_k = 1$; (iv) $|\mathcal{X}| \geq 2$.*

Assumption 1 holds for combinatorial sets including m -sets, spanning forests, bipartite matching, s - t paths. For completeness, we verify the assumption for these action sets in Appendix J. In the design of our MCP algorithm, we will actually need to solve for some $v \in \mathbb{R}^K$ the linear maximization problem $\max_{x \in \mathcal{X}} \langle x, v \rangle$ over $\mathcal{X} \setminus \{i^*(\mu)\}$; in other words, we will probably need to identify the second best action. Fortunately, this can be done in a computationally efficient manner under Assumption 1. The following lemma formalizes this observation. Its proof, presented in Appendix J, is inspired by Lawler's m -best algorithm [Law72].

Lemma 2. *Let $v \in \mathbb{R}^K$ and $x \in \mathcal{X}$. Under Assumption 1, there exists an algorithm that solves $\max_{x' \in \mathcal{X}: x' \neq x} \langle v, x' \rangle$ by only making at most D queries to the LM Oracle.*

3 Solving the lower bound problem: the MCP algorithm

Solving the lower bound problem first requires to evaluate its objective function $F_\mu(\omega)$. A naive approach, enumerating $f_x(\omega, \mu)$ for all $x \in \mathcal{X} \setminus \{i^*\}$, would be computationally infeasible. In this section, we present and analyze MCP, an algorithm that approximates $F_\mu(\omega)$ by calling the LM Oracle a number of times growing at most polynomially in K .

3.1 Lagrangian relaxation

The first step towards the design of MCP consists in considering the Lagrangian relaxation of the optimization problem defining $f_x(\omega, \mu) = \inf_{\lambda \in \mathcal{C}_x} \langle \omega, \frac{(\mu - \lambda)^2}{2} \rangle$ (see e.g., [BV04, Vis21]). For any $(\omega, \mu) \in \Sigma_+ \times \Lambda$ and $x \neq i^*$, the Lagrangian $\mathcal{L}_{\omega, \mu}$ and Lagrange dual function $g_{\omega, \mu}$ of this problem are defined as, $\forall \alpha \geq 0$,

$$\mathcal{L}_{\omega, \mu}(\lambda, x, \alpha) = \left\langle \omega, \frac{(\mu - \lambda)^2}{2} \right\rangle + \alpha \langle i^* - x, \lambda \rangle \quad \text{and} \quad g_{\omega, \mu}(x, \alpha) = \inf_{\lambda \in \mathbb{R}^K} \mathcal{L}_{\omega, \mu}(\lambda, x, \alpha),$$

respectively. The following proposition, proved in Appendix C.1, provides useful properties of $g_{\omega, \mu}$:

Proposition 1. *Let $(\omega, \mu) \in \Sigma_+ \times \Lambda$ and $x \in \mathcal{X} \setminus \{i^*(\mu)\}$.*

(a) *The Lagrange dual function is linear in x . More precisely, $g_{\omega, \mu}(x, \alpha) = c_{\omega, \mu}(\alpha) + \langle \ell_{\omega, \mu}(\alpha), x \rangle$*

- where $c_{\omega, \mu}(\alpha) = \alpha \langle \boldsymbol{\mu} - \frac{\alpha}{2} \boldsymbol{\omega}^{-1}, \mathbf{i}^*(\boldsymbol{\mu}) \rangle$ and $\ell_{\omega, \mu}(\alpha) = -\alpha (\boldsymbol{\mu} + \frac{\alpha}{2} \boldsymbol{\omega}^{-1} \odot (\mathbf{1}_K - 2\mathbf{i}^*(\boldsymbol{\mu})))$.
- (b) $g_{\omega, \mu}(\mathbf{x}, \cdot)$ is strictly concave (for any fixed \mathbf{x}).
- (c) $f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\mu}) = \max_{\alpha \geq 0} g_{\omega, \mu}(\mathbf{x}, \alpha)$ is attained by $\alpha_{\mathbf{x}}^* = \frac{\Delta_{\mathbf{x}}(\boldsymbol{\mu})}{\langle \mathbf{x} \oplus \mathbf{i}^*(\boldsymbol{\mu}), \boldsymbol{\omega}^{-1} \rangle}$.
- (d) $\|\ell_{\omega, \mu}(\alpha_{\mathbf{x}}^*)\|_1 \leq L_{\omega, \mu} = 4D^2K \|\boldsymbol{\mu}\|_{\infty}^2 \|\boldsymbol{\omega}^{-1}\|_{\infty}$.

From Proposition 1 (c), strong duality holds for the program defining $f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\mu})$, and we conclude:

$$F_{\boldsymbol{\mu}}(\boldsymbol{\omega}) = \min_{\mathbf{x} \neq \mathbf{i}^*} \max_{\alpha \geq 0} g_{\omega, \mu}(\mathbf{x}, \alpha). \quad (3)$$

$F_{\boldsymbol{\mu}}(\boldsymbol{\omega})$ can hence be seen as the value in a two-player game. The aforementioned properties of the Lagrange dual function will help to compute this value.

3.2 Solving the two-player game with no regret

There is a rich and growing literature on solving zero-sum games using no-regret algorithms, see for example [RS13, ALLW18, DFG21, ZODS21]. Our game has the particularity that the \mathbf{x} -player has a discrete combinatorial action set whereas the α -player has a convex action set. Importantly, for this game, we wish not only to estimate its value $F_{\boldsymbol{\mu}}(\boldsymbol{\omega})$ but also an *equilibrium* action \mathbf{x}_e such that $F_{\boldsymbol{\mu}}(\boldsymbol{\omega}) = \max_{\alpha \geq 0} g_{\omega, \mu}(\mathbf{x}_e, \alpha)$. Indeed, an estimate of \mathbf{x}_e will be needed when implementing the Frank-Wolfe algorithm and more specifically when estimating the gradient of $F_{\boldsymbol{\mu}}(\boldsymbol{\omega})$. To return such an estimate, one could think of leveraging results from the recent literature on last-iterate convergence, see e.g. [DP19, GPDO20, LNP⁺21, WLZL21, APFS22, AAS⁺23]. However, most of these results concern saddle-point problems only, and are not applicable in our setting. Here, we adopt a much simpler solution, and take advantage of the properties of the Lagrange dual function $g_{\omega, \mu}(\mathbf{x}, \alpha)$ to design an iterative procedure directly leading to estimates of $(F_{\boldsymbol{\mu}}(\boldsymbol{\omega}), \mathbf{x}_e)$. In this procedure, the two players successively update their actions until a stopping criterion is met, say up to the N -th iterations. The procedure generates a sequence $\{(\mathbf{x}^{(n)}, \alpha^{(n)})\}_{1 \leq n \leq N}$, and from this sequence, estimates $(\hat{F}, \hat{\mathbf{x}})$ of $(F_{\boldsymbol{\mu}}(\boldsymbol{\omega}), \mathbf{x}_e)$. The details of the resulting MCP algorithm are presented in Algorithm 1.

\mathbf{x} -player. We use a variant of the Follow-the-Perturbed-Leader (FTPL) algorithm [Han57, KV05]. The \mathbf{x} -player updates her action as follows:

$$\mathbf{x}^{(n)} \in \operatorname{argmin}_{\mathbf{x} \neq \mathbf{i}^*} \left(\sum_{m=1}^{n-1} g_{\omega, \mu}(\mathbf{x}, \alpha^{(m)}) + \left\langle \frac{\mathbf{Z}_n}{\eta_n}, \mathbf{x} \right\rangle \right) = \operatorname{argmin}_{\mathbf{x} \neq \mathbf{i}^*} \left(\left\langle \sum_{m=1}^{n-1} \ell_{\omega, \mu}(\alpha^{(m)}) + \frac{\mathbf{Z}_n}{\eta_n}, \mathbf{x} \right\rangle \right),$$

where \mathbf{Z}_n is a random vector, exponentially distributed and with unit mean ($\{\mathbf{Z}_n\}_{n \geq 1}$ are i.i.d.). Compared to the standard FTPL algorithm, we vary learning rate η_n over time to get *anytime* guarantees (as we do not know a priori when the iterative procedure will stop). This kind of time-varying learning rate was also used in [Neu15] with a similar motivation. Note that thanks to the linearity of $g_{\omega, \mu}$ and Lemma 2, the \mathbf{x} -player update can be computed using at most D calls to the LM Oracle.

α -player and MCP outputs. From Proposition 1, $f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\mu}) = \max_{\alpha \geq 0} g_{\omega, \mu}(\mathbf{x}, \alpha)$. This suggests that the α -player can just implement a best-response strategy: after the \mathbf{x} -player action $\mathbf{x}^{(n)}$ is selected, the α -player chooses $\alpha^{(n)} = \alpha_{\mathbf{x}^{(n)}}^* = \frac{\Delta_{\mathbf{x}^{(n)}}(\boldsymbol{\mu})}{\langle \mathbf{x}^{(n)} \oplus \mathbf{i}^*(\boldsymbol{\mu}), \boldsymbol{\omega}^{-1} \rangle}$. This choice ensures that $f_{\mathbf{x}^{(n)}}(\boldsymbol{\omega}, \boldsymbol{\mu}) = g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)})$, and suggests natural outputs for MCP: should it stops after N iterations, it can return $\hat{F} = \min_{n \in [N]} g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)})$ and $\hat{\mathbf{x}} \in \operatorname{argmin}_{n \in [N]} g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)})$.

Stopping criterion. The design of the MCP stopping criterion relies on the convergence analysis and regret from the \mathbf{x} -player perspective of the above iterative procedure, which we present in the next subsection. This convergence will be controlled by $\ell_{\omega, \mu}(\alpha_{\mathbf{x}}^*)$ and its upper bound $L_{\omega, \mu}$ derived in Proposition 1. Introducing $c_{\theta} = L_{\omega, \mu}(4\sqrt{K}(\ln K + 1) + \sqrt{\ln(\theta^{-1})/2})$, the MCP stopping criterion is: $\sqrt{n} > c_{\theta}(1 + \epsilon)/(\epsilon \hat{F})$. Since \sqrt{n} strictly increases with n and since $\hat{F} \geq F_{\boldsymbol{\mu}}(\boldsymbol{\omega})$, this criterion ensures that the algorithm terminates in a finite number of iterates. Moreover, as shown in the next subsection, it also ensures that \hat{F} returned by MCP is an $(1 + \epsilon)$ -approximation of $F_{\boldsymbol{\mu}}(\boldsymbol{\omega})$ with probability at least $1 - \theta$.

Algorithm 1: (ϵ, θ) -MCP(ω, μ)

initialization: $n = 1, \hat{F} = \infty, c_\theta = L_{\omega, \mu} \left(4\sqrt{K(\ln K + 1)} + \sqrt{\ln(\theta^{-1})/2} \right)$;
while $(n = 1)$ **or** $(n > 1)$ **and** $\sqrt{n} \leq c_\theta(1 + \epsilon)/(\epsilon\hat{F})$ **do**
 Sample $\mathcal{Z}_n \sim \exp(1)^K$ and set $\eta_n = \sqrt{K(\ln K + 1)/(4nL_{\omega, \mu}^2)}$;
 $\mathbf{x}^{(n)} \leftarrow \operatorname{argmin}_{\mathbf{x} \neq \mathbf{i}^*(\mu)} \left(\sum_{m=1}^{n-1} g_{\omega, \mu}(\mathbf{x}, \alpha^{(m)}) + \langle \mathcal{Z}_n, \mathbf{x} \rangle / \eta_n \right)$ (ties broken arbitrarily);
 $\alpha^{(n)} \leftarrow \operatorname{argmax}_{\alpha \geq 0} g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha)$ (uniqueness ensured by Proposition 1 (c));
 if $g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)}) < \hat{F}$ **then** $(\hat{F}, \hat{\mathbf{x}}) \leftarrow (g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)}), \mathbf{x}^{(n)})$;
 $n \leftarrow n + 1$;
end
return $(\hat{F}, \hat{\mathbf{x}})$;

3.3 Performance analysis of the MCP algorithm

We start the analysis by quantifying the regret from the \mathbf{x} -player perspective of MCP before its stops. The following lemma is proved in Appendix C.3.

Lemma 3. Let $N \in \mathbb{N}$. Under (ϵ, θ) -MCP(ω, μ),

$$\mathbb{P} \left[\frac{1}{N} \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)}) - \frac{1}{N} \min_{\mathbf{x} \neq \mathbf{i}^*} \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}, \alpha^{(n)}) \leq \frac{c_\theta}{\sqrt{N}} \right] \geq 1 - \theta.$$

Observe that on the one hand,

$$\frac{1}{N} \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)}) \geq \min_{n \in [N]} g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)}) = \hat{F} \quad (4)$$

always holds. On the other hand, if $\mathbf{x}_e \in \operatorname{argmin}_{\mathbf{x} \neq \mathbf{i}^*} \max_{\alpha \geq 0} g_{\omega, \mu}(\mathbf{x}, \alpha)$, then we have:

$$\frac{1}{N} \min_{\mathbf{x} \neq \mathbf{i}^*} \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}, \alpha^{(n)}) \leq \frac{1}{N} \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}_e, \alpha^{(n)}) \leq \max_{\alpha \geq 0} g_{\omega, \mu}(\mathbf{x}_e, \alpha) = F_\mu(\omega). \quad (5)$$

We conclude that for N such that $\sqrt{N} \geq \frac{c_\theta(1+\epsilon)}{\epsilon\hat{F}}$, Lemma 3 together with the inequalities (4) and (5) imply that $\hat{F} - F_\mu(\omega) \leq \frac{c_\theta}{\sqrt{N}} \leq \frac{\epsilon\hat{F}}{1+\epsilon}$ holds with probability at least $1 - \theta$. Hence $\mathbb{P} \left[\hat{F} \leq (1 + \epsilon)F_\mu(\omega) \right] \geq 1 - \theta$. From this observation, we essentially deduce the following theorem, whose complete proof is given in Appendix C.2.

Theorem 3. Let $\epsilon, \theta \in (0, 1)$. Under Assumption 1, for any $(\omega, \mu) \in \Sigma_+ \times \Lambda$, the (ϵ, θ) -MCP(ω, μ) algorithm outputs $(\hat{F}, \hat{\mathbf{x}})$ satisfying $\mathbb{P} \left[F_\mu(\omega) \leq \hat{F} \leq (1 + \epsilon)F_\mu(\omega) \right] \geq 1 - \theta$ and $\hat{F} = \max_{\alpha \geq 0} g_{\omega, \mu}(\hat{\mathbf{x}}, \alpha)$. Moreover, the number of LM Oracle calls the algorithm does is almost surely at most $\left\lceil \frac{c_\theta^2(1+\epsilon)^2}{\epsilon^2 F_\mu(\omega)^2} \right\rceil = \mathcal{O} \left(\frac{\|\mu\|_\infty^4 \|\omega^{-1}\|_\infty^2 K^3 D^5 \ln K \ln \theta^{-1}}{\epsilon^2 F_\mu(\omega)^2} \right)$.

4 The Perturbed Frank-Wolfe Sampling (P-FWS) algorithm

To identify an optimal sampling strategy, rather than solving the lower-bound problem in each round as a Track-and-Stop algorithm would [GK16], we devise P-FWS, an algorithm that performs a single iteration of the Frank-Wolfe algorithm for the lower-bound problem instantiated with an estimator of μ . This requires us to first smooth the objective function $F_\mu(\omega) = \min_{\mathbf{x} \neq \mathbf{i}^*} f_\mathbf{x}(\omega, \mu)$ (the latter is not differentiable at points ω where the min is achieved for several sub-optimal actions \mathbf{x}). To this aim, we cannot leverage the same technique as in [WTP21], where r -subdifferential subspaces are built from gradients of $f_\mathbf{x}(\omega, \mu)$. These subspaces could indeed be generated by a number of vectors (here gradients) exponentially growing with K . Instead, to cope with the combinatorial

decision sets, P-FWS applies more standard stochastic smoothing techniques as described in the next subsection. All the ingredients of P-FWS are gathered in §4.2. By design, the algorithm just leverages the MCP algorithm as a subroutine, and hence only requires the LM Oracle. In §4.3, we analyze the performance of P-FWS.

4.1 Smoothing the objective function F_μ

Here, we present and analyze a standard stochastic technique to smooth a function Φ . In P-FWS, this technique will be applied to the objective function $\Phi = F_\mu$. Let $\Phi : \mathbb{R}_{>0}^K \mapsto \mathbb{R}$ be a concave and ℓ -Lipschitz function. Assume that the set of points where Φ is not differentiable is of Lebesgue-measure zero. To smooth Φ , we can take its average value in a neighborhood of the point considered, see e.g. [FKM05]. Formally, we define the *stochastic smoothed* approximate of Φ as:

$$\bar{\Phi}_\eta(\omega) = \mathbb{E}_{\mathcal{Z} \sim \text{Uniform}(B_2)}[\Phi(\omega + \eta\mathcal{Z})], \quad (6)$$

where $B_2 = \{v \in \mathbb{R}^K : \|v\|_2 \leq 1\}$ and $\eta \in (0, \min_{k \in [K]} \omega_k)$. The following proposition lists several properties of this smoothed function, and gathers together some of the results from [DBW12], see Appendix H for more details.

Proposition 2. *For any $\omega \in \Sigma_+$ and $\eta \in (0, \min_{k \in [K]} \omega_k)$, $\bar{\Phi}_\eta$ satisfies: (i) $\Phi(\omega) - \eta\ell \leq \bar{\Phi}_\eta(\omega) \leq \Phi(\omega)$; (ii) $\nabla \bar{\Phi}_{\mu, \eta}(\omega) = \mathbb{E}_{\mathcal{Z} \sim \text{Uniform}(B_2)}[\nabla \Phi_\mu(\omega + \eta\mathcal{Z})]$; (iii) $\bar{\Phi}_\eta$ is $\frac{\ell K}{\eta}$ -smooth; (iv) if $\eta > \eta' > 0$, then $\bar{\Phi}_{\eta'}(\omega) \geq \bar{\Phi}_\eta(\omega)$.*

Note that with (i), we may control the approximation error between $\bar{\Phi}_\eta$ and Φ by η . (ii) and (iii) ensure the differentiability and smoothness of $\bar{\Phi}_\eta$ respectively. (iii) is equivalent to $\bar{\Phi}_\eta(\omega') \leq \bar{\Phi}_\eta(\omega) + \langle \nabla \bar{\Phi}_\eta(\omega), \omega' - \omega \rangle + \frac{\ell K}{2\eta} \|\omega - \omega'\|_2^2$, $\forall \omega, \omega' \in \Sigma_+$. Finally, (iv) stems from the concavity of Φ , and implies that the value $\bar{\Phi}_\eta(\omega)$ monotonously increases while η decreases, and it is upper bounded by $\Phi(\omega)$ thanks to (i). The above results hold for $\Phi = F_\mu$. Indeed, first it is clear that the definition (2) of F_μ can be extended to \mathbb{R}^K ; then, it can be shown that F_μ is Lipschitz-continuous and almost-everywhere differentiable – refer to Appendices I and H for formal proofs.

4.2 The algorithm

Before presenting P-FWS, we need to introduce the following notation. For $t \geq 1$, $k \in [K]$, we define $N_k(t) = \sum_{s=1}^t \mathbb{1}\{x_k(s) = 1\}$, $\hat{\omega}_k(t) = N_k(t)/t$, and $\hat{\mu}_k(t) = \sum_{s=1}^t y_k(s) \mathbb{1}\{x_k(s) = 1\} / N_k(t)$ when $N_k(t) > 0$.

Sampling rule. The design of the sampling rule is driven by the following objectives: (i) the empirical allocation should converge to the solution of the lower-bound problem (1), and (ii) the number of calls to the LM Oracle should be controlled. To meet the first objective, we need in the Frank-Wolfe updates to plug an accurate estimator of μ in. The accuracy of our estimator will be guaranteed by alternating between *forced exploration* and *FW update* sampling phases. Now for the second objective, we also use forced exploration phases when in a Frank-Wolfe update, the required number of calls to the LM Oracle predicted by the upper bound presented in Theorem 3 is too large. In view of Lemma 1 and Theorem 3, this happens in round t if $\|\hat{\mu}(t-1)\|_\infty$ or $\Delta_{\min}(\hat{\mu}(t-1))^{-1}$ is too large. Next, we describe the forced exploration and Frank-Wolfe update phases in detail.

Forced exploration. Initially, P-FWS applies the LM Oracle to compute the *forced exploration set* $\mathcal{X}_0 = \{i^*(e_k) : k \in [K]\}$, where e_k is the K -dimensional vector whose k -th component is equal to one and zero elsewhere. P-FWS then selects each action in \mathcal{X}_0 once. Note that Assumption 1 (iii) ensures that the k -th component of $i^*(e_k)$ is equal to one. In turn, this ensures that \mathcal{X}_0 is a $[K]$ -covering set, and that playing actions from \mathcal{X}_0 is enough to estimate μ . P-FWS starts an exploration phase at rounds t such that $\sqrt{t/|\mathcal{X}_0|}$ is an integer or such that the maximum of $\Delta_{\min}(\hat{\mu}(t-1))^{-1}$ and $\|\hat{\mu}(t-1)\|_\infty$ is larger than $\sqrt{t-1}$. Whenever this happens, P-FWS pulls each $x \in \mathcal{X}_0$ once.

Frank-Wolfe updates. When in round t , the algorithm is not in a forced exploration phase, it implements an iteration of the Frank-Wolfe algorithm applied to maximize the smoothed function $\bar{F}_{\hat{\mu}(t-1), \eta_t}(\hat{\omega}(t-1)) = \mathbb{E}_{\mathcal{Z} \sim \text{Uniform}(B_2)}[F_{\hat{\mu}(t-1)}(\hat{\omega}(t-1) + \eta_t\mathcal{Z})]$. The sequence of parameters $\{\eta_t\}_{t \geq 1}$ is chosen to ensure that η_t chosen in $(0, \min_k \hat{\omega}_k(t))$, and hence $\hat{\omega}(t-1) + \eta_t\mathcal{Z} \in \mathbb{R}_{>0}^K$. Also note that in a round t where the algorithm is not in a forced exploration phase,

Algorithm 2: P-FWS $(\{(\epsilon_t, \eta_t, n_t, \rho_t, \theta_t)\}_t)$

initialization:

```
for  $k = 1, \dots, K$  do
  |  $\mathcal{X}_0 \leftarrow \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{e}_k, \mathbf{x} \rangle$  (tie broken arbitrarily)
end
Sample  $\mathbf{x} \in \mathcal{X}_0$  in a round-robin manner for  $4|\mathcal{X}_0|$  rounds; update  $\hat{\boldsymbol{\mu}}(4|\mathcal{X}_0|)$  and  $\hat{\boldsymbol{\omega}}(4|\mathcal{X}_0|)$ ;
for  $t = 4|\mathcal{X}_0| + 1, \dots$  do
  if  $\sqrt{t/|\mathcal{X}_0|} \in \mathbb{N}$  or  $\max\{\Delta_{\min}(\hat{\boldsymbol{\mu}}(t-1))^{-1}, \|\hat{\boldsymbol{\mu}}(t-1)\|_{\infty}\} > \sqrt{t-1}$  then
    | Sample each  $\mathbf{x} \in \mathcal{X}_0$  once, update  $\hat{\boldsymbol{\mu}}(t)$  and  $\hat{\boldsymbol{\omega}}(t)$ , and  $t \leftarrow t + |\mathcal{X}_0| - 1$ ;
  else
    | Compute  $\nabla \tilde{F}_{\hat{\boldsymbol{\mu}}(t-1), \eta_t, n_t}(\hat{\boldsymbol{\omega}}(t-1))$  by  $(\rho_t, \theta_t)$ -MCP algorithm;
    |  $\mathbf{x}(t) \leftarrow \mathbf{i}^*\left(\nabla \tilde{F}_{\hat{\boldsymbol{\mu}}(t-1), \eta_t, n_t}(\hat{\boldsymbol{\omega}}(t-1))\right)$ ;
    | Sample  $\mathbf{x}(t)$  and update  $\hat{\boldsymbol{\mu}}(t)$  and  $\hat{\boldsymbol{\omega}}(t)$ ;
  end
  if  $\max\{\Delta_{\min}(\hat{\boldsymbol{\mu}}(t))^{-1}, \|\hat{\boldsymbol{\mu}}(t)\|_{\infty}\} \leq \sqrt{t}$  then
    |  $\hat{F}_t \leftarrow (\epsilon_t, \delta/t^2)$ -MCP  $(\hat{\boldsymbol{\omega}}(t), \hat{\boldsymbol{\mu}}(t))$ ;
    | if  $t\hat{F}_t > (1 + \epsilon_t)\beta\left(t, \left(1 - \frac{1}{4|\mathcal{X}_0|}\right)\delta\right)$  then break;
  end
return  $\hat{\mathbf{i}} = \mathbf{i}^*(\hat{\boldsymbol{\mu}}(t))$ ;
```

by definition $\Delta_{\min}(\hat{\boldsymbol{\mu}}(t-1)) > 0$. This implies that $\hat{\boldsymbol{\mu}}(t-1) \in \Lambda$ and that $F_{\hat{\boldsymbol{\mu}}(t-1)}$ and $\tilde{F}_{\hat{\boldsymbol{\mu}}(t-1), \eta_t}(\hat{\boldsymbol{\omega}}(t-1))$ are well-defined. Now an ideal FW update would consist in playing an action $\mathbf{i}^*(\nabla \tilde{F}_{\hat{\boldsymbol{\mu}}(t-1), \eta_t}(\hat{\boldsymbol{\omega}}(t-1))) = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \langle \nabla \tilde{F}_{\hat{\boldsymbol{\mu}}(t-1), \eta_t}(\hat{\boldsymbol{\omega}}(t-1)), \mathbf{x} \rangle$, see e.g. [Jag13]. Unfortunately, we do not have access to $\nabla \tilde{F}_{\hat{\boldsymbol{\mu}}(t-1), \eta_t}(\hat{\boldsymbol{\omega}}(t-1))$. But the latter can be approximated, as suggested in Proposition 2 (ii), by $\nabla \tilde{F}_{\hat{\boldsymbol{\mu}}(t-1), \eta_t, n_t}(\hat{\boldsymbol{\omega}}(t-1)) = \frac{1}{n_t} \sum_{m=1}^{n_t} \nabla f_{\hat{\mathbf{x}}_m}(\hat{\boldsymbol{\omega}}(t-1) + \eta_t \mathcal{Z}_m, \hat{\boldsymbol{\mu}}(t-1))$, where $\mathcal{Z}_1, \dots, \mathcal{Z}_{n_t} \stackrel{i.i.d.}{\sim} \operatorname{Uniform}(B_2)$, $\hat{\mathbf{x}}_m$ is the action return by (ρ_t, θ_t) -MCP $(\hat{\boldsymbol{\omega}}(t-1) + \eta_t \mathcal{Z}_m, \hat{\boldsymbol{\mu}}(t-1))$. P-FWS uses this approximation and the LM Oracle to select the action: $\mathbf{x}(t) \in \mathbf{i}^*\left(\nabla \tilde{F}_{\hat{\boldsymbol{\mu}}(t-1), \eta_t, n_t}(\hat{\boldsymbol{\omega}}(t-1))\right)$. The choices of the parameters η_t, n_t, ρ_t and θ_t do matter. η_t impacts the sample complexity and should converge to 0 as $t \rightarrow \infty$ so that $\tilde{F}_{\hat{\boldsymbol{\mu}}, \eta_t}(\boldsymbol{\omega}) \rightarrow F_{\hat{\boldsymbol{\mu}}}(\boldsymbol{\omega})$ at any point $\boldsymbol{\omega} \in \Sigma_+$ (this is a consequence of Proposition 2 (i)(iv)). η_t should not decay too fast however as it would alter the smoothness of $\tilde{F}_{\hat{\boldsymbol{\mu}}, \eta_t}$. We will show that η_t should actually decay as $1/\sqrt{t}$. (n_t, ρ_t, θ_t) impact the trade-off between the sample complexity and the computational complexity of the algorithm. We let $n_t \rightarrow \infty$ and $(\rho_t, \theta_t) \rightarrow 0$ as $t \rightarrow \infty$ so that $\langle \nabla \tilde{F}_{\hat{\boldsymbol{\mu}}, \eta_t, n_t}(\boldsymbol{\omega}) - \nabla \tilde{F}_{\hat{\boldsymbol{\mu}}, \eta_t}(\boldsymbol{\omega}), \mathbf{x} \rangle \rightarrow 0$ for any $(\boldsymbol{\omega}, \mathbf{x}) \in \Sigma_+ \times \mathcal{X}$.

Stopping and decision rule. As often in best arm identification algorithms, the P-FWS stopping rule takes the form of a GLRT:

$$\tau = \inf \left\{ t > 4|\mathcal{X}_0| : \frac{t\hat{F}_t}{1 + \epsilon_t} > \beta\left(t, \left(1 - \frac{1}{4|\mathcal{X}_0|}\right)\delta\right), \max\{\Delta_{\min}(\hat{\boldsymbol{\mu}}(t))^{-1}, \|\hat{\boldsymbol{\mu}}(t)\|_{\infty}\} \leq \sqrt{t} \right\}, \quad (7)$$

where $\epsilon_t \in \mathbb{R}_{>0}$, \hat{F}_t is returned by the $(\epsilon_t, \delta/t^2)$ -MCP $(\hat{\boldsymbol{\omega}}(t), \hat{\boldsymbol{\mu}}(t))$ algorithm. The function β satisfies

$$\forall t \geq 1, \quad (tF_{\hat{\boldsymbol{\mu}}(t)}(\boldsymbol{\omega}(t)) \geq \beta(t, \delta)) \implies (\mathbb{P}_{\boldsymbol{\mu}}[\mathbf{i}^*(\hat{\boldsymbol{\mu}}(t)) \neq \mathbf{i}^*(\boldsymbol{\mu})] \leq \delta), \quad (8)$$

$$\exists c_1, c_2 > 0 : \quad \forall t \geq c_1, \beta(t, \delta) \leq \ln\left(\frac{c_2 t}{\delta}\right). \quad (9)$$

Examples of function β satisfying the above conditions can be found in [GK16, JP20, KK21]. The condition (8) will ensure the δ -correctness of P-FWS, whereas (9) will control its sample complexity. Finally, the action returned by P-FWS is simply defined as $\hat{\mathbf{i}} = \mathbf{i}^*(\hat{\boldsymbol{\mu}}(\tau))$. The complete pseudo-code of P-FWS is presented in Algorithm 2.³

³Our Julia implementation could be found at <https://github.com/rctzeng/NeurIPS2023-PerturbedFWS>.

4.3 Non-asymptotic performance analysis of P-FWS

The following theorem provides an upper bound of the sample complexity of P-FWS valid for any confidence level δ , as well as the computational complexity of the algorithm.

Theorem 4. *Let $\mu \in \Lambda$ and $\delta \in (0, 1)$. If P-FWS is parametrized using*

$$(\epsilon_t, \eta_t, n_t, \rho_t, \theta_t) = \left(t^{-\frac{1}{5}}, \frac{1}{4\sqrt{t|\mathcal{X}_0|}}, \lceil t^{\frac{1}{4}} \rceil, \frac{1}{16tD^2|\mathcal{X}_0|}, \frac{1}{t^{\frac{1}{4}}e^{\sqrt{t}}} \right), \quad (10)$$

then (i) the algorithm finishes in finite time almost surely and $\mathbb{P}_\mu[\hat{i} \neq i^*(\mu)] \leq \delta$; (ii) its sample complexity satisfies $\mathbb{P}_\mu[\limsup_{\delta \rightarrow 0} \frac{\tau}{\ln \frac{1}{\delta}} \leq T^*(\mu)] = 1$ and for any $\epsilon, \tilde{\epsilon} \in (0, 1)$ with $\epsilon < \min\{1, \frac{2D^2\Delta_{\min}^2}{K}, \frac{D^2\|\mu\|_\infty^2}{3}\}$,

$$\mathbb{E}_\mu[\tau] \leq \frac{(1 + \tilde{\epsilon})^2}{T^*(\mu)^{-1} - 6\epsilon} \times H\left(\frac{1}{\delta} \cdot \frac{4c_2}{3} \cdot \frac{(1 + \tilde{\epsilon})^2}{T^*(\mu)^{-1} - 6\epsilon}\right) + \Psi(\epsilon, \tilde{\epsilon}),$$

where $H(x) = \ln x + \ln \ln x + 1$ and $\Psi(\epsilon, \tilde{\epsilon})$ (refer to (34) for a detailed expression) is polynomial in ϵ^{-1} , $\tilde{\epsilon}^{-1}$, K , $\|\mu\|_\infty$, and $\Delta_{\min}(\mu)^{-1}$; (iii) the expected number of LM Oracle calls is upper bounded by a polynomial in $\ln \delta^{-1}$, K , $\|\mu\|_\infty$, and $\Delta_{\min}(\mu)^{-1}$.

The above theorem establishes the statistical asymptotic optimality of P-FWS since it implies that $\limsup_{\delta \rightarrow 0} \mathbb{E}_\mu[\tau] / \ln(1/\delta) \leq (1 + \tilde{\epsilon})^2 / (T^*(\mu)^{-1} - 6\epsilon)$. This upper bound matches the sample complexity lower bound (1) when $\epsilon \rightarrow 0$ and $\tilde{\epsilon} \rightarrow 0$.

Proof sketch. The complete proof of Theorem 4 is presented in Appendix D.

(i) *Correctness.* To establish the δ -correctness of the algorithm, we introduce the event \mathcal{G} under which \hat{F}_t , computed by $(\epsilon_t, \delta/t^2)$ -MCP $(\hat{\omega}(t), \hat{\mu}(t))$, is an $(1 + \epsilon_t)$ -approximation of $F_{\hat{\mu}(t)}(\hat{\omega}(t))$ in each round $t \geq 4|\mathcal{X}_0| + 1$. From Theorem 3, we deduce that $\mathbb{P}_\mu[\mathcal{G}^c] \leq \sum_{t=4|\mathcal{X}_0|+1}^\infty \delta/t^2 \leq \delta/4|\mathcal{X}_0|$. In view of (8), this implies that $\mathbb{P}_\mu[\hat{i} \neq i^*(\mu)] \leq \delta$.

(ii) *Non-asymptotic sample complexity upper bound.*

Step 1. (Concentration and certainty equivalence) We first define two good events, $\mathcal{E}_t^{(1)}$ and $\mathcal{E}_t^{(2)}$. $\mathcal{E}_t^{(2)}$ corresponds to the event where $\hat{\mu}(t)$ is close to μ , and its occurrence probability can be controlled using the forced exploration rounds and concentration inequalities. Under $\mathcal{E}_t^{(1)}$, the selected action $x(t)$ is close to the ideal FW update. Again using concentration results and the performance guarantees of MCP given in Theorem 3, we can control the occurrence probability of $\mathcal{E}_t^{(2)}$. Overall, we show that $\sum_{t=1}^\infty \mathbb{P}_\mu[(\mathcal{E}_t^{(1)} \cap \mathcal{E}_t^{(2)})^c] < \infty$. To this aim, we derive several important continuity results presented in Appendix G. These results essentially allow us to study the convergence of the smoothed FW updates as if the certainty equivalence principle held, i.e., as if $\hat{\mu}(t) = \mu$.

Step 2. (Convergence of the smoothed FW updates) We study the convergence assuming that $(\mathcal{E}_t^{(1)} \cap \mathcal{E}_t^{(2)})$ holds. We first show that \bar{F}_{μ, η_t} is ℓ -Lipschitz and smooth for $\ell = 2D^2\|\mu\|_\infty^2$, see Appendices H and I. Then, in Appendix E, we establish that the dynamics of $\phi_t = \max_{\omega \in \Sigma} \bar{F}_\mu(\omega) - F_\mu(\hat{\omega}(t))$ satisfy $t\phi_t \leq (t-1)\phi_{t-1} + \ell \left(\eta_{t-1} + \frac{K^2}{2t\eta_t} \right)$. Observe that, as mentioned earlier, $1/\sqrt{t}$ is indeed the optimal scaling choice for η_t . We deduce that after a certain finite number T_1 of rounds, ϕ_t is sufficiently small and $\max\{\Delta_{\min}(\hat{\mu}(t))^{-1}, \|\hat{\mu}(t)\|_\infty\} \leq \sqrt{t}$.

Step 3. Finally, we observe that $\mathbb{E}_\mu[\tau] \leq T_1 + \sum_{t=T_1}^\infty \mathbb{P}_\mu \left[t\hat{F}_t \leq (1 + \epsilon_t)\beta\left(t, \left(1 - \frac{1}{4|\mathcal{X}_0|}\right)\delta\right) \right] + \sum_{t=T_1+1}^\infty \mathbb{P}_\mu \left[(\mathcal{E}_t^{(1)} \cap \mathcal{E}_t^{(2)})^c \right]$, and show that the second term in the r.h.s. in this inequality is equivalent to $T^*(\mu) \ln(1/\delta)$ as $\delta \rightarrow 0$ using the property of the function β defining the stopping threshold and similar arguments as those used in [GK16, WTP21].

(iii) *Expected number of LM Oracle calls.* The MCP algorithm is called to compute \hat{F}_t and to perform the FW update only in rounds t such that $\max\{\Delta_{\min}(\hat{\mu}(t))^{-1}, \|\hat{\mu}(t)\|_\infty\} \leq \sqrt{t}$. Thus, from Theorem 3 and Lemma 1, the number of LM Oracle calls per-round is a polynomial in t and K . As the $\mathbb{E}_\mu[\tau]$ is polynomial (in $\ln \delta^{-1}$, K , $\|\mu\|_\infty$ and Δ_{\min}^{-1}), the expected number of LM Oracle calls is also polynomial in the same variables.

5 Related Work

We provide an exhaustive survey of the related literature in Appendix B. To summarize, to the best of our knowledge, CombGame [JMKK21] is the state-of-the-art algorithm for BAI in combinatorial semi-bandits in the high confidence regimes. A complete comparison to P -FWS is presented in Appendix B. CombGame was initially introduced in [DKM19] for classical bandit problems. There, the lower-bound problem is casted as a two-player game and the authors propose to use no-regret algorithms for each player to solve it. [JMKK21] adapts the algorithm for combinatorial semi-bandits, and provides a non-asymptotic sample complexity upper bound matching (1) asymptotically. However, the resulting algorithm requires to call an oracle solving the Most-Confusing-Parameter problem as our MCP algorithm. The authors of [JMKK21] conjectured the existence of such a computationally efficient oracle, and we establish this result here.

6 Conclusion

In this paper, we have presented P -FWS, the first computationally efficient and statistically optimal algorithm for the best arm identification problem in combinatorial semi-bandits. For this problem, we have studied the computational-statistical trade-off through the analysis of the optimization problem leading to instance-specific sample complexity lower bounds. This approach can be extended to study the computational-statistical gap in other learning tasks. Of particular interest are problems with an underlying structure (e.g. linear bandits [DMSV20, JP20], or RL in linear / low rank MDPs [AKKS20]). Most results on these problems are concerned with statistical efficiency, and ignore computational issues.

Acknowledgments and Disclosure of Funding

We thank Aristides Gionis and the anonymous reviewers for their valuable feedback. The research is funded by ERC Advanced Grant REBOUND (834862), the Wallenberg AI, Autonomous Systems and Software Program (WASP), and Digital Futures.

References

- [AAS⁺23] Kenshi Abe, Kaito Ariu, Mitsuki Sakamoto, Kentaro Toyoshima, and Atsushi Iwasaki. Last-iterate convergence with full-and noisy-information feedback in two-player zero-sum games. In *Proc. of AISTATS*, 2023.
- [AKKS20] Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. Flambe: Structural complexity and representation learning of low rank mdps. In *Proc. of NeurIPS*, 2020.
- [ALLW18] Jacob Abernethy, Kevin A Lai, Kfir Y Levy, and Jun-Kun Wang. Faster rates for convex-concave games. In *Proc. of COLT*, 2018.
- [AMP21] Aymen Al Marjani and Alexandre Proutiere. Adaptive sampling for best policy identification in markov decision processes. In *Proc. of ICML*, 2021.
- [APFS22] Ioannis Anagnostides, Ioannis Panageas, Gabriele Farina, and Tuomas Sandholm. On last-iterate convergence beyond zero-sum games. In *Proc. of ICML*, 2022.
- [BBGS11] André Berger, Vincenzo Bonifaci, Fabrizio Grandoni, and Guido Schäfer. Budgeted matching and budgeted matroid intersection via the gasoline puzzle. *Mathematical Programming*, 2011.
- [BGK22] Antoine Barrier, Aurélien Garivier, and Tomáš Kocák. A non-asymptotic approach to best-arm identification for gaussian bandits. In *Proc. of AISTATS*, 2022.
- [BLM13] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, 2013.

- [BV04] Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [CBL06] Nicolo Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006.
- [CBL12] Nicolo Cesa-Bianchi and Gábor Lugosi. Combinatorial bandits. *Journal of Computer and System Sciences*, 2012.
- [CCG21a] Thibaut Cuvelier, Richard Combes, and Eric Gourdin. Asymptotically optimal strategies for combinatorial semi-bandits in polynomial time. In *Proc. of ALT*, 2021.
- [CCG21b] Thibaut Cuvelier, Richard Combes, and Eric Gourdin. Statistically efficient, polynomial-time algorithms for combinatorial semi-bandits. *Proc. of SIGMETRICS*, 2021.
- [CGL16] Lijie Chen, Anupam Gupta, and Jian Li. Pure exploration of multi-armed bandit under matroid constraints. In *Proc. of COLT*, 2016.
- [CGL⁺17] Lijie Chen, Anupam Gupta, Jian Li, Mingda Qiao, and Ruosong Wang. Nearly optimal sampling algorithms for combinatorial pure exploration. In *Proc. of COLT*, 2017.
- [CLK⁺14] Shouyuan Chen, Tian Lin, Irwin King, Michael R Lyu, and Wei Chen. Combinatorial pure exploration of multi-armed bandits. In *Proc. of NeurIPS*, 2014.
- [CMP17] Richard Combes, Stefan Magureanu, and Alexandre Proutiere. Minimal exploration in structured stochastic bandits. In *Proc. of NeurIPS*, 2017.
- [CTMSP⁺15] Richard Combes, Mohammad Sadegh Talebi Mazraeh Shahi, Alexandre Proutiere, et al. Combinatorial bandits revisited. In *Proc. of NeurIPS*, 2015.
- [DBW12] John C Duchi, Peter L Bartlett, and Martin J Wainwright. Randomized smoothing for stochastic optimization. *SIAM Journal on Optimization*, 2012.
- [DFG21] Constantinos Daskalakis, Maxwell Fishelson, and Noah Golowich. Near-optimal no-regret learning in general games. In *Proc. of NeurIPS*, 2021.
- [DKC21] Yihan Du, Yuko Kuroki, and Wei Chen. Combinatorial pure exploration with full-bandit or partial linear feedback. In *Proc. of AAAI*, 2021.
- [DKM19] Rémy Degenne, Wouter M Koolen, and Pierre Ménard. Non-asymptotic pure exploration by solving games. In *Proc. of NeurIPS*, 2019.
- [DMSV20] Rémy Degenne, Pierre Ménard, Xuedong Shang, and Michal Valko. Gamification of pure exploration for linear bandits. In *Proc. of ICML*, 2020.
- [DP19] Constantinos Daskalakis and Ioannis Panageas. Last-iterate convergence: Zero-sum games and constrained min-max optimization. *Proc. of ITCS*, 2019.
- [FKM05] Abraham D Flaxman, Adam Tauman Kalai, and H Brendan McMahan. Online convex optimization in the bandit setting: gradient descent without a gradient. In *Proc. of SODA*, 2005.
- [FKV14] Eugene A Feinberg, Pavlo O Kasyanov, and Mark Voorneveld. Berge’s maximum theorem for noncompact image sets. *Journal of Mathematical Analysis and Applications*, 2014.
- [GC11] A. Garivier and O. Cappé. The KL-UCB algorithm for bounded stochastic bandits and beyond. In *Proc. of COLT*, 2011.
- [GH16] Dan Garber and Elad Hazan. A linearly convergent variant of the conditional gradient algorithm under strong convexity, with applications to online and stochastic optimization. *SIAM Journal on Optimization*, 2016.

- [GK16] Aurélien Garivier and Emilie Kaufmann. Optimal best arm identification with fixed confidence. In *Proc. of COLT*, 2016.
- [GPDO20] Noah Golowich, Sarath Pattathil, Constantinos Daskalakis, and Asuman Ozdaglar. Last iterate is slower than averaged iterate in smooth convex-concave saddle point problems. In *Proc. of COLT*, 2020.
- [H⁺16] Elad Hazan et al. Introduction to online convex optimization. *Foundations and Trends® in Optimization*, 2016.
- [Han57] James Hannan. Approximation to bayes risk in repeated play. *Contributions to the Theory of Games*, 1957.
- [HK12] Elad Hazan and Satyen Kale. Projection-free online learning. In *Proc. of ICML*, 2012.
- [Jag13] Martin Jaggi. Revisiting frank-wolfe: Projection-free sparse convex optimization. In *Proc. of ICML*, 2013.
- [JMKK21] Marc Jourdan, Mojmír Mutný, Johannes Kirschner, and Andreas Krause. Efficient pure exploration for combinatorial bandits with semi-bandit feedback. In *Proc. of ALT*, 2021.
- [JP20] Yassir Jedra and Alexandre Proutiere. Optimal best-arm identification in linear bandits. In *Proc. of NeurIPS*, 2020.
- [KCG16] Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best-arm identification in multi-armed bandit models. *JMLR*, 2016.
- [KK21] Emilie Kaufmann and Wouter M Koolen. Mixture martingales revisited with applications to sequential tests and confidence intervals. *JMLR*, 2021.
- [KLLM22] Daniel Kane, Sihan Liu, Shachar Lovett, and Gaurav Mahajan. Computational-statistical gap in reinforcement learning. In *Proc. of COLT*, 2022.
- [KSJJ⁺20] Julian Katz-Samuels, Lalit Jain, Kevin G Jamieson, et al. An empirical process approach to the union bound: Practical algorithms for combinatorial and linear bandits. In *Proc. of NeurIPS*, 2020.
- [KTAS12] Shivaram Kalyanakrishnan, Ambuj Tewari, Peter Auer, and Peter Stone. Pac subset selection in stochastic multi-armed bandits. In *Proc. of ICML*, 2012.
- [KV05] Adam Kalai and Santosh Vempala. Efficient algorithms for online decision problems. *Journal of Computer and System Sciences*, 2005.
- [KWA⁺14] Branislav Kveton, Zheng Wen, Azin Ashkan, Hoda Eydgahi, and Brian Eriksson. Matroid bandits: fast combinatorial optimization with learning. In *Proc. of UAI*, 2014.
- [Lai87] Tze Leung Lai. Adaptive treatment allocation and the multi-armed bandit problem. *The annals of statistics*, 1987.
- [Law72] Eugene L Lawler. A procedure for computing the k best solutions to discrete optimization problems and its application to the shortest path problem. *Management science*, 1972.
- [LNP⁺21] Qi Lei, Sai Ganesh Nagarajan, Ioannis Panageas, et al. Last iterate convergence in no-regret learning: constrained min-max optimization for convex-concave landscapes. In *Proc. of AISTATS*, 2021.
- [MCP14] Stefan Magureanu, Richard Combes, and Alexandre Proutiere. Lipschitz bandits: Regret lower bounds and optimal algorithms. In *Proc. of COLT*, 2014.
- [Neu15] Gergely Neu. First-order regret bounds for combinatorial semi-bandits. In *Proc. of COLT*, 2015.

- [Oka73] Masashi Okamoto. Distinctness of the eigenvalues of a quadratic form in a multivariate sample. *The Annals of Statistics*, 1973.
- [PBVP20] Pierre Perrault, Etienne Boursier, Michal Valko, and Vianney Perchet. Statistical efficiency of thompson sampling for combinatorial semi-bandits. In *Proc. of NeurIPS*, 2020.
- [Per22] Pierre Perrault. When combinatorial thompson sampling meets approximation regret. In *Proc. of NeurIPS*, 2022.
- [PPV19] Pierre Perrault, Vianney Perchet, and Michal Valko. Exploiting structure of uncertainty for efficient matroid semi-bandits. In *Proc. of ICML*, 2019.
- [RG96] Ram Ravi and Michel X Goemans. The constrained minimum spanning tree problem. In *Scandinavian Workshop on Algorithm Theory*. Springer, 1996.
- [RS13] Sasha Rakhlin and Karthik Sridharan. Optimization, learning, and games with predictable sequences. In *Proc. of NeurIPS*, 2013.
- [S⁺03] Alexander Schrijver et al. *Combinatorial optimization: polyhedra and efficiency*, volume 24. Springer, 2003.
- [SN20] Arun Suggala and Praneeth Netrapalli. Follow the perturbed leader: Optimism and fast parallel algorithms for smooth minimax games. In *Proc. of NeurIPS*, 2020.
- [Vis21] Nisheeth K. Vishnoi. *Algorithms for Convex Optimization*. Cambridge University Press, 2021.
- [WLZL21] Chen-Yu Wei, Chung-Wei Lee, Mengxiao Zhang, and Haipeng Luo. Linear last-iterate convergence in constrained saddle-point optimization. In *Prof. of ICLR*, 2021.
- [WTP21] Po-An Wang, Ruo-Chun Tzeng, and Alexandre Proutiere. Fast pure exploration via frank-wolfe. In *Proc. of NeurIPS*, 2021.
- [ZODS21] Tom Zahavy, Brendan O’Donoghue, Guillaume Desjardins, and Satinder Singh. Reward is enough for convex mdps. In *Proc. of NeurIPS*, 2021.

Contents

1	Introduction	1
2	Preliminaries	3
2.1	The lower-bound problem	4
2.2	The Linear Maximization Oracle	4
3	Solving the lower bound problem: the MCP algorithm	4
3.1	Lagrangian relaxation	4
3.2	Solving the two-player game with no regret	5
3.3	Performance analysis of the MCP algorithm	6
4	The Perturbed Frank-Wolfe Sampling (P-FWS) algorithm	6
4.1	Smoothing the objective function F_μ	7
4.2	The algorithm	7
4.3	Non-asymptotic performance analysis of P-FWS	9
5	Related Work	10
6	Conclusion	10
A	Table of Notation	16
B	Further related work	17
C	Results related to our (ϵ, θ)-MCP algorithm	19
C.1	Properties of Lagrangian dual of f_x	19
C.2	Analysis of MCP	20
C.3	Regret analysis of Follow-the-Perturbed-Leader	21
D	Analysis of P-FWS	24
D.1	δ -correctness (Theorem 4 (i))	24
D.2	Almost-sure upper bound (Theorem 4 (ii))	24
D.3	Non-asymptotic sample complexity (Theorem 4 (ii))	25
D.3.1	Good events	26
D.3.2	Proof of non-asymptotic sample complexity	27
D.3.3	Technical lemmas	28
D.4	Computational complexity (Theorem 4 (iii))	29
E	Convergence of P-FWS under the good events	31
F	Upper bound of $\sum_{T=M+1}^{\infty} \mathbb{P}_\mu[(\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T))^c]$ under P-FWS	34
F.1	Lemmas for bounding $\mathbb{P}_\mu[\mathcal{E}_{1,\epsilon}(T)^c]$	35

F.2	Lemmas for bounding $\mathbb{P}_\mu[\mathcal{E}_{2,\epsilon}(T)^c]$	36
F.3	Technical lemmas	36
G	Continuity arguments	38
G.1	An application of the maximum theorem	39
G.2	The length of gradients	40
H	Stochastic smoothing	43
I	Lipschitzness of F_μ and boundness of F_μ on $\Sigma_K \cap \mathbb{R}_{>0}^K$	44
J	Proofs related to combinatorial sets	45
K	Sample complexity lower bound	47
L	Extension to the transductive setting	48

A Table of Notation

Problem setting	
K	Number of arms
$[m]$ for any $m \in \mathbb{N}$	The set $\{1, 2, \dots, m\}$
δ	Required uncertainty
$\boldsymbol{\mu} \in \mathbb{R}^K$	Vector of the expected rewards of the various arms
\mathbb{E}_μ and \mathbb{P}_μ	The expectation and probability measure corresponding to $\boldsymbol{\mu}$
Λ	$\{\boldsymbol{\mu} \in \mathbb{R}^K : \mathbf{i}^*(\boldsymbol{\mu}) = 1\}$
$\mathbf{i}^*(\boldsymbol{\mu})$	Best arm under parameter $\boldsymbol{\mu}$
\mathcal{X}	Set of actions in $\{0, 1\}^K$
D	$\max_{\mathbf{x} \in \mathcal{X}} \ \mathbf{x}\ _1$
$\Delta_{\mathbf{x}}(\boldsymbol{\mu})$	$\langle \mathbf{i}^*(\boldsymbol{\mu}) - \mathbf{x}, \boldsymbol{\mu} \rangle$
$\Delta_{\min}(\boldsymbol{\mu})$	$\min_{\mathbf{x} \neq \mathbf{i}^*(\boldsymbol{\mu})} \Delta_{\mathbf{x}}(\boldsymbol{\mu})$
Notation related to a given algorithm	
$N_k(t)$	Number of pulls of arm k up to time t
$\hat{\omega}_k(t)$	$N_k(t)/t$
$\mathbf{x}(t)$	The action taken in time t
$y_k(t)$	Random reward received if $x_k(t) = 1$
$\hat{\mu}_k(t)$	$\sum_{s=1}^t y_k(s) \mathbb{1}\{x_k(s) = 1\} / N_k(t)$
τ	Stopping time
$\hat{\mathbf{i}}$	Recommended action
Notation used for sets and vectors	
\odot	Elementwise product
\oplus	Elementwise sum over \mathbb{Z}_2
\mathbf{x}^i	The i -th elementwise power of $\mathbf{x} \in \mathbb{R}^K$, i.e., $(x_k^i)_{k \in [K]}$
$\text{cl}(\mathcal{S})$	The closure of set \mathcal{S}
\mathbf{e}_k	the K -dimensional vector whose k -th component is equal to one and zero elsewhere
Properties for lower bound	
$d(\mu, \mu')$	KL divergence between the distributions parametrized by μ and μ'
$\text{kl}(a, b)$	KL divergence between two Bernoulli distributions of means a and b
$\text{Alt}(\boldsymbol{\mu})$	$\{\boldsymbol{\lambda} \in \Lambda : \mathbf{i}^*(\boldsymbol{\lambda}) \neq \mathbf{i}^*(\boldsymbol{\mu})\}$
Σ	$\{\sum_{\mathbf{x} \in \mathcal{X}} w_{\mathbf{x}} \mathbf{x} : \mathbf{w} \in \Sigma_{ \mathcal{X} }\}$ where Σ_N is a $(N - 1)$ -dimensional simplex
Σ_+	$\Sigma \cap \mathbb{R}_{>0}^K$
Notation for MCP	
$F_\mu(\boldsymbol{\omega})$	$\min_{\mathbf{x} \neq \mathbf{i}^*(\boldsymbol{\mu})} f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\mu})$
$f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\mu})$	$\inf_{\boldsymbol{\lambda} \in \mathcal{C}_{\mathbf{x}}} \langle \boldsymbol{\omega}, \frac{(\boldsymbol{\mu} - \boldsymbol{\lambda})^2}{2} \rangle$, where $\mathcal{C}_{\mathbf{x}} = \{\boldsymbol{\lambda} \in \mathbb{R}^K : \langle \boldsymbol{\lambda}, \mathbf{i}^*(\boldsymbol{\mu}) - \mathbf{x} \rangle < 0\}$
$\mathcal{L}_{\boldsymbol{\omega}, \boldsymbol{\mu}}(\boldsymbol{\lambda}, \mathbf{x}, \alpha)$	$\langle \boldsymbol{\omega}, \frac{(\boldsymbol{\mu} - \boldsymbol{\lambda})^2}{2} \rangle + \alpha \langle \mathbf{i}^*(\boldsymbol{\mu}) - \mathbf{x}, \boldsymbol{\lambda} \rangle$
$g_{\boldsymbol{\omega}, \boldsymbol{\mu}}(\mathbf{x}, \alpha)$	$\inf_{\boldsymbol{\lambda} \in \mathbb{R}^K} \mathcal{L}_{\boldsymbol{\omega}, \boldsymbol{\mu}}(\boldsymbol{\lambda}, \mathbf{x}, \alpha)$
Notation for P-FWS	
\mathcal{X}_0	A $[K]$ -covering set
\hat{F}_t	MCP-approximated value of $F_{\hat{\mu}(t)}(\hat{\boldsymbol{\omega}}(t))$ for stopping rule
$\bar{F}_{\boldsymbol{\mu}, \eta}(\cdot)$	$\mathbb{E}_{\boldsymbol{Z} \sim \text{Uniform}(B_2)}[\nabla F_{\boldsymbol{\mu}}(\cdot + \eta \boldsymbol{Z})]$ where $B_2 = \{\mathbf{v} \in \mathbb{R}^K : \ \mathbf{v}\ _2 \leq 1\}$
$\bar{F}_{\boldsymbol{\mu}, \eta, n}$	The empirical n -sample estimate of $\bar{F}_{\boldsymbol{\mu}, \eta}$
ℓ	Lipschitz constant of F_μ

B Further related work

Combinatorial semi-bandits [CBL12] have found numerous applications including online ranking [DKC21], network routing [CLK⁺14, KWA⁺14], loan assignment [KWA⁺14], path planning problem [JMKK21], and influence marketing [Per22]). We do not discuss these applications here, but rather focus the literature that is the most relevant to our analysis and results.

Solving the lower-bound problem in combinatorial semi-bandits. We are not aware of any computationally efficient algorithm to solve the lower-bound problem, or to compute its objective function. To the best of our knowledge, MCP is the first algorithm to do so. A work closed to ours is [CCG21a] for combinatorial semi-bandits but in the regret minimization. Regret minimization yields a different lower-bound problem. There exists a statistically optimal algorithm [CMP17], called OSSB, that matches the regret lower bound by [CTMSP⁺15]. OSSB requires to solve the lower-bound problem in each round, and the authors [CCG21a] are the first to investigate whether this is at all possible in a computationally efficient way. They establish that if budgeted-linear maximization (BLM) [RG96, BBGS11] can be solved within an ε -approximation factor for the combinatorial set \mathcal{X} , then the lower-bound problem can be approximately solved with a precision depending on ε . As a consequence, the approach leads to an algorithm with asymptotically minimal regret only if one has access to an exact BLM solver. This is the case for m -sets and s - t paths but this is not the case for spanning trees and perfect matchings. For the latter case, as mentioned [CCG21a], an algorithm using an approximately correct BLM solver would not be statistically optimal.

Best arm identification in combinatorial semi-bandits. Many tasks related to combinatorial best arm identification are formulated in the *transductive* setting [JMKK21], where the set $\mathcal{A} \subseteq \{0, 1\}^K$ available for exploration is not necessarily the same as the set $\mathcal{X} \subseteq \{0, 1\}^K$ for decision. The minimal sample complexity in the transductive setting is exactly (1) with Σ replaced with $\{\sum_{x \in \mathcal{A}} w_x x : \omega \in \Sigma_{|\mathcal{A}|}\}$ - see (58) in Appendix L for details. Two most studied tasks are combinatorial multi-arm bandit (C-MB) where $\mathcal{A} = \{e_k\}_{k \in [K]}$ and the best action identification (C-BAI) where $\mathcal{A} = \mathcal{X}$. The former is arguably simpler than the latter if we compare the corresponding minimal sample complexities (note that $\Sigma_K \supseteq \Sigma$). We note that our results for C-BAI can be easily generalized to the transductive setting (see Appendix L).

Prior works mainly focus on the C-MB task. UCB-based [KTAS12, CLK⁺14] and elimination-based [CGL16, CGL⁺17, KSJJ⁺20] approaches are popular. Among these, `EfficientGapElim` [CGL⁺17] achieves the lowest sample complexity $\mathcal{O}(T^*(\mu)(\ln \delta^{-1} + \ln^2 \Delta_{\min}^{-1} (\ln \ln \Delta_{\min}^{-1} + \ln |\mathcal{X}|)))$ with high probability⁴, but its computational complexity is hard to analyze. `Peace` [KSJJ⁺20], another elimination-based approach by experimental design, requires with high probability a polynomial number of the LM Oracle calls in total. The sample complexity of `Peace` has a δ -dependent term (scaling as $KT^*(\mu) \ln \delta^{-1}$) worse than `EfficientGapElim`. Overall, none of these are statistically optimal when $\delta \rightarrow 0$. Note that algorithms for linear best-arm identification [DMSV20, WTP21] are applicable to C-MB but not to C-BAI and the general transductive setting.

For the task of C-BAI, we are only aware of two works: `GCB-PE` [DKC21] and `CombGame` [JMKK21]. `GCB-PE` is a UCB-based algorithm with guarantees on the sample complexity and computational complexity valid with high probability only. `CombGame` [JMKK21] is proposed for the transductive setting, and its design inherits from [DKM19] that interprets the lower-bound problem and more precisely $T^*(\mu)^{-1}$ as the value of a two-player game (a ω -player and a λ -player)⁵. Assuming that an MCP oracle is available, `CombGame` leverages Frank-Wolfe algorithms, namely OFW [HK12] and LLOO [GH16], for the ω -player and the MCP algorithm for the λ -player. [JMKK21] leaves the existence of such an oracle running in polynomial time as an open problem. Our MCP algorithm resolves this issue. `CombGame` is statistically optimal in the high confidence regime but has no clear guarantees in the moderate regime [BGK22].

We wish to finally mention an algorithm that has inspired the design of `P-FWS`. This algorithm is referred to as Frank-Wolfe Sampling (`FWS`) [WTP21]. `FWS` is optimal in high confidence regime

⁴In Section 4.5 in [CGL⁺17], the authors provide a lemma stating that: if *parallel simulation* is additionally allowed, then any high-probability sample complexity upper bound can be converted to an upper bound in expectation.

⁵Note that this two-player game is different than the two-player game involved in our algorithm MCP.

but is not computationally efficient for combinatorial semi-bandits. For example, to deal with the non-smoothness issue of the objective function F_μ , FWS needs to construct the so-called r -subdifferentiable spaces and to optimize a linear function on these spaces. Unfortunately, these spaces can be generated by a number of vectors exponentially increasing with K in combinatorial semi-bandits. Moreover, in moderate confidence regime, the sample complexity upper bound derived in [WTP21] has an exponential dependence in K .

All the relevant algorithms, their sample complexity guarantees and computational complexity are summarized in Table 1.

Table 1: Algorithms for best-arm identification in combinatorial semi-bandits with fixed confidence and their performance.

Algorithm	Task	Instance-specific Sample Complexity		Computational Complexity	
		Non-asympt.	Asympt. Opt.	Needed (Provided)	Total LM oracle calls
Peace	C-MB	$\text{poly}(K, \Delta_{\min}^{-1}, \ln \delta^{-1})$ w.h.p.	\times	LP solver (\checkmark)	$\text{poly}(K, \Delta_{\min}^{-1}, \delta^{-1})$ w.h.p.
GCB-PE	C-BAI	$\text{poly}(K, \Delta_{\min}^{-1}, \ln \delta^{-1})$ w.h.p.	\times	-	$\text{poly}(K, \Delta_{\min}^{-1}, \ln \delta^{-1})$ w.h.p.
CombGame	Trans.	\times (incomparable)	\checkmark	MCP (\times)	\times
P-FWS	Trans.	$\text{poly}(K, \Delta_{\min}^{-1}, \ln \delta^{-1})$	\checkmark	MCP (\checkmark)	$\text{poly}(K, \Delta_{\min}^{-1}, \ln \delta^{-1})$

C Results related to our (ϵ, θ) -MCP algorithm

C.1 Properties of Lagrangian dual of f_x

Proposition 1. Let $(\omega, \mu) \in \Sigma_+ \times \Lambda$ and $x \in \mathcal{X} \setminus \{i^*(\mu)\}$.

(a) The Lagrange dual function is linear in x . More precisely, $g_{\omega, \mu}(x, \alpha) = c_{\omega, \mu}(\alpha) + \langle \ell_{\omega, \mu}(\alpha), x \rangle$ where $c_{\omega, \mu}(\alpha) = \alpha \langle \mu - \frac{\alpha}{2} \omega^{-1}, i^*(\mu) \rangle$ and $\ell_{\omega, \mu}(\alpha) = -\alpha (\mu + \frac{\alpha}{2} \omega^{-1} \odot (\mathbf{1}_K - 2i^*(\mu)))$.

(b) $g_{\omega, \mu}(x, \cdot)$ is strictly concave (for any fixed x).

(c) $f_x(\omega, \mu) = \max_{\alpha \geq 0} g_{\omega, \mu}(x, \alpha)$ is attained by $\alpha_x^* = \frac{\Delta_x(\mu)}{\langle x \oplus i^*(\mu), \omega^{-1} \rangle}$.

(d) $\|\ell_{\omega, \mu}(\alpha_x^*)\|_1 \leq L_{\omega, \mu} = 4D^2K \|\mu\|_\infty^2 \|\omega^{-1}\|_\infty$.

Proof Fix any $(\omega, \mu) \in \Sigma_+ \times \Lambda$ and let $i^* = i^*(\mu)$ for short. For convenience, the definition of $\mathcal{L}_{\omega, \mu}$ and $g_{\omega, \mu}$ are restated:

$$\mathcal{L}_{\omega, \mu}(\lambda, x, \alpha) = \left\langle \omega, \frac{(\mu - \lambda)^2}{2} \right\rangle + \alpha \langle i^* - x, \lambda \rangle \quad \text{and} \quad g_{\omega, \mu}(x, \alpha) = \inf_{\lambda \in \mathbb{R}^K} \mathcal{L}_{\omega, \mu}(\lambda, x, \alpha).$$

Proof of (a): linearity of $g_{\omega, \mu}(\cdot, \alpha)$: Let $\lambda_{\omega, \mu}^*(x, \alpha) \in \arg \inf_{\lambda \in \mathbb{R}^K} \mathcal{L}_{\omega, \mu}(\lambda, x, \alpha)$. The first-order condition implies that $\mathbf{0}_K = \nabla_\lambda \mathcal{L}_{\omega, \mu}(\lambda_{\omega, \mu}^*(x, \alpha), x, \alpha) = \omega \odot (\lambda_{\omega, \mu}^*(x, \alpha) - \mu) + \alpha(i^* - x)$, which directly yields (as $\omega > \mathbf{0}_K$)

$$\lambda_{\omega, \mu}^*(x, \alpha) = \mu + \alpha \omega^{-1} \odot (x - i^*). \quad (11)$$

We plug (11) into $\mathcal{L}_{\omega, \mu}(\lambda_{\omega, \mu}^*(x, \alpha), x, \alpha)$ and directly obtain that

$$\begin{aligned} g_{\omega, \mu}(x, \alpha) &= \left\langle \omega, \frac{\alpha^2}{2} \omega^{-2} \odot (x - i^*)^2 \right\rangle + \alpha \langle \mu, i^* - x \rangle - \alpha^2 \langle \omega^{-1}, (x - i^*)^2 \rangle \\ &= \alpha \langle \mu, i^* - x \rangle - \frac{\alpha^2}{2} \langle \omega^{-1}, (x - i^*)^2 \rangle \end{aligned} \quad (12)$$

$$= c_{\omega, \mu}(\alpha) + \langle \ell_{\omega, \mu}(\alpha), x \rangle, \quad (13)$$

where (13) follows from a fact that $(x - i^*)^2 = i^* - 2x \odot i^* + x = i^* + x \odot (\mathbf{1}_K - 2i^*)$.

Proof of (b): strict concavity of $g_{\omega, \mu}(x, \cdot)$: This is trivial from (12).

Proof of (c): $f_x(\omega, \mu) = \max_{\alpha \geq 0} g_{\omega, \mu}(x, \alpha)$ is attained by $\alpha_x^* = \frac{\Delta_x(\mu)}{\langle x \oplus i^*, \omega^{-1} \rangle}$: For a fixed $x \neq i^*$, by the first-order condition of (12), we find that the maximum of $g_{\omega, \mu}(x, \cdot)$ is reached at

$$\alpha_x^* = \frac{\Delta_x(\mu)}{\langle \omega^{-1}, (x - i^*)^2 \rangle} = \frac{\Delta_x(\mu)}{\langle x \oplus i^*, \omega^{-1} \rangle}, \quad (14)$$

where for the second equality, we use the assumption that i^* and x are binary vectors and hence $(x - i^*)^2 = x \oplus i^*$. We now verify that $(\alpha_x^*, \lambda_x^*)$, where $\lambda_x^* = \lambda_{\omega, \mu}^*(x, \alpha_x^*)$ (see (11)), satisfies KKT conditions, which is equivalent to strong duality (refer to [Vis21, BV04]) under Slater's condition (there exists a $\lambda \in \mathbb{R}^K$ such that the constraint is strict). Since x is a suboptimal action, $\Delta_x(\mu)$ is positive, so is α_x^* (dual feasibility). To verify $\langle \lambda_x^*, i^* - x \rangle \leq 0$ (primal feasibility), the definition of $\lambda_{\omega, \mu}^*(\cdot, \cdot)$, (11), yields

$$\begin{aligned} \langle \lambda_x^*, i^* - x \rangle &= \Delta_x(\mu) + \alpha_x^* \langle \omega^{-1} \odot (x - i^*), (i^* - x) \rangle \\ &= \Delta_x(\mu) - \alpha_x^* \langle \omega^{-1}, x \oplus i^* \rangle = 0, \end{aligned}$$

which implies that $\alpha_x^* \langle i^* - x, \lambda_x^* \rangle = 0$ (complementary slackness). Finally, stationarity holds automatically as $\nabla_\lambda \mathcal{L}_{\omega, \mu}(\lambda_x^*, x, \alpha) = 0$ for all α .

Proof of (d): $\|\ell_{\omega, \mu}(\alpha_x^*)\|_1 \leq L_{\omega, \mu} = 4D^2K \|\mu\|_\infty^2 \|\omega^{-1}\|_\infty$: Following from the expression of $\ell_{\omega, \mu}(\alpha)$, we have $\ell_{\omega, \mu}(\alpha_x^*) = -\alpha_x^* \mu + \frac{\alpha_x^{*2}}{2} \omega^{-1} \odot (\mathbf{1}_K - 2i^*)$. Observe that $\|\mu\|_1 \leq K \|\mu\|_\infty \leq$

$K \|\omega^{-1}\|_\infty \|\mu\|_\infty$ (as $\omega \in \Sigma_+$) and the coordinate of $\mathbf{1}_K - 2\mathbf{i}^*$ is either 1 or -1 , a simple application of triangle inequality leads to

$$\|\ell_{\omega, \mu}(\alpha_{\mathbf{x}}^*)\|_1 \leq K \|\omega^{-1}\|_\infty \left(\|\mu\|_\infty + \frac{\alpha_{\mathbf{x}}^*}{2} \right) \alpha_{\mathbf{x}}^*.$$

As for $\alpha_{\mathbf{x}}^*$ (see (14)), $\Delta_{\mathbf{x}}(\mu) \leq 2D \|\mu\|_\infty$ and $\langle \omega^{-1}, \mathbf{x} \oplus \mathbf{i}^* \rangle \geq \min_k \omega_k^{-1} \geq 1$, hence we conclude that $\|\ell_{\omega, \mu}(\alpha_{\mathbf{x}}^*)\|_1 \leq 2D(D+1)K \|\mu\|_\infty^2 \|\omega^{-1}\|_\infty \leq L_{\omega, \mu}$. \square

C.2 Analysis of MCP

Theorem 3. *Let $\epsilon, \theta \in (0, 1)$. Under Assumption 1, for any $(\omega, \mu) \in \Sigma_+ \times \Lambda$, the (ϵ, θ) -MCP(ω, μ) algorithm outputs $(\hat{F}, \hat{\mathbf{x}})$ satisfying*

$$\mathbb{P} \left[F_\mu(\omega) \leq \hat{F} \leq (1 + \epsilon) F_\mu(\omega) \right] \geq 1 - \theta \quad \text{and} \quad \hat{F} = \max_{\alpha \geq 0} g_{\omega, \mu}(\hat{\mathbf{x}}, \alpha).$$

Moreover, the number of LM Oracle calls the algorithm does is almost surely at most

$$\left\lceil \frac{c_\theta^2 (1 + \epsilon)^2}{\epsilon^2 F_\mu(\omega)^2} \right\rceil = \mathcal{O} \left(\frac{\|\mu\|_\infty^4 \|\omega^{-1}\|_\infty^2 K^3 D^5 \ln K \ln \theta^{-1}}{\epsilon^2 F_\mu(\omega)^2} \right).$$

Proof Fix any $(\omega, \mu) \in \Sigma_+ \times \Lambda$ and denote by $\mathbf{i}^* = \mathbf{i}^*(\mu)$. Suppose Algorithm 1 reaches the stopping criterion at the N -th iteration.

Guarantees on the outputs of MCP: By Proposition 1 (a),

$$\sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)}) - \min_{\mathbf{x} \neq \mathbf{i}^*} \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}, \alpha^{(n)}) = \sum_{n=1}^N \langle \ell_{\omega, \mu}(\alpha^{(n)}), \mathbf{x}^{(n)} \rangle - \min_{\mathbf{x} \neq \mathbf{i}^*} \sum_{n=1}^N \langle \ell_{\omega, \mu}(\alpha^{(n)}), \mathbf{x} \rangle.$$

The regret of \mathbf{x} -player can be bounded by applying Lemma 3, resulting in:

$$\mathbb{P} \left[\sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)}) - \min_{\mathbf{x} \neq \mathbf{i}^*} \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}, \alpha^{(n)}) \leq c_\theta \sqrt{N} \right] \geq 1 - \theta. \quad (15)$$

To relate $F_\mu(\omega)$ with (15), let \mathbf{x}_e be the minimizer attaining $F_\mu(\omega) = f_{\mathbf{x}_e}(\omega, \mu)$. Then,

$$\min_{\mathbf{x} \neq \mathbf{i}^*} \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}, \alpha^{(n)}) \leq \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}_e, \alpha^{(n)}) \leq N \max_{\alpha \geq 0} g_{\omega, \mu}(\mathbf{x}_e, \alpha) = N F_\mu(\omega). \quad (16)$$

Recall that $\alpha^{(n)}$ is chosen as the best response $\max_{\alpha \geq 0} g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha) = g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)})$ and that $\hat{F} = \min_{n \in [N]} g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)})$. These together with (16) imply that

$$N(\hat{F} - F_\mu(\omega)) \leq \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)}) - \min_{\mathbf{x} \neq \mathbf{i}^*} \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}, \alpha^{(n)}). \quad (17)$$

A simple rearrangement on (15) and (17) implies that: with probability at least $1 - \theta$,

$$\hat{F} - F_\mu(\omega) \leq \frac{c_\theta}{\sqrt{N}} \leq \epsilon(\hat{F} - \frac{c_\theta}{\sqrt{N}}) \leq \epsilon F_\mu(\omega),$$

where the second inequality follows from the stopping criterion that $\sqrt{N} > c_\theta(1 + \epsilon)/\epsilon\hat{F}$, and the last inequality simply comes from the rearrangement of the first inequality.

Computational cost: From the stopping criterion of MCP, we know that

$$N = \left\lceil \frac{c_\theta^2 (1 + \epsilon)^2}{\epsilon^2 \hat{F}^2} \right\rceil \leq \left\lceil \frac{c_\theta^2 (1 + \epsilon^{-1})^2}{F_\mu(\omega)^2} \right\rceil = \mathcal{O} \left(\frac{L_{\omega, \mu}^2 \left(\sqrt{K \ln K} + \sqrt{\ln \theta^{-1}} \right)^2}{\epsilon^2 F_\mu(\omega)^2} \right)$$

since $\hat{F} \geq F_\mu(\omega)$ and $c_\theta = L_{\omega, \mu} \left(4\sqrt{K(\ln K + 1)} + \sqrt{\ln(\theta^{-1})/2} \right)$. Finally, as computing each $\mathbf{x}^{(n)}$ takes at most D calls to LM Oracle, the total number of LM Oracle calls is

$$\mathcal{O} \left(\frac{L_{\omega, \mu}^2 D \left(\sqrt{K \ln K} + \sqrt{\ln \theta^{-1}} \right)^2}{\epsilon^2 F_\mu(\omega)^2} \right) = \mathcal{O} \left(\frac{\|\mu\|_\infty^4 \|\omega^{-1}\|_\infty^2 K^3 D^5 \ln K \ln \theta^{-1}}{\epsilon^2 F_\mu(\omega)^2} \right)$$

by recalling $L_{\omega, \mu} = 4D^2 K \|\mu\|_\infty^2 \|\omega^{-1}\|_\infty$ from Proposition 1 (d) and $(\sqrt{\ln K} + \sqrt{\ln \theta^{-1}})^2 = \mathcal{O}(\ln K \ln \theta^{-1})$. \square

C.3 Regret analysis of Follow-the-Perturbed-Leader

In this subsection, we aim at proving Lemma 3, which is a direct consequence of Lemma 4. One can find similar proofs in e.g. [KV05, Neu15, SN20]. However, the parameter η_n in our MCP algorithm is varying and carefully chosen (without the knowledge of the last round), which makes the proof slightly more complicated.

Lemma 3. *Let $N \in \mathbb{N}$. Under (ϵ, θ) -MCP(ω, μ), then*

$$\mathbb{P} \left[\frac{1}{N} \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}^{(n)}, \alpha^{(n)}) - \frac{1}{N} \min_{\mathbf{x} \neq \mathbf{i}^*} \sum_{n=1}^N g_{\omega, \mu}(\mathbf{x}, \alpha^{(n)}) \leq \frac{c_\theta}{\sqrt{N}} \right] \geq 1 - \theta.$$

Lemma 4. *Let $\theta \in (0, 1)$ and $\mathcal{M} \subseteq \{0, 1\}^K$. Given an arbitrary sequence $\{\ell_n\}_{n \geq 1}$ of vectors in \mathbb{R}^K whose length $\|\ell_n\|_1$ is bounded by $L > 0$ for all $n \in \mathbb{N}$. Suppose $\{\mathbf{x}^{(n)}\}_{n \geq 1}$ is generated by*

$$\mathbf{x}^{(n)} \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{M}} \left(\sum_{m=1}^{n-1} \langle \ell_m, \mathbf{x} \rangle + \left\langle \frac{\mathcal{Z}_n}{\eta_n}, \mathbf{x} \right\rangle \right),$$

where $\mathcal{Z}_n = (\mathcal{Z}_{1,n}, \dots, \mathcal{Z}_{K,n})$ is a random vector with uncorrelated exponentially distributed (with unit mean) components, and $\eta_n = \sqrt{K(\ln K + 1)/(4nL^2)}$. Then, for any $N \in \mathbb{N}$,

$$\mathbb{P} \left[\sum_{n=1}^N \langle \ell_n, \mathbf{x}^{(n)} \rangle - \min_{\mathbf{x} \in \mathcal{M}} \sum_{n=1}^N \langle \ell_n, \mathbf{x} \rangle \leq L\sqrt{N} \left(4\sqrt{K(\ln K + 1)} + \sqrt{\frac{\ln \theta^{-1}}{2}} \right) \right] \geq 1 - \theta.$$

Proof of Lemma 4: We will prove this lemma as if $\{\ell_n\}_n$ is chosen in advance since there exists a standard technique for extending regret against oblivious player to the one against nonoblivious one (see Lemma 4.1 in [CBL06]). For convenience, we introduce the following notation. Let $\mathbf{m}^*(\cdot) = \operatorname{argmin}_{\mathbf{x} \in \mathcal{M}} \langle \cdot, \mathbf{x} \rangle$. Finally, further define global minimizer $\mathbf{x}_* = \mathbf{m}^* \left(\sum_{n=1}^N \ell_n \right)$ and an auxiliary vector $\mathbf{b}^{(n)} = \mathbf{m}^* \left(\sum_{m=1}^n \ell_m + \mathcal{Z}_1/\eta_n \right)$.

It suffices to show the expected regret bound (18).

$$\mathbb{E} \left[\sum_{n=1}^N \langle \ell_n, \mathbf{x}^{(n)} \rangle \right] - \min_{\mathbf{x} \in \mathcal{M}} \sum_{n=1}^N \langle \ell_n, \mathbf{x} \rangle \leq 4L\sqrt{NK(\ln K + 1)}. \quad (18)$$

This is because $\{\langle \ell_n, \mathbf{x}^{(n)} \rangle - \mathbb{E}[\langle \ell_n, \mathbf{x}^{(n)} \rangle]\}_n$ forms a sequence of bounded martingale difference, so an application of a concentration inequality (Lemma 6) with $V_n = \langle \ell_n, \mathbf{x}^{(n)} \rangle - \mathbb{E}[\langle \ell_n, \mathbf{x}^{(n)} \rangle]$, $r_n = L$ for $n \in [N]$, and $s = L\sqrt{N \ln \theta^{-1}/2}$ gives that

$$\mathbb{P} \left[\sum_{n=1}^N \langle \ell_n, \mathbf{x}^{(n)} \rangle - \mathbb{E} \left[\sum_{n=1}^N \langle \ell_n, \mathbf{x}^{(n)} \rangle \right] > L\sqrt{\frac{N \ln \theta^{-1}}{2}} \right] \leq \theta$$

and combining this with (18) completes the proof.

Proof of (18): We decompose the regret into two terms:

$$\sum_{n=1}^N \langle \ell_n, \mathbf{x}^{(n)} - \mathbf{x}_* \rangle = \sum_{n=1}^N \langle \mathbf{b}^{(n)} - \mathbf{x}_*, \ell_n \rangle + \sum_{n=1}^N \langle \mathbf{x}^{(n)} - \mathbf{b}^{(n)}, \ell_n \rangle.$$

(i). We show that $\mathbb{E} \left[\sum_{n=1}^N \langle \mathbf{b}^{(n)} - \mathbf{x}_*, \ell_n \rangle \right] \leq \frac{K(\ln K + 1)}{\eta_N}$. Invoking Lemma 5 with $\mathbf{x} = \mathbf{x}_*$ results in

$$\begin{aligned} \mathbb{E} \left[\sum_{n=1}^N \langle \mathbf{b}^{(n)} - \mathbf{x}_*, \ell_n \rangle \right] &\leq \mathbb{E} \left[\left\langle \frac{\mathbf{x}_*}{\eta_N} - \left(\frac{\mathbf{b}^{(1)}}{\eta_1} + \sum_{n=2}^N \left(\frac{1}{\eta_n} - \frac{1}{\eta_{n-1}} \right) \mathbf{b}^{(n)} \right), \mathcal{Z}_1 \right\rangle \right] \\ &\leq \mathbb{E} \left[\left\| \mathcal{Z}_1 \right\|_\infty \left\| \frac{\mathbf{x}_*}{\eta_N} - \left(\frac{\mathbf{b}^{(1)}}{\eta_1} + \sum_{n=2}^N \left(\frac{1}{\eta_n} - \frac{1}{\eta_{n-1}} \right) \mathbf{b}^{(n)} \right) \right\| \right], \end{aligned}$$

where the last inequality uses Hölder's inequality. As all the components of \mathbf{x}_* and $\frac{\eta_N \mathbf{b}^{(1)}}{\eta_1} + \sum_{n=2}^N \eta_N \left(\frac{1}{\eta_n} - \frac{1}{\eta_{n-1}} \right) \mathbf{b}^{(n)}$ are nonnegative and bounded by 1, the 1-norm of their difference is bounded by K . It remains to show

$$\mathbb{E}[\|\mathcal{Z}_1\|_\infty] = \int_0^\infty \mathbb{P}[\max_i \mathcal{Z}_{1,i} \geq x] dx \leq \int_0^{\ln K} \mathbb{P}[\max_i \mathcal{Z}_{1,i} \geq x] dx + \int_{\ln K}^\infty K e^{-x} dx \leq \ln K + 1.$$

(ii). We show that $\mathbb{E} \left[\sum_{n=1}^N \langle \mathbf{x}^{(n)} - \mathbf{b}^{(n)}, \ell_n \rangle \right] \leq 2L^2 \sum_{n=1}^N \eta_n$. Let the pdf of $\exp(1)$ be $\pi(\cdot) = e^{-\|\cdot\|_1}$.

$$\begin{aligned} \mathbb{E} \left[\langle \mathbf{b}^{(n)}, \ell_n \rangle \right] &= \int_{\mathbf{z} \in \mathbb{R}^K} \left\langle \mathbf{m}^* \left(\eta_n \sum_{m=1}^n \ell_m + \mathbf{z} \right), \ell_n \right\rangle d\pi(\mathbf{z}) \\ &= \int_{\mathbf{y} \in \mathbb{R}^K} \left\langle \mathbf{m}^* \left(\eta_n \sum_{m=1}^{n-1} \ell_m + \mathbf{y} \right), \ell_n \right\rangle d\pi(\mathbf{y} - \eta_n \ell_n) \\ &= \int_{\mathbf{y} \in \mathbb{R}^K} \left\langle \mathbf{m}^* \left(\eta_n \sum_{m=1}^{n-1} \ell_m + \mathbf{y} \right), \ell_n \right\rangle e^{-\|\mathbf{y} - \eta_n \ell_n\|_1 + \|\mathbf{y}\|_1} d\pi(\mathbf{y}). \end{aligned}$$

Notice that the triangular inequality implies $-\|\mathbf{y} - \eta_n \ell_n\|_1 + \|\mathbf{y}\|_1 \leq \|\eta_n \ell_n\|_1 \leq \eta_n L$ and $e^x \leq 1 + 2x$ for all $x \in (0, 1)$ (Taylor expansion), so recalling $\mathbf{x}^{(n)} = \mathbf{m}^* \left(\eta_n \sum_{m=1}^{n-1} \ell_m + \mathcal{Z}_1 \right)$, we deduce that

$$\begin{aligned} \sum_{n=1}^N \mathbb{E} \left[\langle \mathbf{x}^{(n)} - \mathbf{b}^{(n)}, \ell_n \rangle \right] &\leq \sum_{n=1}^N 2\eta_n L \int_{\mathbf{z} \in \mathbb{R}^K} \left\langle \mathbf{m}^* \left(\eta_n \sum_{m=1}^n \ell_m + \mathbf{z} \right), \ell_n \right\rangle d\pi(\mathbf{z}) \\ &\leq \sum_{n=1}^N 2\eta_n L \int_{\mathbf{z} \in \mathbb{R}^K} \left\| \mathbf{m}^* \left(\eta_n \sum_{m=1}^n \ell_m + \mathbf{z} \right) \right\|_\infty \|\ell_n\|_1 d\pi(\mathbf{z}) \leq 2L^2 \sum_{n=1}^N \eta_n. \end{aligned}$$

Finally, plugging $\eta_n = \sqrt{\frac{K(\ln K + 1)}{4nL^2}}$ into (i). and (ii). directly concludes the proof. \square

The following lemma is a result that can be found in [CBL06, H⁺16], we rewrite it here for completeness.

Lemma 5. According to $\mathbf{b}^{(n)} = \mathbf{m}^* \left(\eta_n \sum_{m=1}^n \ell_m + \mathcal{Z}_1 \right)$, we can have

$$\forall \mathbf{x} \in \mathcal{M}, \quad \sum_{n=1}^N \langle \mathbf{b}^{(n)} - \mathbf{x}, \ell_n \rangle \leq \left\langle \frac{\mathbf{x}}{\eta_N}, \mathcal{Z}_1 \right\rangle - \left\langle \frac{\mathbf{b}^{(1)}}{\eta_1} + \sum_{n=2}^N \left(\frac{1}{\eta_n} - \frac{1}{\eta_{n-1}} \right) \mathbf{b}^{(n)}, \mathcal{Z}_1 \right\rangle. \quad (19)$$

Proof This is done by induction. For the base case, $N = 1$, as $\mathbf{b}^{(1)} = \mathbf{m}^*(\boldsymbol{\ell}_1 + \boldsymbol{\mathcal{Z}}_1/\eta_1)$

$$\left\langle \mathbf{b}^{(1)}, \boldsymbol{\ell}_1 + \boldsymbol{\mathcal{Z}}_1/\eta_1 \right\rangle \leq \left\langle \mathbf{x}, \boldsymbol{\ell}_1 + \boldsymbol{\mathcal{Z}}_1/\eta_1 \right\rangle$$

for any $\mathbf{x} \in \mathcal{M}$. A simple rearrangement yields (19). While considering $N + 1$, we suppose (19) holds for all integers smaller than $N + 1$. For an arbitrary $\mathbf{x} \in \mathcal{M}$, the fact $\mathbf{b}^{(N+1)} = \mathbf{m}^*\left(\sum_{n=1}^{N+1} \boldsymbol{\ell}_n + \boldsymbol{\mathcal{Z}}_1/\eta_{N+1}\right)$ directly implies that

$$\begin{aligned} \left\langle \mathbf{x}, \sum_{n=1}^{N+1} \boldsymbol{\ell}_n + \frac{\boldsymbol{\mathcal{Z}}_1}{\eta_{N+1}} \right\rangle &\geq \left\langle \mathbf{b}^{(N+1)}, \sum_{n=1}^{N+1} \boldsymbol{\ell}_n + \frac{\boldsymbol{\mathcal{Z}}_1}{\eta_{N+1}} \right\rangle \\ &= \left\langle \mathbf{b}^{(N+1)}, \boldsymbol{\ell}_{N+1} + \left(\frac{1}{\eta_{N+1}} - \frac{1}{\eta_N}\right) \boldsymbol{\mathcal{Z}}_1 \right\rangle + \left\langle \mathbf{b}^{(N+1)}, \sum_{n=1}^N \boldsymbol{\ell}_n + \frac{\boldsymbol{\mathcal{Z}}_1}{\eta_N} \right\rangle \\ &\geq \sum_{n=1}^{N+1} \left\langle \mathbf{b}^{(n)}, \boldsymbol{\ell}_n \right\rangle + \left\langle \frac{\mathbf{b}^{(1)}}{\eta_1} + \sum_{n=2}^{N+1} \left(\frac{1}{\eta_n} - \frac{1}{\eta_{n-1}}\right) \mathbf{b}^{(n)}, \boldsymbol{\mathcal{Z}}_1 \right\rangle, \end{aligned}$$

where the last inequality comes from applying the hypothesis (19) with $\mathbf{x} = \mathbf{b}^{(N+1)}$ on the second inner product. Rearrange the above inequality, our induction is completed. \square

Lemma 6 (Hoeffding-Azuma). *Let $N \in \mathbb{N}$, V_1, V_2, \dots, V_N be a bounded martingale difference sequence w.r.t. X_1, X_2, \dots, X_N such that for any $n \in [N]$ $V_n \in [A_n, A_n + r_n]$ for some random variable A_n , measurable w.r.t. X_1, \dots, X_{n-1} and a positive constant r_n . Then, for any $s > 0$,*

$$\mathbb{P} \left[\sum_{n \in [N]} V_n > s \right] \leq \exp \left(-\frac{2s^2}{\sum_{n \in [N]} r_n^2} \right) \quad \text{and} \quad \mathbb{P} \left[\sum_{n \in [N]} V_n < -s \right] \leq \exp \left(-\frac{2s^2}{\sum_{n \in [N]} r_n^2} \right).$$

D Analysis of P-FWS

In this appendix, we prove our main theorem.

Theorem 4. *Let $\mu \in \Lambda$ and $\delta \in (0, 1)$. If P-FWS is parametrized using*

$$(\epsilon_t, \eta_t, n_t, \rho_t, \theta_t) = \left(t^{-\frac{1}{5}}, \frac{1}{4\sqrt{t|\mathcal{X}_0|}}, \lceil t^{\frac{1}{4}} \rceil, \frac{1}{16tD^2|\mathcal{X}_0|}, \frac{1}{t^{\frac{1}{4}}e^{\sqrt{t}}} \right), \quad (10)$$

then (i) the algorithm finishes in finite time almost surely and $\mathbb{P}_\mu[\hat{\mathbf{i}} \neq \mathbf{i}^*(\mu)] \leq \delta$; (ii) its sample complexity satisfies $\mathbb{P}_\mu[\limsup_{\delta \rightarrow 0} \frac{\tau}{\ln \frac{1}{\delta}} \leq T^*(\mu)] = 1$ and for any $\epsilon, \tilde{\epsilon} \in (0, 1)$ with $\epsilon < \min\{1, \frac{2D^2\Delta_{\min}^2}{K}, \frac{D^2\|\mu\|_\infty^2}{3}\}$,

$$\mathbb{E}_\mu[\tau] \leq \frac{(1 + \tilde{\epsilon})^2}{T^*(\mu)^{-1} - 6\epsilon} \times H\left(\frac{1}{\delta} \cdot \frac{4c_2}{3} \cdot \frac{(1 + \tilde{\epsilon})^2}{T^*(\mu)^{-1} - 6\epsilon}\right) + \Psi(\epsilon, \tilde{\epsilon}),$$

where $H(x) = \ln x + \ln \ln x + 1$ and $\Psi(\epsilon, \tilde{\epsilon})$ (refer to (34) for a detailed expression) is polynomial in ϵ^{-1} , $\tilde{\epsilon}^{-1}$, K , $\|\mu\|_\infty$ and $\Delta_{\min}(\mu)^{-1}$; (iii) the expected number of LM Oracle calls is upper bounded by a polynomial in $\ln \delta^{-1}$, K , $\|\mu\|_\infty$ and $\Delta_{\min}(\mu)^{-1}$.

D.1 δ -correctness (Theorem 4 (i))

Recall that P-FWS stopping rule is:

$$\tau = \inf \left\{ t > 4|\mathcal{X}_0| : \frac{t\hat{F}_t}{1 + \epsilon_t} > \beta\left(t, \left(1 - \frac{1}{4|\mathcal{X}_0|}\right)\delta\right), \max\left\{\frac{1}{\Delta_{\min}(\hat{\mu}(t))}, \|\hat{\mu}(t)\|_\infty\right\} \leq \sqrt{t} \right\}, \quad (7)$$

where \hat{F}_t is computed by $(\epsilon_t, \delta/t^2)$ -MCP($\hat{\omega}(t), \hat{\mu}(t)$). Let $\hat{\mathbf{i}} = \mathbf{i}^*(\hat{\mu}(\tau))$ be the output of P-FWS. Define the good event $\mathcal{G} = \bigcap_{t=4|\mathcal{X}_0|+1}^\infty \{\hat{F}_t \leq (1 + \epsilon_t)F_{\hat{\mu}(t)}(\hat{\omega}(t))\}$. Hence, it follows from the guarantee of $(\epsilon_t, \delta/t^2)$ -MCP algorithm that

$$\mathbb{P}_\mu[\mathcal{G}^c] \leq \delta \sum_{t=4|\mathcal{X}_0|+1}^\infty t^{-2} \leq \delta \int_{4|\mathcal{X}_0|}^\infty x^{-2} dx \leq \frac{\delta}{4|\mathcal{X}_0|}.$$

Besides, under the event \mathcal{G} ,

$$(1 + \epsilon_\tau)\tau F_{\hat{\mu}(\tau)}(\hat{\omega}(\tau)) \geq \tau\hat{F}_\tau \geq (1 + \epsilon_\tau)\beta\left(\tau, \left(1 - \frac{1}{4|\mathcal{X}_0|}\right)\delta\right)$$

holds, implying that $\tau F_{\hat{\mu}(\tau)}(\hat{\omega}(\tau)) \geq \beta\left(\tau, \left(1 - \frac{1}{4|\mathcal{X}_0|}\right)\delta\right)$. So, by (8)-(9), $\hat{\mathbf{i}} = \mathbf{i}^*(\hat{\mu}(t))$ satisfies:

$$\mathbb{P}_\mu[\hat{\mathbf{i}} \neq \mathbf{i}^*(\mu), \mathcal{G}] \leq \left(1 - \frac{1}{4|\mathcal{X}_0|}\right)\delta,$$

and thus $\mathbb{P}_\mu[\hat{\mathbf{i}} \neq \mathbf{i}^*(\mu)] \leq \mathbb{P}_\mu[\hat{\mathbf{i}} \neq \mathbf{i}^*(\mu), \mathcal{G}] + \mathbb{P}_\mu[\mathcal{G}^c] \leq \delta$.

D.2 Almost-sure upper bound (Theorem 4 (ii))

In this section, we show Theorem 4 (ii) an almost-sure upper bound on the sample complexity for P-FWS. Our proof is based on the continuity of F_μ in μ (as in [GK16, WTP21]) and also on the following observations:

- (a) $\{\hat{\mu}(t) \xrightarrow{t \rightarrow \infty} \mu\}$ and $\{\nabla \tilde{F}_{\mu, \eta_t, n_t}(\omega) \xrightarrow{t \rightarrow \infty} \nabla \bar{F}_{\mu, \eta_t}(\omega), \forall \omega \in \Sigma_+\}$ happen almost surely,
- (b) $\hat{F}_t \geq F_{\hat{\mu}(t)}(\hat{\omega}(t))$.

For (a), by the law of large numbers, $\hat{\mu}(t) \xrightarrow{t \rightarrow \infty} \mu$ as $N_k(t) \xrightarrow{t \rightarrow \infty} \infty$ for all $k \in [K]$ yielded by forced exploration rounds involved in P-FWS (Lemma 14 in Appendix F), $\nabla \tilde{F}_{\mu, \eta_t, n_t}(\omega) \xrightarrow{t \rightarrow \infty}$

$\nabla \bar{F}_{\mu, \eta_t}(\omega)$, $\forall \omega \in \Sigma_+$ is a direct consequence that $n_t \xrightarrow{t \rightarrow \infty} \infty$. (b) is immediately derived from the definition of \hat{F}_t as $\hat{F}_t = f_{\hat{x}}(\hat{\omega}(t), \hat{\mu}(t))$ for some action $\hat{x} \neq i^*(\hat{\mu}(t))$ and $F_{\hat{\mu}(t)}(\hat{\omega}(t)) = \min_{x \in \mathcal{X} \setminus i^*(\hat{\mu}(t))} f_x(\hat{\omega}(t), \hat{\mu}(t))$.

Introduce the event

$$\mathcal{E} = \left\{ F_{\mu}(\hat{\omega}(t)) \xrightarrow{t \rightarrow \infty} \max_{\omega \in \Sigma} F_{\mu}(\omega) \text{ and } \hat{\mu}(t) \xrightarrow{t \rightarrow \infty} \mu \right\}.$$

Because of (a), Theorem 5 in Appendix E ensures that $\mathbb{P}_{\mu}[\mathcal{E}] = 1$. Also, by the uniform continuity of $F_{\mu}(\omega)$ in μ for an arbitrary $\omega \in \Sigma_+$ (Lemma 7 in D.3.3),

$$\max_{\omega \in \Sigma_+} |F_{\hat{\mu}(t)}(\omega) - F_{\mu}(\omega)| \xrightarrow{t \rightarrow \infty} 0$$

almost surely, and hence by the triangle inequality, this implies that

$$\mathbb{P}_{\mu} \left[F_{\hat{\mu}(t)}(\hat{\omega}(t)) \xrightarrow{t \rightarrow \infty} \max_{\omega \in \Sigma} F_{\mu}(\omega) \right] = 1.$$

For any $\epsilon \in (0, 1)$, under \mathcal{E} , there exists a positive integer $T_{\epsilon} > \max\{c_1, 4|\mathcal{X}_0|\}$ such that for any $t \geq T_{\epsilon}$, we have

$$F_{\hat{\mu}(t)}(\hat{\omega}(t)) \geq (1 - \epsilon) \max_{\omega \in \Sigma} F_{\mu}(\omega), \quad \max \left\{ \frac{1}{\Delta_{\min}(\hat{\mu}(t))}, \|\hat{\mu}(t)\|_{\infty} \right\} \leq \sqrt{t}, \text{ and } \epsilon_t \leq \epsilon, \quad (20)$$

where the second inequality is due to (a) and the third is because $\epsilon_t \rightarrow 0$. So, the stopping time (7) can be upper bounded by

$$\begin{aligned} \tau &\leq T_{\epsilon} + \inf \left\{ t > T_{\epsilon} : t\hat{F}_t > (1 + \epsilon) \beta \left(t, \frac{(4|\mathcal{X}_0| - 1)\delta}{4|\mathcal{X}_0|} \right) \right\} \\ &\leq T_{\epsilon} + \inf \left\{ t > T_{\epsilon} : tF_{\hat{\mu}(t)}(\hat{\omega}(t)) > (1 + \epsilon) \beta \left(t, \frac{(4|\mathcal{X}_0| - 1)\delta}{4|\mathcal{X}_0|} \right) \right\} \\ &\leq T_{\epsilon} + \inf \left\{ t > T_{\epsilon} : t(1 - \epsilon) \max_{\omega \in \Sigma} F_{\mu}(\omega) > (1 + \epsilon) \beta \left(t, \frac{(4|\mathcal{X}_0| - 1)\delta}{4|\mathcal{X}_0|} \right) \right\} \\ &\leq T_{\epsilon} + \inf \left\{ t > T_{\epsilon} : \frac{(1 - \epsilon)t}{(1 + \epsilon)T^*(\mu)} > \ln \left(\frac{c_2 t}{\delta} \cdot \frac{4|\mathcal{X}_0|}{4|\mathcal{X}_0| - 1} \right) \right\} \\ &\leq 2T_{\epsilon} + \left(\frac{1 + \epsilon}{1 - \epsilon} \right) T^*(\mu) H \left(\frac{1}{\delta} \cdot \frac{8c_2}{7} \left(\frac{1 + \epsilon}{1 - \epsilon} \right) T^*(\mu) \right). \end{aligned} \quad (21)$$

where the first inequality uses the last two inequalities of (20), the second inequality uses (b), the third inequality is based on the first inequality of (20), the fourth uses $T^*(\mu)^{-1} = \max_{\omega \in \Sigma} F_{\mu}(\omega)$ and (9), and the last inequality results from $(4|\mathcal{X}_0|)/(4|\mathcal{X}_0| - 1) \leq 8/7$ (as $|\mathcal{X}_0| \geq 2$), and an application of Lemma 9 with

$$\alpha = 1, \quad b_1 = \frac{1 - \epsilon}{1 + \epsilon} \cdot \frac{1}{T^*(\mu)} \quad \text{and} \quad b_2 = \frac{8c_2}{7} \cdot \frac{1}{\delta}.$$

Finally, as $\epsilon \in (0, 1)$ can be arbitrarily small, (21) implies that

$$\mathbb{P}_{\mu} \left[\limsup_{\delta \rightarrow 0} \frac{\tau}{\ln \delta^{-1}} \leq T^*(\mu) \right] = 1.$$

D.3 Non-asymptotic sample complexity (Theorem 4 (ii))

We establish the following non-asymptotic upper bound on $\mathbb{E}_{\mu}[\tau]$: for any $\tilde{\epsilon}, \epsilon \in (0, 1)$ small enough,

$$\mathbb{E}_{\mu}[\tau] \leq \frac{(1 + \tilde{\epsilon})^2}{T^*(\mu)^{-1} - 6\epsilon} H \left(\frac{1}{\delta} \cdot \frac{8c_2}{7} \cdot \frac{(1 + \tilde{\epsilon})^2}{T^*(\mu)^{-1} - 6\epsilon} \right) + \Psi(\epsilon, \tilde{\epsilon}),$$

where $H(x) = \ln x + \ln \ln x + 1$ and $\Psi(\epsilon, \tilde{\epsilon})$ is defined in (34).

Note that this directly implies the asymptotic optimality. Indeed, when $\delta \rightarrow 0$, we get:

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\boldsymbol{\mu}}[\tau]}{\ln \delta^{-1}} \leq \frac{(1 + \tilde{\epsilon})^2}{T^*(\boldsymbol{\mu})^{-1} - 6\epsilon}.$$

As $\epsilon, \tilde{\epsilon}$ can be set arbitrarily small and $\text{kl}(\delta, 1 - \delta) \approx \ln \delta^{-1}$ as $\delta \rightarrow 0$, it matches the sample complexity lower bound (1) (Theorem 7 in Appendix K) asymptotically.

Throughout this section, we assume $\boldsymbol{\mu} \in \Lambda$ is given and take any $\epsilon \in (0, 1)$ satisfying the following:

$$\epsilon < \min \left\{ 1, \frac{2D^2 \Delta_{\min}^2}{K}, \frac{1}{6T^*(\boldsymbol{\mu})} \right\} \leq \min \left\{ 1, \frac{2D^2 \Delta_{\min}^2}{K}, \frac{D^2 \|\boldsymbol{\mu}\|_{\infty}^2}{3} \right\}, \quad (22)$$

where the second inequality is because $T^*(\boldsymbol{\mu})^{-1} \leq \ell = 2D^2 \|\boldsymbol{\mu}\|_{\infty}^2$ by Lemma 22 in Appendix I. The assumption of $\epsilon < \min\{1, \frac{2D^2 \Delta_{\min}^2}{K}, \frac{D^2 \|\boldsymbol{\mu}\|_{\infty}^2}{3}\}$ is used to define the good events introduced in D.3.1 as well as to derive several necessary technical lemmas summarized in D.3.3.

D.3.1 Good events

Since in early rounds, the estimation of $\hat{\boldsymbol{\mu}}(t)$ is noisy, we introduce two threshold functions, \underline{h} and \bar{h} , on the round index T :

$$\begin{cases} \underline{h}(T) &= \min\{t \in \mathbb{N} : t \geq T^a, \sqrt{t/|\mathcal{X}_0|} \in \mathbb{N}\} \\ \bar{h}(T) &= \min\{t \in \mathbb{N} : t \geq T^b \underline{h}(T), \sqrt{t/|\mathcal{X}_0|} \in \mathbb{N}\}, \end{cases} \quad (23)$$

where $a, b \in (0, 1)$ and $a + b < 1$ will be explained later in (27). Now, we define our good events:

$$\mathcal{E}_{1,\epsilon}(T) = \bigcap_{t=\underline{h}(T)}^T \mathcal{E}_{1,\epsilon}^{(t)} \quad \text{and} \quad \mathcal{E}_{2,\epsilon}(T) = \bigcap_{t=\underline{h}(T)}^T \mathcal{E}_{2,\epsilon}^{(t)}, \quad (24)$$

where $\mathcal{E}_{1,\epsilon}^{(t)} = \{\langle \nabla \bar{F}_{\hat{\boldsymbol{\mu}}(t-1), \eta_t}(\hat{\boldsymbol{\omega}}(t-1)), \mathbf{x}(t) \rangle \geq \max_{\mathbf{x} \in \mathcal{X}} \langle \nabla \bar{F}_{\hat{\boldsymbol{\mu}}(t-1), \eta_t}(\hat{\boldsymbol{\omega}}(t-1)), \mathbf{x} \rangle - \epsilon\}$ and $\mathcal{E}_{2,\epsilon}^{(t)} = \{\|\hat{\boldsymbol{\mu}}(t-1) - \boldsymbol{\mu}\|_{\infty} < \frac{\epsilon}{24D^3 \|\boldsymbol{\mu}\|_{\infty}}\}$.

$\mathcal{E}_{1,\epsilon}^{(t)}$ is the event when the solution of FW update is bounded by at most ϵ , and $\mathcal{E}_{2,\epsilon}^{(t)}$ is the event when the empirical estimate of $\boldsymbol{\mu}$ is sufficiently accurate. Under $\mathcal{E}_{2,\epsilon}^{(t)}$, the uniform continuity shown in Lemma 7 in D.3.3 ensures that:

$$\begin{aligned} |F_{\hat{\boldsymbol{\mu}}(t-1)}(\boldsymbol{\omega}) - F_{\boldsymbol{\mu}}(\boldsymbol{\omega})| &< \epsilon, \quad \forall \boldsymbol{\omega} \in \Sigma_+, \\ |\langle \nabla \bar{F}_{\hat{\boldsymbol{\mu}}(t-1), \eta}(\boldsymbol{\omega}) - \nabla \bar{F}_{\boldsymbol{\mu}, \eta}(\boldsymbol{\omega}), \mathbf{x} - \boldsymbol{\omega} \rangle| &< \epsilon, \quad \forall (\boldsymbol{\omega}, \mathbf{x}) \in \Sigma_+ \times \mathcal{X}, \forall \eta \in (0, \min_{k \in [K]} \omega_k). \end{aligned}$$

The second inequality enables the duality gap of FW algorithm to be controlled, leading to the convergence of P-FWS. Let

$$\begin{aligned} M &= \max \left\{ (4|\mathcal{X}_0|)^{\frac{1}{a}}, \left(\frac{4K^2}{\epsilon^2 D^2 |\mathcal{X}_0|} \right)^{\frac{1}{a}}, \left(\frac{2}{\Delta_{\min}(\boldsymbol{\mu})} \right)^{\frac{2}{a}}, \left(\frac{3 \|\boldsymbol{\mu}\|_{\infty}}{2} \right)^{\frac{2}{a}} \right\} \\ &\quad + \max \left\{ \left(\frac{\ell}{\epsilon} \right)^{\frac{1}{b}}, \left(\frac{5\ell K^2}{\epsilon \sqrt{|\mathcal{X}_0|}} \right)^{\frac{2}{a+b}} \right\}, \end{aligned} \quad (25)$$

then overall, we have (Theorem 5 in D.3.3): for any $t \geq \bar{h}(M)$,

$$\max_{\boldsymbol{\omega} \in \Sigma} F_{\boldsymbol{\mu}}(\boldsymbol{\omega}) - F_{\boldsymbol{\mu}}(\hat{\boldsymbol{\omega}}(t)) \leq 5\epsilon, \quad \Delta_{\min}(\hat{\boldsymbol{\mu}}(t)) \geq \frac{\Delta_{\min}(\boldsymbol{\mu})}{2}, \quad \text{and} \quad \|\hat{\boldsymbol{\mu}}(t)\|_{\infty} \leq \frac{3 \|\boldsymbol{\mu}\|_{\infty}}{2}. \quad (26)$$

Finally, the values of a, b are set to the following:

$$a = \frac{7}{9} \quad \text{and} \quad b = \frac{1}{9}. \quad (27)$$

This choices will balance the leading order between ϵ^{-1} and $\tilde{\epsilon}^{-1}$ in the δ -independent terms (34) of the non-asymptotic upper bound (which will be shown later).

D.3.2 Proof of non-asymptotic sample complexity

Let $\delta \in (0, 1)$. We claim that:

$$\mathbb{E}_{\boldsymbol{\mu}}[\tau] \leq \sum_{T=1}^{\infty} \mathbb{P}_{\boldsymbol{\mu}}[\tau \geq T] \leq T_0(\delta) + \sum_{T=M+1}^{\infty} \mathbb{P}_{\boldsymbol{\mu}}[(\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T))^c], \quad (28)$$

where $T_0(\delta) = \inf \left\{ T \geq M : \bar{h}(T) + \frac{(1+\epsilon_T)}{T^*(\boldsymbol{\mu})^{-1}-6\epsilon} \beta \left(T, \frac{(4|\mathcal{X}_0|-1)\delta}{4|\mathcal{X}_0|} \right) \leq T \right\}$. The proof is completed by bounding each term in the right-hand side of (28).

Proof of (28): Suppose $T \geq M$ and $\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T)$ holds. Observe that

$$\min\{\tau, T\} \leq \bar{h}(T) + \sum_{t=\lceil \bar{h}(T) \rceil}^T \mathbb{1}\{\tau > t\}.$$

To derive an upper bound of $\sum_{t=\lceil \bar{h}(T) \rceil}^T \mathbb{1}\{\tau > t\}$, recall the stopping rule (7) that

$$\begin{aligned} \tau &= \inf \left\{ t > 4|\mathcal{X}_0| : \frac{t\hat{F}_t}{1+\epsilon_t} > \beta \left(t, \frac{(4|\mathcal{X}_0|-1)\delta}{4|\mathcal{X}_0|} \right), \max \left\{ \frac{1}{\Delta_{\min}(\hat{\boldsymbol{\mu}}(t))}, \|\hat{\boldsymbol{\mu}}(t)\|_{\infty} \right\} \leq \sqrt{t} \right\} \\ &\leq \inf \left\{ t \geq \bar{h}(M) : t\hat{F}_t > (1+\epsilon_t) \beta \left(t, \frac{(4|\mathcal{X}_0|-1)\delta}{4|\mathcal{X}_0|} \right) \right\} \\ &\leq \inf \left\{ t \geq \bar{h}(M) : t(T^*(\boldsymbol{\mu})^{-1} - 6\epsilon) > (1+\epsilon_t) \beta \left(t, \frac{(4|\mathcal{X}_0|-1)\delta}{4|\mathcal{X}_0|} \right) \right\}, \end{aligned}$$

where the first inequality uses (26), and the second follows from Lemma 7 and Theorem 5 in D.3.3:

$$|F_{\hat{\boldsymbol{\mu}}(t-1)}(\hat{\boldsymbol{\omega}}(t-1)) - F_{\boldsymbol{\mu}}(\hat{\boldsymbol{\omega}}(t-1))| < \epsilon \quad \text{and} \quad T^*(\boldsymbol{\mu})^{-1} - F_{\boldsymbol{\mu}}(\hat{\boldsymbol{\omega}}(t)) \leq 5\epsilon, \quad (29)$$

and the fact that $\hat{F}_t \geq F_{\hat{\boldsymbol{\mu}}(t)}(\hat{\boldsymbol{\omega}}(t))$. Hence, $\sum_{t=\lceil \bar{h}(T) \rceil}^T \mathbb{1}\{\tau > t\}$ is upper bounded by

$$\sum_{t=\lceil \bar{h}(T) \rceil}^T \mathbb{1} \left\{ t(T^*(\boldsymbol{\mu})^{-1} - 6\epsilon) \leq (1+\epsilon_t) \beta \left(t, \frac{(4|\mathcal{X}_0|-1)\delta}{4|\mathcal{X}_0|} \right) \right\} \leq \frac{(1+\epsilon_T)}{T^*(\boldsymbol{\mu})^{-1} - 6\epsilon} \beta \left(T, \frac{(4|\mathcal{X}_0|-1)\delta}{4|\mathcal{X}_0|} \right).$$

By defining $T_0(\delta)$ as done in (28), we get (28), i.e.,

$$\mathbb{E}_{\boldsymbol{\mu}}[\tau] \leq \sum_{T=1}^{\infty} \mathbb{P}_{\boldsymbol{\mu}}[\tau \geq T] \leq T_0(\delta) + \sum_{T=M+1}^{\infty} \mathbb{P}_{\boldsymbol{\mu}}[(\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T))^c]$$

because $\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T) \subseteq \{\tau \leq T\}$ for any $T \geq T_0(\delta)$.

Now, we proceed with the proof by upper-bounding each term in the right-hand side of (28).

Bounding $T_0(\delta)$: Introduce $\tilde{\epsilon} \in (0, 1)$ that can be chosen arbitrarily small. Notice that

$$T - \bar{h}(T) = T - T^{a+b} \geq \frac{T}{1+\tilde{\epsilon}} \quad \text{when} \quad T \geq \left(1 + \frac{1}{\tilde{\epsilon}}\right)^{\frac{1}{1-(a+b)}}, \quad (30)$$

$$\epsilon_T = T^{-\frac{1}{9}} \leq \tilde{\epsilon} \quad \text{when} \quad T \geq \left(\frac{1}{\tilde{\epsilon}}\right)^9, \quad (31)$$

where the first inequality results from a simple rearrangement, and the second substitutes $\epsilon_t = t^{-1/9}$. Then, it follows from (9) that:

$$\begin{aligned} T_0(\delta) &\leq \inf \left\{ T \geq \max \left\{ M, (1+\tilde{\epsilon}^{-1})^{\frac{1}{1-(a+b)}}, \tilde{\epsilon}^{-9} \right\} : \frac{(1+\tilde{\epsilon}) \beta \left(T, \frac{3\delta}{4} \right)}{T^*(\boldsymbol{\mu})^{-1} - 6\epsilon} \leq \frac{T}{1+\tilde{\epsilon}} \right\} \\ &\leq \inf \left\{ T \geq \max \left\{ M, (1+\tilde{\epsilon}^{-1})^{\frac{1}{1-(a+b)}}, \tilde{\epsilon}^{-9}, c_1 \right\} : \ln \left(\frac{4c_2 T}{3\delta} \right) \leq \frac{T^*(\boldsymbol{\mu})^{-1} - 6\epsilon}{(1+\tilde{\epsilon})^2} \cdot T \right\} \\ &\leq \max \left\{ M, (1+\tilde{\epsilon}^{-1})^{\frac{1}{1-(a+b)}}, \tilde{\epsilon}^{-9}, c_1 \right\} + \frac{(1+\tilde{\epsilon})^2}{T^*(\boldsymbol{\mu})^{-1} - 6\epsilon} \times H \left(\frac{4c_2}{3\delta} \cdot \frac{(1+\tilde{\epsilon})^2}{T^*(\boldsymbol{\mu})^{-1} - 6\epsilon} \right), \end{aligned} \quad (32)$$

where the first inequality uses (30)-(31) and $\frac{4|\mathcal{X}_0|-1}{4|\mathcal{X}_0|} \geq \frac{3}{4}$ (as $|\mathcal{X}_0| \geq 2$ is shown in Lemma 23 in Appendix J), the second inequality is due to (9), and the last results from an application of Lemma 9 in Appendix D.3.3 with

$$\alpha = 1, \quad b_1 = \frac{T^*(\boldsymbol{\mu})^{-1} - 6\epsilon}{(1 + \tilde{\epsilon})^2}, \quad \text{and} \quad b_2 = \frac{4c_2}{3\delta}.$$

Bounding $\sum_{T=M+1}^{\infty} \mathbb{P}_{\boldsymbol{\mu}}[(\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T))^c]$: By Lemma 8 in Appendix D.3.3, it is upper bounded by

$$2K \left(\frac{3}{\min\{1, \frac{\epsilon^2}{8\ell^2 K^3 D^2}\}} \right)^{2+\frac{2}{a}} + 2 \left(\frac{2304D^6 \|\boldsymbol{\mu}\|_{\infty}^2 \sqrt{|\mathcal{X}_0|}}{\epsilon^2} \right)^{2+\frac{2}{a}} \Gamma \left(2 + \frac{2}{a} \right). \quad (33)$$

Putting things together: Finally, substituting $(a, b) = (\frac{7}{9}, \frac{1}{9})$ into (25)-(32)-(33) yields that:

- $T_0(\delta) \leq M + (1 + \frac{1}{\tilde{\epsilon}})^9 + (\frac{1}{\tilde{\epsilon}})^9 + c_1 + \frac{(1+\tilde{\epsilon})^2}{T^*(\boldsymbol{\mu})^{-1}-6\epsilon} \times H \left(\frac{4c_2}{3\delta} \cdot \frac{(1+\tilde{\epsilon})^2}{T^*(\boldsymbol{\mu})^{-1}-6\epsilon} \right)$
- $M \leq \max\{(4|\mathcal{X}_0|)^{\frac{9}{7}}, (\frac{4K^2}{\epsilon^2 D^2 |\mathcal{X}_0|})^{\frac{9}{7}}, (\frac{4}{\Delta_{\min}^2})^{\frac{9}{7}}, (\frac{9\|\boldsymbol{\mu}\|_{\infty}^2}{4})^{\frac{9}{7}}\} + \max\{(\frac{\ell}{\epsilon})^9, (\frac{5\ell K^2}{\epsilon \sqrt{|\mathcal{X}_0|}})^{2.25}\}$
- $\sum_{T=M}^{\infty} \mathbb{P}_{\boldsymbol{\mu}}[(\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T))^c] < \frac{78K}{\epsilon^{10}} \left(2^{15} K^{15} D^{10} \ell^{10} + 2^{41} 3^9 D^{30} \|\boldsymbol{\mu}\|_{\infty}^{10} |\mathcal{X}_0|^{2.5} \right)$

where simplifications are obtained remarking that $\Gamma(2 + \frac{2}{a}) \leq 13$ and $2 + \frac{2}{a} < 5$. Therefore, substituting $\ell = 2D^2 \|\boldsymbol{\mu}\|_{\infty}^2$ (defined in Appendix I) and $78 < 2^7$, $3^9 \leq 2^{15}$, and $4^{9/7} < 6$ lead to:

$$\mathbb{E}_{\boldsymbol{\mu}}[\tau] \leq \frac{(1 + \tilde{\epsilon})^2}{T^*(\boldsymbol{\mu})^{-1} - 6\epsilon} \times H \left(\frac{4c_2}{3\delta} \cdot \frac{(1 + \tilde{\epsilon})^2}{T^*(\boldsymbol{\mu})^{-1} - 6\epsilon} \right) + \Psi(\epsilon, \tilde{\epsilon}),$$

where

$$\begin{aligned} \Psi(\epsilon, \tilde{\epsilon}) = & 6 \max \left\{ |\mathcal{X}_0|, \frac{K^2}{\epsilon^2 D^2 |\mathcal{X}_0|}, \frac{1}{\Delta_{\min}^2}, \|\boldsymbol{\mu}\|_{\infty}^2 \right\}^{\frac{9}{7}} + \max \left\{ \frac{2^9 D^{18} \|\boldsymbol{\mu}\|_{\infty}^{18}}{\epsilon^9}, \frac{10^{2.25} D^{4.5} \|\boldsymbol{\mu}\|_{\infty}^{4.5} K^2}{\epsilon^{2.25} |\mathcal{X}_0|^{1.125}} \right\} \\ & + \left(1 + \frac{1}{\tilde{\epsilon}} \right)^9 + \left(\frac{1}{\tilde{\epsilon}} \right)^9 + c_1 + \frac{2^{32} K D^{30} \|\boldsymbol{\mu}\|_{\infty}^{10} \left(K^{15} \|\boldsymbol{\mu}\|_{\infty}^{10} + 2^{31} |\mathcal{X}_0|^{2.5} \right)}{\epsilon^{10}}. \end{aligned} \quad (34)$$

D.3.3 Technical lemmas

The most important step in Theorem 4 is to bound the term $\sum_{T=M+1}^{\infty} \mathbb{P}_{\boldsymbol{\mu}}[(\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T))^c]$ in (28) explicitly in terms of K , $\|\boldsymbol{\mu}\|_{\infty}$ and ϵ . For this purpose, inspired by Assumption 3 in [WTP21], we developed Proposition 4 (see Appendix G.2 for the proof) and combine it with the mean-valued theorem to derive our main continuity results in Lemma 7 (see Appendix G for the proof). Throughout this section, we fix $\boldsymbol{\mu} \in \Lambda$ and denote $\Delta_{\min}(\boldsymbol{\mu})$ by Δ_{\min} .

Lemma 7. *Let $\epsilon \in (0, \frac{2D^2 \Delta_{\min}^2}{K})$. Then, any $\boldsymbol{\pi} \in \mathbb{R}^K$ with $\|\boldsymbol{\pi} - \boldsymbol{\mu}\|_{\infty} < \frac{\epsilon}{24D^3 \|\boldsymbol{\mu}\|_{\infty}}$ satisfies the following:*

$$|F_{\boldsymbol{\mu}}(\boldsymbol{\omega}) - F_{\boldsymbol{\pi}}(\boldsymbol{\omega})| < \epsilon, \quad \forall \boldsymbol{\omega} \in \Sigma_+ \quad (35)$$

$$|\langle \nabla \bar{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}) - \nabla \bar{F}_{\boldsymbol{\mu}, \eta}(\boldsymbol{\omega}), \boldsymbol{x} - \boldsymbol{\omega} \rangle| < \epsilon, \quad \forall (\boldsymbol{\omega}, \boldsymbol{x}) \in \Sigma_+ \times \mathcal{X}, \forall \eta \in (0, \min_{k \in [K]} \omega_k). \quad (36)$$

Our main concentration result with error specified explicitly in terms of ϵ is (see Appendix F for the proof):

Lemma 8. *Let $\epsilon \in (0, \frac{2D^2 \Delta_{\min}^2}{K})$ and M be defined as in (25). Then,*

$$\sum_{T=M+1}^{\infty} \mathbb{P}_{\boldsymbol{\mu}}[(\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T))^c] < 2K \left(\frac{3}{A_1(\epsilon)^{2+\frac{2}{a}}} + \frac{2}{A_2(\epsilon)^{2+\frac{2}{a}}} \right) \Gamma \left(2 + \frac{2}{a} \right),$$

where $A_1(\epsilon) = \min\{1, \frac{\epsilon^2}{8\ell^2 K^3 D^2}\}$, $A_2(\epsilon) = \frac{\epsilon^2}{2304D^6 \|\boldsymbol{\mu}\|_{\infty}^2 \sqrt{|\mathcal{X}_0|}}$, and Γ denotes the gamma function $\Gamma(z) = \int_0^{\infty} t^{z-1} e^{-t} dt$ for any $z > 0$.

We remark that Lemma 8 sharpens a similar result, Lemma 2 in [WTP21], by a factor of e^K after performing a more careful analysis.

Under good events $\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T)$, we show Theorem 5, the convergence of P-FWS when $\hat{\boldsymbol{\mu}}(t)$ is replaced with $\boldsymbol{\mu}$. As shown in Appendix D.3.2, the extra error due to this replacement is controlled, thanks to Lemma 7.

Theorem 5. *Let $\epsilon \in (0, \min\{1, \frac{2D^2\Delta_{\min}^2}{K}\})$ and T be an integer at least larger than*

$$\max \left\{ (4|\mathcal{X}_0|)^{\frac{1}{a}}, \left(\frac{4K^2}{\epsilon^2 D^2 |\mathcal{X}_0|} \right)^{\frac{1}{a}}, \left(\frac{2}{\Delta_{\min}} \right)^{\frac{2}{a}}, \left(\frac{3\|\boldsymbol{\mu}\|_{\infty}}{2} \right)^{\frac{2}{a}} \right\} + \max \left\{ \left(\frac{\ell}{\epsilon} \right)^{\frac{1}{b}}, \left(\frac{5\ell K^2}{\epsilon \sqrt{|\mathcal{X}_0|}} \right)^{\frac{2}{a+b}} \right\}.$$

Under $\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T)$, Algorithm 2 with (10) satisfies that: for any $t = \bar{h}(T), \bar{h}(T) + 1, \dots, T$,

$$(i) \max_{\boldsymbol{\omega} \in \Sigma} F_{\boldsymbol{\mu}}(\boldsymbol{\omega}) - F_{\boldsymbol{\mu}}(\hat{\boldsymbol{\omega}}(t)) \leq 5\epsilon, (ii) \Delta_{\min}(\hat{\boldsymbol{\mu}}(t)) \geq \frac{\Delta_{\min}}{2}, \text{ and } (iii) \|\hat{\boldsymbol{\mu}}(t)\|_{\infty} \leq \frac{3\|\boldsymbol{\mu}\|_{\infty}}{2}.$$

The proof of Theorem 5 is given in Appendix E.

Finally, the last ingredient is Lemma 9.

Lemma 9 (Lemma 18 in [GK16]). *Let $\alpha \in [1, \frac{\epsilon}{2}]$ and $b_1, b_2 > 0$. Then,*

$$x = \frac{1}{b_1} \left(\ln \left(\frac{b_2 e}{b_1^\alpha} \right) + \ln \ln \left(\frac{b_2}{b_1^\alpha} \right) \right)$$

satisfies $b_1 x \geq \ln(b_2 x^\alpha)$.

D.4 Computational complexity (Theorem 4 (iii))

In this section, we analyze the computational complexity of P-FWS running with (10) in terms of the number of calls to LM Oracle. We will show that the expected number of LM Oracle calls is upper bounded by a polynomial in $\ln \delta^{-1}$, K , $\|\boldsymbol{\mu}\|_{\infty}$ and $\Delta_{\min}(\boldsymbol{\mu})^{-1}$.

Proof The construction of \mathcal{X}_0 and computation of $\hat{\boldsymbol{i}}$ merely takes $\mathcal{O}(KD)$ calls to LM Oracle. The overall complexity is dominated by the LM Oracle calls performed from $4|\mathcal{X}_0| + 1$ to round τ , analyzed as follows.

Per-round complexity: Fix $t \in \{4|\mathcal{X}_0| + 1, \dots, \tau\}$. Recall from P-FWS that the FW update in round t and the stopping rule in round $t - 1$ are computed only if:

$$\max\{\Delta_{\min}(\hat{\boldsymbol{\mu}}(t-1))^{-1}, \|\hat{\boldsymbol{\mu}}(t-1)\|_{\infty}\} \leq \sqrt{t-1}. \quad (37)$$

Otherwise, forced-exploration procedure is invoked. Verifying (37) takes at most $D + 1$ calls.⁶ The computation of \hat{F}_{t-1} and $\hat{\boldsymbol{i}}^*(\nabla \hat{F}_{\hat{\boldsymbol{\mu}}(t-1), \eta_t, n_t}(\hat{\boldsymbol{\omega}}(t-1)))$ by Theorem 3 in Appendix C.2 takes at most

$$\mathcal{O} \left(D + \frac{\|\hat{\boldsymbol{\mu}}(t-1)\|_{\infty}^4 \|\hat{\boldsymbol{\omega}}(t-1)\|_{\infty}^2 K^3 D^5 \ln K}{F_{\hat{\boldsymbol{\mu}}(t-1)}(\hat{\boldsymbol{\omega}}(t-1))^2} \left(\frac{\ln(t^2 \delta^{-1})}{\epsilon_t^2} + \frac{n_t \ln \theta_t^{-1}}{\rho_t^2} \right) \right) \quad (38)$$

calls to LM Oracle. To evaluate (38), we need a lower bound on $F_{\hat{\boldsymbol{\mu}}(t-1)}(\hat{\boldsymbol{\omega}}(t-1))$. By Proposition 1 (c) in Appendix C.1, one evaluates $F_{\hat{\boldsymbol{\mu}}(t-1)}(\hat{\boldsymbol{\omega}}(t-1))$ in closed-form:

$$F_{\hat{\boldsymbol{\mu}}(t-1)}(\hat{\boldsymbol{\omega}}(t-1)) = \min_{\boldsymbol{x} \neq \hat{\boldsymbol{i}}^*(\hat{\boldsymbol{\mu}}(t-1))} \frac{\Delta_{\boldsymbol{x}}(\hat{\boldsymbol{\mu}}(t-1))^2}{2 \langle \boldsymbol{x} \oplus \hat{\boldsymbol{i}}^*(\hat{\boldsymbol{\mu}}(t-1)), \hat{\boldsymbol{\omega}}(t-1)^{-1} \rangle} \geq \frac{\min_{k \in [K]} \hat{\omega}_k(t-1)}{4D(t-1)},$$

where the inequality results from (37) that $\Delta_{\min}(\hat{\boldsymbol{\mu}}(t-1)) \geq \frac{1}{\sqrt{t-1}}$, $\|\boldsymbol{x} \oplus \boldsymbol{x}'\|_1 \leq 2D$ for any $\boldsymbol{x}, \boldsymbol{x}' \in \mathcal{X}$, and $\langle \boldsymbol{y}, \boldsymbol{z} \rangle \leq \|\boldsymbol{y}\|_1 \|\boldsymbol{z}\|_{\infty}$ for any $\boldsymbol{y}, \boldsymbol{z} \in \mathbb{R}^K$. Further, combining with Lemma 14 (which states $\min_{k \in [K]} \hat{\omega}_k(t-1) \geq \frac{1}{2\sqrt{(t-1)|\mathcal{X}_0|}}$) in Appendix F yields

$$F_{\hat{\boldsymbol{\mu}}(t-1)}(\hat{\boldsymbol{\omega}}(t-1)) \geq \frac{1}{8D\sqrt{|\mathcal{X}_0|}(t-1)^{1.5}}. \quad (39)$$

⁶For any $\boldsymbol{\pi} \in \Lambda$, $\Delta_{\min}(\boldsymbol{\pi})$ requires to compute $\hat{\boldsymbol{i}}^*(\boldsymbol{\pi})$ and solve $\max_{\boldsymbol{x} \neq \hat{\boldsymbol{i}}^*(\boldsymbol{\pi})} \langle \boldsymbol{\pi}, \boldsymbol{x} \rangle$, where the latter requires at most D calls to the LM Oracle by Lemma 2 in § 2.2.

From (39), $\|\hat{\boldsymbol{\mu}}(t-1)\|_\infty \leq \sqrt{t-1}$, Lemma 14, and substituting the parameters (10) into (38), we know that the number of LM Oracle calls performed at any round $t \geq 4|\mathcal{X}_0| + 1$ is at most

$$\begin{aligned} & \mathcal{O}\left(t^6 |\mathcal{X}_0|^2 K^3 D^7 \ln K \left(\frac{\ln(t^2 \delta^{-1})}{\epsilon_t^2} + \frac{n_t \ln \theta_t^{-1}}{\rho_t^2} \right)\right) \\ &= \mathcal{O}\left(t^6 |\mathcal{X}_0|^2 K^3 D^7 \ln K \left(t^{\frac{2}{5}} \ln \left(\frac{t}{\delta} \right) + t^{2.75} D^4 |\mathcal{X}_0|^2 \right)\right) \\ &= \mathcal{O}\left(t^{8.75} \ln \left(\frac{t}{\delta} \right) |\mathcal{X}_0|^4 D^{11} K^3 \ln K\right). \end{aligned} \quad (40)$$

Overall complexity: Invoking Theorem 4 in D.3 with $\tilde{\epsilon} = 0.1$ and $\epsilon = \frac{1}{12T^*(\boldsymbol{\mu})}$ results in

$$\mathbb{E}_{\boldsymbol{\mu}}[\tau] = \mathcal{O}\left(T^*(\boldsymbol{\mu}) \ln \left(\frac{T^*(\boldsymbol{\mu})}{\delta} \right) + \frac{1}{\Delta_{\min}^{\frac{18}{7}}} + K^{16} D^{30} \|\boldsymbol{\mu}\|_\infty^{20} T^*(\boldsymbol{\mu})^{10}\right)$$

which after using $T^*(\boldsymbol{\mu}) \leq 4KD/\Delta_{\min}^2$ (Lemma 1 in §2.1) becomes

$$\mathbb{E}_{\boldsymbol{\mu}}[\tau] = \mathcal{O}\left(\frac{KD}{\Delta_{\min}^2} \ln \left(\frac{KD}{\delta \Delta_{\min}^2} \right) + \frac{K^{26} D^{40} \|\boldsymbol{\mu}\|_\infty^{20}}{\Delta_{\min}^{20}}\right). \quad (41)$$

Hence, by a summation of (40) over $t = 4|\mathcal{X}_0| + 1$ to $\mathbb{E}_{\boldsymbol{\mu}}[\tau]$, the expected total number of the LM Oracle calls is upper bounded by

$$\mathcal{O}\left(\mathbb{E}_{\boldsymbol{\mu}}[\tau]^{9.75} \ln \left(\frac{\mathbb{E}_{\boldsymbol{\mu}}[\tau]}{\delta} \right) |\mathcal{X}_0|^4 D^{11} K^3 \ln K\right), \quad (42)$$

where the inequality uses integral by parts $\int t^{8.75} \ln t dt = \mathcal{O}(t^{9.75} \ln t)$. Remind that $\max\{D, |\mathcal{X}_0|\} \leq K$. Thus, we conclude that (42) is bounded by a polynomial function in $\ln \delta^{-1}$, $\|\boldsymbol{\mu}\|_\infty$, Δ_{\min}^{-1} , and K (due to (41), $\mathbb{E}_{\boldsymbol{\mu}}[\tau]$ is bounded by a polynomial function in the same variables). \square

E Convergence of P-FWS under the good events

Throughout this section, we assume that $\boldsymbol{\mu}$ is fixed and drop $\boldsymbol{\mu}$ from the notation, e.g., $F = F_{\boldsymbol{\mu}}$, $\bar{F}_{\eta} = \bar{F}_{\boldsymbol{\mu}, \eta}$, $\hat{F}_{\eta, t} = \hat{F}_{\boldsymbol{\mu}, \eta, t}$, and $\Delta_{\min} = \Delta_{\min}(\boldsymbol{\mu})$. Also, we will use $\boldsymbol{\omega}^* \in \operatorname{argmax}_{\boldsymbol{\omega} \in \Sigma} F(\boldsymbol{\omega})$ to denote any optimal allocation and let $\hat{\boldsymbol{i}}^* = \hat{\boldsymbol{i}}^*(\boldsymbol{\mu})$. Recall that $\underline{h}(T) \geq T^a$ and $\bar{h}(T) \geq T^{a+b}$ is defined in (23) in Appendix D.3.1 for some $a, b \in (0, 1)$.

Theorem 5. *Let $\epsilon \in (0, \min\{1, \frac{2D^2\Delta_{\min}^2}{K}\})$ and T be an integer at least larger than*

$$\max \left\{ (4|\mathcal{X}_0|)^{\frac{1}{a}}, \left(\frac{4K^2}{\epsilon^2 D^2 |\mathcal{X}_0|} \right)^{\frac{1}{a}}, \left(\frac{2}{\Delta_{\min}} \right)^{\frac{2}{a}}, \left(\frac{3\|\boldsymbol{\mu}\|_{\infty}}{2} \right)^{\frac{2}{a}} \right\} + \max \left\{ \left(\frac{\ell}{\epsilon} \right)^{\frac{1}{b}}, \left(\frac{5\ell K^2}{\epsilon \sqrt{|\mathcal{X}_0|}} \right)^{\frac{2}{a+b}} \right\}.$$

Under $\mathcal{E}_{1, \epsilon}(T) \cap \mathcal{E}_{2, \epsilon}(T)$, Algorithm 2 with (10) satisfies that: for any $t = \bar{h}(T), \bar{h}(T) + 1, \dots, T$,

$$(i) F(\boldsymbol{\omega}^*) - F(\hat{\boldsymbol{\omega}}(t)) \leq 5\epsilon, (ii) \Delta_{\min}(\hat{\boldsymbol{\mu}}(t)) \geq \frac{\Delta_{\min}}{2}, \text{ and } (iii) \|\hat{\boldsymbol{\mu}}(t)\|_{\infty} \leq \frac{3\|\boldsymbol{\mu}\|_{\infty}}{2}.$$

Proof Fix arbitrary ϵ and T that satisfy the conditions in the statement, and suppose $\mathcal{E}_{1, \epsilon}(T) \cap \mathcal{E}_{2, \epsilon}(T)$ holds. (ii)(iii) directly follows from Lemma 10 (where one can verified that its assumption of Lemma 10 on T is satisfied). With (ii)(iii), the analysis of FW convergence will be greatly simplified as (ii)(iii) ensure that

$$\max \left\{ \frac{1}{\Delta_{\min}(\hat{\boldsymbol{\mu}}(t-1))}, \|\hat{\boldsymbol{\mu}}(t-1)\| \right\} \leq \sqrt{t-1}.$$

This means that the forced-exploration procedure will only be invoked by the condition of $\sqrt{t/|\mathcal{X}_0|}$ when $t \geq \underline{h}(T) = T^a$.

Proof of (i) $F(\boldsymbol{\omega}^*) - F(\hat{\boldsymbol{\omega}}(t)) \leq 5\epsilon$: Fix $t \geq \underline{h}(T)$. As mentioned above, for such t , the forced-exploration procedure will be invoked only when $\sqrt{t/|\mathcal{X}_0|} \in \mathbb{N}$. To specify the rounds performing FW updates, introduce $s(t) = \lfloor \sqrt{t/|\mathcal{X}_0|} \rfloor - 1$ and define

$$p(t) = (s(t)^2 + 1)|\mathcal{X}_0| \quad \text{and} \quad q(t) = (s(t) + 1)^2 |\mathcal{X}_0| - 1.$$

Notice that $p(t)$ and $q(t)$ are respectively the starting (including) and the ending (including) round of a successive FW update rounds with no forced exploration in between. Let $\phi_t = F(\boldsymbol{\omega}^*) - \bar{F}_{\eta_t}(\hat{\boldsymbol{\omega}}(t))$ be the error. By a careful analysis, we derive a recursive relationship satisfied by ϕ_t (Lemma 11):

$$\begin{cases} t\phi_t \leq (t - |\mathcal{X}_0|)\phi_{t-|\mathcal{X}_0|} + 2\ell\sqrt{D}|\mathcal{X}_0|^2 & \text{if } t = p(t) - 1, \\ t\phi_t \leq (t-1)\phi_{t-1} + 3\epsilon + \ell \left(\eta_{t-1} + \frac{K^2}{2t\eta_t} \right) & \text{if } t \in [p(t), q(t)]. \end{cases} \quad (43)$$

The first case (in round $t = p(t) - 1$) is exactly the ending round of a forced-exploration procedure (from $t - |\mathcal{X}_0|, \dots, t$), and the second case (in round $t \in [p(t), q(t)]$) is a FW-update round. By repeatedly applying (43), we have

$$\begin{aligned} \bar{h}(T)\phi_{\bar{h}(T)} &\leq \underline{h}(T)\phi_{\underline{h}(T)} + 2\ell\sqrt{D}|\mathcal{X}_0|^2 (s(\bar{h}(T)) - s(\underline{h}(T))) + 3 \sum_{t=\underline{h}(T)}^{\bar{h}(T)} \left(\frac{\ell K^2}{\sqrt{t}|\mathcal{X}_0|} + \epsilon \right) \\ &\leq \underline{h}(T)\ell + \ell \left(2\sqrt{D}|\mathcal{X}_0|^{1.5} + \frac{3K^2}{\sqrt{|\mathcal{X}_0|}} \right) \left(\sqrt{\bar{h}(T)} - \sqrt{\underline{h}(T)} \right) + 3\epsilon(\bar{h}(T) - \underline{h}(T)), \end{aligned}$$

where the second inequality follows from $\phi_{\bar{h}(T)} \leq \max_{\boldsymbol{\omega} \in \Sigma} F_{\boldsymbol{\mu}}(\boldsymbol{\omega}) \leq \ell$ (Lemma 22 in Appendix I), $s(\bar{h}(T)) - s(\underline{h}(T)) \leq \frac{\sqrt{\bar{h}(T)} - \sqrt{\underline{h}(T)}}{\sqrt{|\mathcal{X}_0|}}$ and $\sum_{t=\underline{h}(T)}^{\bar{h}(T)} \frac{1}{\sqrt{t}} \leq \sqrt{\bar{h}(T)} - \sqrt{\underline{h}(T)}$. Substituting $\underline{h}(T)$ and $\bar{h}(T)$ from (23) and simplifying the terms, we get:

$$F(\boldsymbol{\omega}^*) - F(\hat{\boldsymbol{\omega}}(t)) \leq \phi_{\bar{h}(T)} \leq \ell T^{-b} + \frac{5\ell K^2}{\sqrt{|\mathcal{X}_0|}} T^{-\frac{a+b}{2}} + 3\epsilon(1 - T^{-b}) \leq 5\epsilon$$

when

$$T \geq \max \left\{ (4|\mathcal{X}_0|)^{\frac{1}{a}}, \left(\frac{4K^2}{\epsilon^2 D^2 |\mathcal{X}_0|} \right)^{\frac{1}{a}}, \left(\frac{2}{\Delta_{\min}} \right)^{\frac{2}{a}}, \left(\frac{3\|\boldsymbol{\mu}\|_{\infty}}{2} \right)^{\frac{2}{a}} \right\} + \max \left\{ \left(\frac{\ell}{\epsilon} \right)^{\frac{1}{b}}, \left(\frac{5\ell K^2}{\epsilon \sqrt{|\mathcal{X}_0|}} \right)^{\frac{2}{a+b}} \right\},$$

where the first inequality is due to $F(\hat{\boldsymbol{\omega}}(t)) \geq \bar{F}_{\eta_t}(\hat{\boldsymbol{\omega}}(t))$ by Proposition 2 (i) in §4.1. \square

Lemma 10. *Let $\epsilon \in (0, \min\{1, 2D^2\Delta_{\min}^2/K\})$ and T be a positive integer s.t.*

$$T \geq \max \left\{ (4|\mathcal{X}_0|)^{\frac{1}{a}}, \left(\frac{4K^2}{\epsilon^2 D^2 |\mathcal{X}_0|} \right)^{\frac{1}{a}}, \left(\frac{2}{\Delta_{\min}} \right)^{\frac{2}{a}}, \left(\frac{3\|\boldsymbol{\mu}\|_{\infty}}{2} \right)^{\frac{2}{a}} \right\}. \quad (44)$$

Suppose $\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T)$ holds. Then, for any $t \geq \underline{h}(T)$,

$$\Delta_{\min}(\hat{\boldsymbol{\mu}}(t)) \geq \frac{\Delta_{\min}}{2} \quad \text{and} \quad \|\hat{\boldsymbol{\mu}}(t)\|_{\infty} \leq \frac{3\|\boldsymbol{\mu}\|_{\infty}}{2}. \quad (45)$$

Proof Fix any T satisfying (45) and suppose $\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T)$ holds. Consider any $t \geq T^a$ and hence $t \geq 4|\mathcal{X}_0|$. To show the first inequality of (45), from $\mathcal{E}_{2,\epsilon}(T)$ and $\epsilon < 2D^2\Delta_{\min}^2/K$, we have

$$\Delta_{\min}(\hat{\boldsymbol{\mu}}(t-1)) \geq \Delta_{\min} - \frac{2D\epsilon}{24D^3\|\boldsymbol{\mu}\|_{\infty}} > \Delta_{\min} - \frac{\Delta_{\min}^2}{6K\|\boldsymbol{\mu}\|_{\infty}} > \frac{\Delta_{\min}}{2},$$

where the last inequality is because $\Delta_{\min} \leq 2D\|\boldsymbol{\mu}\|_{\infty}$ and $D \leq K$. To show the second inequality of (45), observe that

$$\|\hat{\boldsymbol{\mu}}(t)\|_{\infty} \leq \|\boldsymbol{\mu}\|_{\infty} + \frac{\epsilon}{24D^3\|\boldsymbol{\mu}\|_{\infty}} < \|\boldsymbol{\mu}\|_{\infty} + \frac{\Delta_{\min}^2}{12KD\|\boldsymbol{\mu}\|_{\infty}} < \frac{3\|\boldsymbol{\mu}\|_{\infty}}{2},$$

where the first inequality is because of $\mathcal{E}_{2,\epsilon}(T)$, the second is due to $\epsilon < 2D^2\Delta_{\min}^2/K$, and the last uses $\Delta_{\min} \leq 2D\|\boldsymbol{\mu}\|_{\infty}$. \square

Lemma 11. *Let $\epsilon > 0$ and $t \in \mathbb{N}$ be such that (45) holds. Then, under the event $\mathcal{E}_{1,\epsilon}^{(t)} \cap \mathcal{E}_{2,\epsilon}^{(t)}$,*

$$\begin{cases} t\phi_t \leq (t - |\mathcal{X}_0|)\phi_{t-|\mathcal{X}_0|} + 2\ell\sqrt{D}|\mathcal{X}_0|^2 & \text{if } t = p(t) - 1 \\ t\phi_t \leq (t-1)\phi_{t-1} + 3\epsilon + \ell \left(\eta_{t-1} + \frac{K^2}{2t\eta_t} \right) & \text{if } t \in [p(t), q(t)] \end{cases}, \quad (43)$$

where $p(t) = (s(t)^2 + 1)|\mathcal{X}_0|$, $q(t) = (s(t) + 1)^2|\mathcal{X}_0| - 1$, and $s(t) = \lfloor \sqrt{t/|\mathcal{X}_0|} \rfloor - 1$.

Proof The first case basically follows from the Lipschitzness of F (Appendix I), whereas the second relies on results on stochastic smoothing (Appendix H).

Case $t = p(t) - 1$: In this case, round t is exactly the end (including) round of a forced-exploration procedure. By ℓ -Lipschitzness of F (Lemma 21 in Appendix I),

$$F(\hat{\boldsymbol{\omega}}(t)) - F(\hat{\boldsymbol{\omega}}(t - |\mathcal{X}_0|)) \geq -\ell \|\hat{\boldsymbol{\omega}}(t) - \hat{\boldsymbol{\omega}}(t - |\mathcal{X}_0|)\|_2 \geq -\frac{\ell\sqrt{D}|\mathcal{X}_0|^2}{t},$$

where the second inequality stems from $\hat{\boldsymbol{\omega}}(t) = \frac{t-|\mathcal{X}_0|}{t}\hat{\boldsymbol{\omega}}(t-|\mathcal{X}_0|) + \frac{|\mathcal{X}_0|}{t}\sum_{\mathbf{x} \in \mathcal{X}_0} \mathbf{x}$ after performing the forced exploration. It then follows that $\|\hat{\boldsymbol{\omega}}(t) - \hat{\boldsymbol{\omega}}(t - |\mathcal{X}_0|)\|_2 \leq \frac{\sqrt{D}|\mathcal{X}_0|^2}{t}$. By $\max_{\boldsymbol{\omega} \in \Sigma} F(\boldsymbol{\omega}) \leq \ell$ (Lemma 22 in Appendix I) and a rearrangement of the above yields

$$t\phi_t \leq t\phi_{t-|\mathcal{X}_0|} + \ell\sqrt{D}|\mathcal{X}_0|^2 \leq (t - |\mathcal{X}_0|)\phi_{t-|\mathcal{X}_0|} + \ell|\mathcal{X}_0|(\sqrt{D}|\mathcal{X}_0| + 1).$$

The proof is completed after simplifying the terms.

Case: $t \in [p(t), q(t)]$: In this case, round t performs a FW update. For brevity, let $\mathbf{z} = \hat{\boldsymbol{\omega}}(t)$ and $\mathbf{y} = \hat{\boldsymbol{\omega}}(t-1)$. By $\frac{\ell K}{\eta_t}$ -smoothness of \bar{F}_{η_t} (Proposition 2 (iii) in §4.1) and $\mathbf{z} - \mathbf{y} = \frac{1}{t}(\mathbf{x}(t) - \mathbf{y})$,

$$\bar{F}_{\eta_t}(\mathbf{z}) \geq (*) - \frac{\ell K}{2\eta_t} \|\mathbf{z} - \mathbf{y}\|_2^2 \geq (*) - \frac{\ell K}{2t^2\eta_t} \|\mathbf{x}(t) - \mathbf{y}\|_2^2 \geq (*) - \frac{\ell K^2}{2t^2\eta_t},$$

where $(*) = \bar{F}_{\eta_t}(\mathbf{y}) + \langle \nabla \bar{F}_{\eta_t}(\mathbf{y}), \mathbf{z} - \mathbf{y} \rangle = \bar{F}_{\eta_t}(\mathbf{y}) + \frac{1}{t} \langle \nabla \bar{F}_{\eta_t}(\mathbf{y}), \mathbf{x}(t) - \mathbf{y} \rangle$. It follows from $\mathcal{E}_{1,\epsilon}^{(t)} \cap \mathcal{E}_{2,\epsilon}^{(t)}$ and the continuity argument (Lemma 7 in Appendix G.1) that

$$\begin{aligned} \langle \nabla \bar{F}_{\eta_t}(\mathbf{y}), \mathbf{x}(t) - \mathbf{y} \rangle &\geq \langle \nabla \bar{F}_{\hat{\mu}(t-1), \eta_t}(\mathbf{y}), \mathbf{x}(t) - \mathbf{y} \rangle - \epsilon \\ &\geq \max_{\mathbf{x} \in \mathcal{X}} \langle \nabla \bar{F}_{\hat{\mu}(t-1), \eta_t}(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle - 2\epsilon \geq \max_{\mathbf{x} \in \mathcal{X}} \langle \nabla \bar{F}_{\eta_t}(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle - 3\epsilon. \end{aligned}$$

Then, the duality gap [Jag13] and the ℓ -Lipschitzness of F (Lemma 21 in Appendix I) yield

$$\begin{aligned} \max_{\mathbf{x} \in \mathcal{X}} \langle \nabla \bar{F}_{\eta_t}(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle &\geq \max_{\boldsymbol{\omega} \in \Sigma} \bar{F}_{\eta_t}(\boldsymbol{\omega}) - \bar{F}_{\eta_t}(\mathbf{y}) \\ &\geq (F(\boldsymbol{\omega}^*) - \eta_t \ell) - (\bar{F}_{\eta_{t-1}}(\mathbf{y}) + \ell(\eta_{t-1} - \eta_t)) = \phi_{t-1} - \ell \eta_{t-1}. \end{aligned}$$

Therefore, $\bar{F}_{\eta_t}(\mathbf{z}) \geq \bar{F}_{\eta_t}(\mathbf{y}) + \frac{\phi_{t-1} - \ell \eta_{t-1} - 3\epsilon}{t} - \frac{\ell K^2}{2t^2 \eta_t}$ and subtracting $F(\boldsymbol{\omega}^*)$ on both sides,

$$\begin{aligned} \phi_t &= F(\boldsymbol{\omega}^*) - \bar{F}_{\eta_t}(\mathbf{z}) \\ &\leq (F(\boldsymbol{\omega}^*) - \bar{F}_{\eta_t}(\mathbf{y})) + \frac{-\phi_{t-1} + \ell \eta_{t-1} + 3\epsilon}{t} + \frac{\ell K^2}{2t^2 \eta_t} \\ &= \frac{t-1}{t} \phi_{t-1} + \frac{1}{t} \left(3\epsilon + \ell \left(\eta_{t-1} + \frac{K^2}{2t \eta_t} \right) \right), \end{aligned}$$

which completes the proof. \square

F Upper bound of $\sum_{T=M+1}^{\infty} \mathbb{P}_{\mu}[(\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T))^c]$ under P-FWS

Recall from (24) that $\mathcal{E}_{1,\epsilon}(T) = \cap_{t=\underline{h}(T)}^T \mathcal{E}_{1,\epsilon}^{(t)}$ and $\mathcal{E}_{2,\epsilon}(T) = \cap_{t=\underline{h}(T)}^T \mathcal{E}_{2,\epsilon}^{(t)}$, where

$$\begin{aligned} \mathcal{E}_{1,\epsilon}^{(t)} &= \left\{ \langle \nabla \bar{F}_{\hat{\mu}(t-1), \eta_t}(\hat{\omega}(t-1)), \mathbf{x}(t) \rangle \geq \max_{\mathbf{x} \in \mathcal{X}} \langle \nabla \bar{F}_{\hat{\mu}(t-1), \eta_t}(\hat{\omega}(t-1)), \mathbf{x} \rangle - \epsilon \right\}, \\ \mathcal{E}_{2,\epsilon}^{(t)} &= \left\{ \|\hat{\mu}(t-1) - \mu\|_{\infty} < \frac{\epsilon}{24D^3 \|\mu\|_{\infty}} \right\}, \end{aligned}$$

$T \geq M$ and M is defined in (25). Also, recall $\mathbf{x}(t) \in \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \langle \nabla \tilde{F}_{\hat{\mu}(t-1), \eta_t, n_t}(\hat{\omega}(t-1)), \mathbf{x} \rangle$, where $\nabla \tilde{F}_{\hat{\mu}(t-1), \eta_t, n_t}(\hat{\omega}(t-1))$ is computed by (ρ_t, θ_t) -MCP algorithm with

$$(\eta_t, n_t, \rho_t, \theta_t) = \left(\frac{1}{4\sqrt{t|\mathcal{X}_0|}}, \lceil t^{\frac{1}{4}} \rceil, \frac{1}{16tD^2\mathcal{X}_0}, \frac{1}{t^{\frac{1}{4}}e^{\sqrt{t}}} \right). \quad (10)$$

Our main result Lemma 8 is built by bounding $\mathbb{P}_{\mu}[\mathcal{E}_{1,\epsilon}(T)]$ and $\mathbb{P}_{\mu}[\mathcal{E}_{2,\epsilon}(T)]$ separately with Lemma 12 in F.1 and Lemma 14 in F.2.

Lemma 8. *Let $\epsilon \in (0, 2D^2\Delta_{\min}^2/K)$ and M be defined as in (25) Then,*

$$\sum_{T=M+1}^{\infty} \mathbb{P}_{\mu}[(\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T))^c] < 2K \left(\frac{3}{A_1(\epsilon)^{2+\frac{2}{a}}} + \frac{2}{A_2(\epsilon)^{2+\frac{2}{a}}} \right) \Gamma\left(2 + \frac{2}{a}\right),$$

where $A_1(\epsilon) = \min\{1, \frac{\epsilon^2}{8\ell^2 K^3 D^2}\}$, $A_2(\epsilon) = \frac{\epsilon^2}{2304D^6 \|\mu\|_{\infty}^2 \sqrt{|\mathcal{X}_0|}}$, and Γ denotes the gamma function $\Gamma(z) = \int_0^{\infty} t^{z-1} e^{-t} dt$ for any $z > 0$.

Proof Fix $\epsilon > 0$. For all $T \geq M$, we have

$$\mathbb{P}_{\mu}[(\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T))^c] \leq \mathbb{P}_{\mu}[\mathcal{E}_{1,\epsilon}(T)^c] + \mathbb{P}_{\mu}[\mathcal{E}_{2,\epsilon}(T)^c].$$

Bounding $\mathbb{P}_{\mu}[\mathcal{E}_{1,\epsilon}(T)^c]$: This is done by using Lemma 12 with $\mathbf{v} = \hat{\omega}(t)$ and $(\eta, n) = (\eta_t, n_t)$. Before applying Lemma 12, we verify that our chosen ρ_t in (10) satisfies the assumption of Lemma 12:

$$\frac{(\min_{k \in [K]} \hat{\omega}_k(t) - \eta_t)^2}{D^2} \geq \frac{1}{D^2} \left(\frac{1}{2\sqrt{t|\mathcal{X}_0|}} - \frac{1}{4\sqrt{t|\mathcal{X}_0|}} \right)^2 = \frac{1}{16tD^2|\mathcal{X}_0|} = \rho_t,$$

where the inequality is because of $\min_{k \in [K]} \hat{\omega}_k(t) \geq \frac{1}{2\sqrt{t|\mathcal{X}_0|}}$ (Lemma 14) and that $\eta_t = \frac{1}{4\sqrt{t|\mathcal{X}_0|}}$ in (10). Then, applying Lemma 12 with $\mathbf{v} = \hat{\omega}(t)$, $(v, \eta, n) = (\frac{1}{2\sqrt{t|\mathcal{X}_0|}}, \frac{1}{4\sqrt{t|\mathcal{X}_0|}}, \lceil t^{\frac{1}{4}} \rceil)$ yields that:

$$\mathbb{P}_{\mu}[\mathcal{E}_{1,\epsilon}(T)^c] \leq \sum_{t=\underline{h}(T)}^T K \left(2 \exp\left(-\frac{\epsilon^2 \sqrt{t}}{8\ell^2 K^3 D^2}\right) + \exp(-\sqrt{t}) \right) \leq \sum_{t=\underline{h}(T)}^T 3K \exp(-\sqrt{t}A_1(\epsilon)),$$

where $A_1(\epsilon) = \min\{1, \frac{\epsilon^2}{8\ell^2 K^3 D^2}\}$.

Bounding $\mathbb{P}_{\mu}[\mathcal{E}_{2,\epsilon}(T)^c]$: As Lemma 14 provides a lower bound on the number of pulls, $\min_{k \in [K]} N_k(t) \geq \frac{1}{2} \sqrt{\frac{t}{|\mathcal{X}_0|}}$, for all arms, using this lower bound of $N_k(t)$ as the number of i.i.d. samples in the application of Chernoff bound leads to:

$$\mathbb{P}_{\mu} \left[|\hat{\mu}_k(t) - \mu_k| \geq \frac{\epsilon}{24D^3 \|\mu\|_{\infty}} \right] \leq 2 \exp(-\sqrt{t}A_2(\epsilon)).$$

Hence, $\mathbb{P}_{\mu}[\mathcal{E}_{2,\epsilon}(T)^c] \leq 2K \sum_{t=\underline{h}(T)}^T \exp(-\sqrt{t}A_2(\epsilon))$. Then, we have

$$\begin{aligned} \sum_{T=M+1}^{\infty} \mathbb{P}_{\mu}[(\mathcal{E}_{1,\epsilon}(T) \cap \mathcal{E}_{2,\epsilon}(T))^c] &\leq \int_{M+1}^{\infty} \int_{T^a}^{\infty} (3K e^{-\sqrt{t}A_1(\epsilon)} + 2K e^{-\sqrt{t}A_2(\epsilon)}) dt dT \\ &\leq 2K \left(\frac{3}{A_1(\epsilon)^{2+\frac{2}{a}}} + \frac{2}{A_2(\epsilon)^{2+\frac{2}{a}}} \right) \Gamma\left(2 + \frac{2}{a}\right), \end{aligned}$$

where the second inequality uses Lemma 15. \square

F.1 Lemmas for bounding $\mathbb{P}_\mu[\mathcal{E}_{1,\epsilon}(T)^c]$

The following lemma is a result of two concentration inequalities, one bounds how much the empirical average deviates from the expectation (Proposition 3), and the other bounds the error incurred by MCP (Lemma 13).

Lemma 12. *Let $(\boldsymbol{\pi}, \boldsymbol{\omega}, \theta) \in \Lambda \times \Sigma_+ \times (0, 1)$, $v \in (0, \min_{k \in [K]} \omega_k)$, and $\eta \in (0, v)$. Then, $\forall \epsilon \in (0, 4K(v - \eta)/D)$,*

$$\mathbb{P} \left[\left\langle \nabla \bar{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}), \tilde{\boldsymbol{x}}^* - \boldsymbol{\omega} \right\rangle \geq \max_{\boldsymbol{x} \in \mathcal{X}} \left\langle \nabla \bar{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}), \boldsymbol{x} - \boldsymbol{\omega} \right\rangle - \epsilon \right] \geq 1 - K \left(2 \exp \left(-\frac{\epsilon^2 n^2}{8\ell^2 K^3 D^2} \right) + n\theta \right),$$

where $\nabla \tilde{F}_{\boldsymbol{\pi}, \eta, n}(\boldsymbol{\omega})$ is computed by $\left(\frac{(v-\eta)^2}{D^2}, \theta \right)$ -MCP, and $\tilde{\boldsymbol{x}}^* \in \operatorname{argmax}_{\boldsymbol{x} \in \mathcal{X}} \left\langle \nabla \tilde{F}_{\boldsymbol{\pi}, \eta, n}(\boldsymbol{\omega}), \boldsymbol{x} \right\rangle$.

Proof Let $\boldsymbol{x}^* \in \operatorname{argmax}_{\boldsymbol{x} \in \mathcal{X}} \left\langle \nabla \bar{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}), \boldsymbol{x} \right\rangle$. From $\tilde{\boldsymbol{x}}^* \in \operatorname{argmax}_{\boldsymbol{x} \in \mathcal{X}} \left\langle \nabla \tilde{F}_{\boldsymbol{\pi}, \eta, n}(\boldsymbol{\omega}), \boldsymbol{x} \right\rangle$,

$$\begin{aligned} \left\langle \nabla \bar{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}), \boldsymbol{x}^* - \tilde{\boldsymbol{x}}^* \right\rangle &\leq \left\langle \nabla \bar{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}), \boldsymbol{x}^* - \tilde{\boldsymbol{x}}^* \right\rangle + \left\langle \nabla \tilde{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}), \tilde{\boldsymbol{x}}^* - \boldsymbol{x}^* \right\rangle \\ &= \left\langle \nabla \bar{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}) - \nabla \tilde{F}_{\boldsymbol{\pi}, \eta, n}(\boldsymbol{\omega}), \boldsymbol{x}^* - \tilde{\boldsymbol{x}}^* \right\rangle. \end{aligned}$$

Fix $\epsilon > 0$. Recall that $\nabla \tilde{F}_{\boldsymbol{\pi}, \eta, n}(\boldsymbol{\omega}) = \frac{1}{n} \sum_{m=1}^n \nabla_{\boldsymbol{\omega}} f_{\hat{\boldsymbol{x}}_m}(\boldsymbol{\omega} + \eta \boldsymbol{Z}_m, \boldsymbol{\pi})$ where each $\hat{\boldsymbol{x}}_m$ is computed by $\left(\frac{(v-\eta)^2}{D^2}, \theta \right)$ -MCP $(\boldsymbol{\omega} + \eta \boldsymbol{Z}_m, \boldsymbol{\pi})$, and each \boldsymbol{Z}_m is independently sampled from $\text{Uniform}(B_2)$. Now, consider any fixed $\boldsymbol{x} = \boldsymbol{e}_k$ for any $k \in [K]$. Invoking Proposition 3 with $\epsilon = \frac{\epsilon}{4K}$ and $\boldsymbol{x} = \boldsymbol{e}_k$, we get:

$$\mathbb{P} \left[\left| \left\langle \nabla \bar{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}) - \frac{1}{n} \sum_{m=1}^n \nabla F_{\boldsymbol{\pi}}(\boldsymbol{\omega} + \eta \boldsymbol{Z}_m), \boldsymbol{e}_k \right\rangle \right| \geq \frac{\epsilon}{4K} \right] \leq 2 \exp \left(-\frac{\epsilon^2 n^2}{8\ell^2 K^3 D^2} \right).$$

Also, for $\nabla \tilde{F}_{\boldsymbol{\pi}, \eta, n}(\boldsymbol{\omega})$ computed by the $\left(\frac{(v-\eta)^2}{D^2}, \theta \right)$ -MCP algorithm, Lemma 13 with $\boldsymbol{x} = \boldsymbol{e}_k$, and $\theta = \theta$ and the assumption that $\epsilon \in (0, 4K(v - \eta)/D)$ implies that:

$$\mathbb{P} \left[\left| \left\langle \frac{1}{n} \sum_{m=1}^n \nabla F_{\boldsymbol{\pi}}(\boldsymbol{\omega} + \eta \boldsymbol{Z}_m) - \nabla \tilde{F}_{\boldsymbol{\pi}, \eta, n}(\boldsymbol{\omega}), \boldsymbol{e}_k \right\rangle \right| \geq \frac{\epsilon}{4K} \right] \leq n\theta.$$

Combining the two inequalities leads to:

$$\mathbb{P} \left[\left| \left\langle \nabla \bar{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}) - \nabla \tilde{F}_{\boldsymbol{\pi}, \eta, n}(\boldsymbol{\omega}), \boldsymbol{e}_k \right\rangle \right| \leq \frac{\epsilon}{2K} \right] \geq 1 - \left(2 \exp \left(-\frac{\epsilon^2 n^2}{8\ell^2 K^3 D^2} \right) + n\theta \right).$$

Then, an application of a union bound over all $\{\boldsymbol{e}_k\}_{k \in [K]}$ gives

$$\mathbb{P} \left[\left\langle \nabla \bar{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}) - \nabla \tilde{F}_{\boldsymbol{\pi}, \eta, n}(\boldsymbol{\omega}), \boldsymbol{x}^* - \tilde{\boldsymbol{x}}^* \right\rangle \leq \epsilon \right] \geq 1 - K \left(2 \exp \left(-\frac{\epsilon^2 n^2}{8\ell^2 K^3 D^2} \right) + n\theta \right). \quad (46)$$

Observe $\left\langle -\nabla \tilde{F}_{\boldsymbol{\pi}, \eta, n}(\boldsymbol{\omega}), \boldsymbol{x}^* - \tilde{\boldsymbol{x}}^* \right\rangle \geq 0$ implies

$$\left\{ \left\langle \nabla \bar{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}) - \nabla \tilde{F}_{\boldsymbol{\pi}, \eta, n}(\boldsymbol{\omega}), \boldsymbol{x}^* - \tilde{\boldsymbol{x}}^* \right\rangle \leq \epsilon \right\} \subseteq \left\{ \left\langle \nabla \bar{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}), \boldsymbol{x}^* - \tilde{\boldsymbol{x}}^* \right\rangle \leq \epsilon \right\}. \quad (47)$$

From (46)-(47), we conclude that the r.h.s. of (47) happens with probability at least $1 - K \left(2 \exp \left(-\frac{\epsilon^2 n^2}{8\ell^2 K^3 D^2} \right) + n\theta \right)$. The proof is completed by simply rearranging the r.h.s. of (47). \square

Lemma 13. *Let $(\boldsymbol{\pi}, \boldsymbol{\omega}, \boldsymbol{x}, \theta) \in \Lambda \times \Sigma_+ \times \{0, 1\}^K \times (0, 1)$ with $\|\boldsymbol{x}\|_1 \leq D$ and $v \in (0, \min_{k \in [K]} \omega_k)$.*

$$\forall (\eta, \boldsymbol{z}) \in (0, v) \times B_2, \quad \mathbb{P} \left[\left| \left\langle \nabla_{\boldsymbol{\omega}} f_{\boldsymbol{x}_*}(\boldsymbol{\omega} + \eta \boldsymbol{z}) - \nabla_{\boldsymbol{\omega}} f_{\hat{\boldsymbol{x}}}(\boldsymbol{\omega} + \eta \boldsymbol{z}), \boldsymbol{x} \right\rangle \right| \leq \frac{v - \eta}{D} \right] \geq 1 - \theta,$$

where \boldsymbol{x}_* is some action satisfying $f_{\boldsymbol{x}_*}(\boldsymbol{\omega} + \eta \boldsymbol{z}) = F_{\boldsymbol{\pi}}(\boldsymbol{\omega} + \eta \boldsymbol{z})$, and $\hat{\boldsymbol{x}}$ is the returned action of $\left(\frac{(v-\eta)^2}{D^2}, \theta \right)$ -MCP $(\boldsymbol{\omega} + \eta \boldsymbol{z}, \boldsymbol{\pi})$.

Proof This basically follows from a direct calculation. Let $\epsilon > 0$ and fix any $(\boldsymbol{\pi}, \boldsymbol{\omega}, \boldsymbol{x}) \in \Lambda \times \Sigma_+ \times \{0, 1\}^K$, $\|\boldsymbol{x}\|_1 \leq D$, and any $(\eta, \boldsymbol{z}) \in (0, v) \times B_2$. Then, for $\hat{\boldsymbol{x}}$ computed by (ρ, θ) -MCP($\boldsymbol{\omega} + \eta\boldsymbol{z}, \boldsymbol{\pi}$) with $\rho = (v - \eta)^2/D^2$, we have with probability at least $1 - \theta$

$$\begin{aligned} \rho &\geq |\langle \nabla_{\boldsymbol{\omega}} f_{\boldsymbol{x}_*}(\boldsymbol{\omega} + \eta\boldsymbol{z}) - \nabla_{\boldsymbol{\omega}} f_{\hat{\boldsymbol{x}}}(\boldsymbol{\omega} + \eta\boldsymbol{z}), \boldsymbol{\omega} + \eta\boldsymbol{z} \rangle| \\ &\geq \min_{k \in [K]} (\boldsymbol{\omega} + \eta\boldsymbol{z})_k \|\nabla_{\boldsymbol{\omega}} f_{\boldsymbol{x}_*}(\boldsymbol{\omega} + \eta\boldsymbol{z}) - \nabla_{\boldsymbol{\omega}} f_{\hat{\boldsymbol{x}}}(\boldsymbol{\omega} + \eta\boldsymbol{z})\|_{\infty}. \end{aligned}$$

Hence, remarking that $\min_{k \in [K]} (\boldsymbol{\omega} + \eta\boldsymbol{z})_k \geq v - \eta > 0$, we get: with probability at least $1 - \theta$,

$$|\langle \nabla_{\boldsymbol{\omega}} f_{\boldsymbol{x}_*}(\boldsymbol{\omega} + \eta\boldsymbol{z}) - \nabla_{\boldsymbol{\omega}} f_{\hat{\boldsymbol{x}}}(\boldsymbol{\omega} + \eta\boldsymbol{z}), \boldsymbol{x} \rangle| \leq \frac{\rho D}{v - \eta} = \frac{v - \eta}{D},$$

where we used the fact that $\|\boldsymbol{x}\|_1 \leq D$ and Hölder's inequality. \square

Proposition 3. Let $(\boldsymbol{\pi}, \boldsymbol{\omega}, \boldsymbol{x}) \in \Lambda \times \Sigma_+ \times \{0, 1\}^K$, $\eta \in (0, \min_{k \in [K]} \omega_k)$, and $\|\boldsymbol{x}\|_1 \leq D$. Then,

$$\forall \epsilon > 0, \quad \mathbb{P} \left[\left| \left\langle \nabla \bar{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}) - \frac{1}{n} \sum_{m=1}^n \nabla F_{\boldsymbol{\pi}}(\boldsymbol{\omega} + \eta \boldsymbol{z}_m), \boldsymbol{x} \right\rangle \right| \geq \epsilon \right] \leq 2 \exp \left(-\frac{2\epsilon^2 n^2}{\ell^2 K D^2} \right),$$

where $\boldsymbol{z}_1, \dots, \boldsymbol{z}_n$ are independently sampled from $\text{Uniform}(B_2)$.

Proof Fix $(\boldsymbol{\pi}, \boldsymbol{\omega}, \boldsymbol{x}) \in \Lambda \times \Sigma_+ \times \{0, 1\}^K$ where $\|\boldsymbol{x}\|_1 \leq D$, and fix $\eta \in (0, \min_{k \in [K]} \omega_k)$. Define

$$\phi(\boldsymbol{z}_1, \dots, \boldsymbol{z}_n) = \left\langle \nabla \bar{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}) - \frac{1}{n} \sum_{m=1}^n \nabla F_{\boldsymbol{\pi}}(\boldsymbol{\omega} + \eta \boldsymbol{z}_m), \boldsymbol{x} \right\rangle.$$

Note that $\mathbb{E}_{\boldsymbol{z}_1, \dots, \boldsymbol{z}_n} [\phi(\boldsymbol{z}_1, \dots, \boldsymbol{z}_n)] = 0$ by definition. Now we also observe that:

$$\max_{\boldsymbol{z}_1, \dots, \boldsymbol{z}_n, \boldsymbol{z}' \in B_2, m \in [n]} |\phi(\boldsymbol{z}_1, \dots, \boldsymbol{z}_n) - \phi(\boldsymbol{z}_1, \dots, \boldsymbol{z}_{m-1}, \boldsymbol{z}', \boldsymbol{z}_{m+1}, \dots, \boldsymbol{z}_n)| \leq \frac{\ell D}{n}$$

due to the ℓ -Lipschitzness of $F_{\hat{\boldsymbol{\mu}}}$ (Lemma 21 in Appendix I) and $\max_{\boldsymbol{x} \in \mathcal{X}} \|\boldsymbol{x}\|_1 \leq D$. Hence it follows from McDiarmid's inequality (Lemma 16 in F.3) that

$$\forall \epsilon > 0, \quad \mathbb{P}[|\phi(\boldsymbol{z}_1, \dots, \boldsymbol{z}_n)| \geq \epsilon] \leq 2 \exp \left(-\frac{2\epsilon^2}{K \left(\frac{\ell D}{n}\right)^2} \right) = 2 \exp \left(-\frac{2\epsilon^2 n^2}{\ell^2 K D^2} \right).$$

\square

F.2 Lemmas for bounding $\mathbb{P}_{\mu}[\mathcal{E}_{2, \epsilon}(T)^c]$

Lemma 14 (forced exploration). Let $\mathcal{X}_0 \subseteq \mathcal{X}$ be any set covering all arms $[K]$ and $t \geq 4|\mathcal{X}_0|$. Any algorithm with forced-exploration procedure satisfies

$$\hat{\boldsymbol{\omega}}(t) \in \Sigma_{\sqrt{\frac{1}{t|\mathcal{X}_0|}} - \frac{1}{t}} \subset \Sigma_{\frac{1}{2}\sqrt{\frac{1}{t|\mathcal{X}_0|}}}, \quad \forall t \geq 4|\mathcal{X}_0|.$$

Proof Fix any $k \in [K]$. By merely counting the rounds before t performing forced exploration,

$$N_k(t) \geq \sum_{s \in [t]: \lfloor \sqrt{\frac{s}{|\mathcal{X}_0|}} \rfloor \in \mathbb{N}} \sum_{\boldsymbol{x} \in \mathcal{X}_0} \boldsymbol{x}_k \geq \sqrt{\frac{t}{|\mathcal{X}_0|}} - 1 \geq \frac{1}{2} \sqrt{\frac{t}{|\mathcal{X}_0|}},$$

where the last inequality holds for any $t \geq 4|\mathcal{X}_0|$. \square

F.3 Technical lemmas

Lemma 15 ([WTP21]). Let $\alpha \in (0, 1)$ and $A, \beta > 0$. Then,

$$\int_0^{\infty} \left(\int_{T^{\alpha}}^{\infty} e^{-At^{\beta}} dt \right) dT = \frac{\Gamma(\frac{1}{\alpha\beta} + \frac{1}{\beta})}{\beta A^{\frac{1}{\alpha\beta} + \frac{1}{\beta}}}.$$

Proof The result of Lemma 5 [WTP21] is stated for $\alpha, \beta \in (0, 1)$ but it actually applies for the case of $\beta > 0$ as well. Here we provide a proof for completeness.

$$\int_0^\infty \left(\int_{T^\alpha}^\infty e^{-At^\beta} dt \right) dT = \int_0^\infty \alpha T^\alpha e^{-AT^{\alpha\beta}} dT = \frac{1}{\beta} \int_0^\infty x^{\frac{1}{\alpha\beta} + \frac{1}{\beta} - 1} e^{-Ax} dx = \frac{\Gamma(\frac{1}{\alpha\beta} + \frac{1}{\beta})}{\beta A^{\frac{1}{\alpha\beta} + \frac{1}{\beta}}}.$$

□

The below Lemma 16, also known as bounded different inequality, can be found in many textbooks, e.g., Theorem 6.2 in [BLM13].

Lemma 16 (McDiarmid's inequality). *Let $\mathcal{Z} = (\mathcal{Z}_1, \dots, \mathcal{Z}_n)$ be independent random variables, and $\phi : \mathbb{R}^n \mapsto \mathbb{R}$ be a measurable function. Suppose $\phi(\mathbf{z})$ changes by at most $c_i > 0$ under an arbitrary change of the i -th coordinate. Then,*

$$\forall \epsilon > 0, \quad \mathbb{P}[\phi(\mathcal{Z}) - \mathbb{E}[\phi(\mathcal{Z})] \geq \epsilon] \leq \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^n c_i^2}\right).$$

G Continuity arguments

In this section, we establish the continuity of $F_\pi(\boldsymbol{\omega})$ and $\nabla \bar{F}_{\pi,\eta}(\boldsymbol{\omega})$ in π for any fixed $\boldsymbol{\omega} \in \Sigma_+$, where $\nabla \bar{F}_{\pi,\eta}(\boldsymbol{\omega})$ denotes the gradient $\nabla_{\boldsymbol{\omega}} \bar{F}_{\pi,\eta}(\boldsymbol{\omega})$ taken w.r.t. the input $\boldsymbol{\omega}$. As the consequence of the continuity of F_π and $\nabla \bar{F}_{\pi,\eta}$ in π , we can show the point-wise convergence of $F_{\hat{\boldsymbol{\mu}}(t)} \rightarrow F_\mu$ and $\nabla \bar{F}_{\hat{\boldsymbol{\mu}}(t),\eta} \rightarrow \nabla \bar{F}_{\mu,\eta}$ given that $\hat{\boldsymbol{\mu}}(t) \rightarrow \mu$ almost surely.

Notation. Throughout this section, we define $\nabla \bar{F}_{\pi,\eta}(\boldsymbol{\omega}) = \mathbf{0}_K$ if $\eta \geq \min_{k \in [K]} \omega_k$ for any $(\pi, \boldsymbol{\omega}) \in \Lambda \times \Sigma_+$. Moreover, for any $(\mathbf{v}, \boldsymbol{\omega}) \in \mathbb{R}^K \times \Sigma_+$, we will use $\nabla_\pi F_{\mathbf{v}}(\boldsymbol{\omega})$ (resp. $\nabla_\pi \left(\frac{\partial \bar{F}_{\mathbf{v},\eta}(\boldsymbol{\omega})}{\partial \omega_k} \right)$) to denote the gradient of the function $\pi \mapsto F_\pi(\boldsymbol{\omega})$ (resp. $\pi \mapsto \frac{\partial \bar{F}_{\pi,\eta}(\boldsymbol{\omega})}{\partial \omega_k}$) evaluated at the point \mathbf{v} .

The main result in this section, Lemma 7, is derived based on Lemma 17 in Appendix G.1 (which asserts the continuity of the function $\psi_{\boldsymbol{\omega},\mathbf{x},\eta}(\boldsymbol{\pi}) = \langle \nabla \bar{F}_{\pi,\eta}(\boldsymbol{\omega}), \mathbf{x} - \boldsymbol{\omega} \rangle$ on \mathbb{R}^K) and Proposition 4 in Appendix G.2 (which upper bounds the length of $\nabla f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\mu})$).

Lemma 7. *Let $\boldsymbol{\mu} \in \Lambda$ and $\epsilon \in (0, \frac{2D^2 \Delta_{\min}(\boldsymbol{\mu})^2}{K})$. Then, any $\boldsymbol{\pi} \in \mathbb{R}^K$ with $\|\boldsymbol{\pi} - \boldsymbol{\mu}\|_\infty < \frac{\epsilon}{24D^3 \|\boldsymbol{\mu}\|_\infty}$ satisfies the following:*

$$|F_\mu(\boldsymbol{\omega}) - F_\pi(\boldsymbol{\omega})| < \epsilon, \quad \forall \boldsymbol{\omega} \in \Sigma_+ \quad (35)$$

$$|\langle \nabla \bar{F}_{\pi,\eta}(\boldsymbol{\omega}) - \nabla \bar{F}_{\mu,\eta}(\boldsymbol{\omega}), \mathbf{x} - \boldsymbol{\omega} \rangle| < \epsilon, \quad \forall (\boldsymbol{\omega}, \mathbf{x}) \in \Sigma_+ \times \mathcal{X}, \forall \eta \in (0, \min_{k \in [K]} \omega_k). \quad (36)$$

Proof Inspired by Lemma 14 in [WTP21], we prove this lemma using Proposition 4 and applying the mean-value theorem to $\psi_{\boldsymbol{\omega},\mathbf{x},\eta}$.

Fix $(\boldsymbol{\omega}, \boldsymbol{\mu}) \in \Sigma_+ \times \Lambda$, and let $\mathbf{i}^* = \mathbf{i}^*(\boldsymbol{\mu})$ and $\Delta_{\mathbf{x}} = \Delta_{\mathbf{x}}(\boldsymbol{\mu})$ for any $\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{i}^*\}$. Fix $\epsilon \in (0, \frac{2D^2 \Delta_{\min}^2}{K})$ and $\boldsymbol{\pi} \in \mathbb{R}^K$ such that $\|\boldsymbol{\pi} - \boldsymbol{\mu}\|_\infty < \frac{\epsilon}{24D^3 \|\boldsymbol{\mu}\|_\infty}$. One may check that this $\boldsymbol{\pi}$ satisfies the assumption of Proposition 4 as

$$\|\boldsymbol{\pi} - \boldsymbol{\mu}\|_\infty < \frac{\epsilon}{24D^3 \|\boldsymbol{\mu}\|_\infty} < \frac{2D^2 \Delta_{\min}^2}{24KD^3 \|\boldsymbol{\mu}\|_\infty} = \frac{\Delta_{\min}^2}{12KD \|\boldsymbol{\mu}\|_\infty} \leq \frac{\Delta_{\min}}{6K} < \frac{\Delta_{\min}}{\sqrt{2KD}},$$

where the second inequality stems from the choice of ϵ and the second last is because $\Delta_{\min} \leq 2D \|\boldsymbol{\mu}\|_\infty$. In what follows, we will be applying the mean-value theorem to $\psi_{\boldsymbol{\omega},\mathbf{x},\eta}$ (whose continuity is stated in Lemma 17). For convenience, introduce the function $\mathbf{r}(\beta) = (1 - \beta)\boldsymbol{\mu} + \beta\boldsymbol{\pi}$ for any $\beta \in (0, 1)$.

Proof of (35): For any $\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{i}^*\}$, by the mean-value theorem, there exists a $\beta \in (0, 1)$ such that

$$\begin{aligned} |f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\pi}) - f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\mu})| &= |\langle \nabla_\pi f_{\mathbf{x}}(\boldsymbol{\omega}, \mathbf{r}(\beta)), \boldsymbol{\pi} - \boldsymbol{\mu} \rangle| \\ &= \left| \sum_{k \in [K]} \omega_k \left\langle \nabla_\pi \left(\frac{\partial f_{\mathbf{x}}(\boldsymbol{\omega}, \mathbf{r}(\beta))}{\partial \omega_k} \right), \boldsymbol{\pi} - \boldsymbol{\mu} \right\rangle \right| \\ &\leq \sum_{k \in [K]} \omega_k \left\| \nabla_\pi \left(\frac{\partial f_{\mathbf{x}}(\boldsymbol{\omega}, \mathbf{r}(\beta))}{\partial \omega_k} \right) \right\|_1 \|\boldsymbol{\pi} - \boldsymbol{\mu}\|_\infty < \epsilon, \end{aligned} \quad (48)$$

where the last inequality uses $\boldsymbol{\omega} \in \Sigma_+$, $\|\boldsymbol{\pi} - \boldsymbol{\mu}\|_\infty < \frac{\epsilon}{24D^3 \|\boldsymbol{\mu}\|_\infty}$, $\left\| \nabla_\pi \left(\frac{\partial f_{\mathbf{x}}(\boldsymbol{\omega}, \mathbf{r}(\beta))}{\partial \omega_k} \right) \right\|_1 \leq 12D^2 \|\boldsymbol{\mu}\|_\infty$ (Proposition 4). Hence, from a substitution of \mathbf{x} in (48) with $\mathbf{x}_e \in \operatorname{argmin}_{\mathbf{x} \neq \mathbf{i}^*} f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\mu})$ and the fact that $F_\pi(\boldsymbol{\omega}) \leq f_{\mathbf{x}_e}(\boldsymbol{\omega}, \boldsymbol{\pi})$, we derive

$$F_\pi(\boldsymbol{\omega}) - F_\mu(\boldsymbol{\omega}) \leq f_{\mathbf{x}_e}(\boldsymbol{\omega}, \boldsymbol{\pi}) - f_{\mathbf{x}_e}(\boldsymbol{\omega}, \boldsymbol{\mu}) < \epsilon.$$

The other inequality of $F_\mu(\boldsymbol{\omega}) - F_\pi(\boldsymbol{\omega}) < \epsilon$ can be derived similarly. This proves (35).

Proof of (36): Recall that $\psi_{\boldsymbol{\omega},\mathbf{x},\eta}(\boldsymbol{\pi}) = \langle \nabla \bar{F}_{\pi,\eta}(\boldsymbol{\omega}), \mathbf{x} - \boldsymbol{\omega} \rangle$ is continuous on \mathbb{R}^K (Lemma 17). By the mean-value theorem, there exists $\beta \in (0, 1)$ such that

$$\begin{aligned} |\psi_{\boldsymbol{\omega},\mathbf{x},\eta}(\boldsymbol{\pi}) - \psi_{\boldsymbol{\omega},\mathbf{x},\eta}(\boldsymbol{\mu})| &= |\langle \nabla_\pi \psi_{\boldsymbol{\omega},\mathbf{x},\eta}(\mathbf{r}(\beta)), \boldsymbol{\pi} - \boldsymbol{\mu} \rangle| \\ &\leq \|\nabla_\pi \psi_{\boldsymbol{\omega},\mathbf{x},\eta}(\mathbf{r}(\beta))\|_1 \|\boldsymbol{\pi} - \boldsymbol{\mu}\|_\infty. \end{aligned} \quad (49)$$

To bound $\|\nabla_{\pi}\psi_{\omega,\mathbf{x},\eta}(\mathbf{r}(\beta))\|_1$, we write

$$\nabla_{\pi}\psi_{\omega,\mathbf{x},\eta}(\mathbf{r}(\beta)) = \sum_{k \in [K]} \nabla_{\pi} \left(\frac{\partial \bar{F}_{\mathbf{r}(\beta),\eta}(\boldsymbol{\omega})}{\partial \omega_k} \right) (x_k - \omega_k).$$

Then it follows from the fundamental theorem of calculus that: the gradient ∇_{π} and the expectation operators are exchangeable, i.e.,

$$\forall k \in [K], \quad \nabla_{\pi} \left(\frac{\partial \bar{F}_{\mathbf{r}(\beta),\eta}(\boldsymbol{\omega})}{\partial \omega_k} \right) = \mathbb{E}_{\mathcal{Z} \sim \text{Uniform}(B_2)} \left[\nabla_{\pi} \left(\frac{\partial F_{\mathbf{r}(\beta)}(\boldsymbol{\omega} + \eta \mathcal{Z})}{\partial \omega_k} \right) \right].$$

As shown in Appendix H, $\frac{\partial F_{\mathbf{r}(\beta)}(\boldsymbol{\omega} + \eta \mathcal{Z})}{\partial \omega_k}$ exists almost surely. When such gradient exists, Proposition 4 bounds its 1-norm length by

$$\left\| \nabla_{\pi} \left(\frac{\partial F_{\mathbf{r}(\beta)}(\boldsymbol{\omega} + \eta \mathcal{Z})}{\partial \omega_k} \right) \right\|_1 \leq 12D^2 \|\boldsymbol{\mu}\|_{\infty},$$

so it follows that $\left\| \nabla_{\pi} \left(\frac{\partial \bar{F}_{\mathbf{r}(\beta),\eta}(\boldsymbol{\omega})}{\partial \omega_k} \right) \right\|_1 \leq 12D^2 \|\boldsymbol{\mu}\|_{\infty}$ as well. Hence, substituting the above back to $\nabla_{\pi}\psi_{\omega,\mathbf{x},\eta}(\mathbf{r}(\beta))$ yields:

$$\|\nabla_{\pi}\psi_{\omega,\mathbf{x},\eta}(\mathbf{r}(\beta))\|_1 \leq \max_{k \in [K]} \left\| \nabla_{\pi} \left(\frac{\partial \bar{F}_{\mathbf{r}(\beta),\eta}(\boldsymbol{\omega})}{\partial \omega_k} \right) \right\|_1 \|\mathbf{x} - \boldsymbol{\omega}\|_1 \leq 24D^3 \|\boldsymbol{\mu}\|_{\infty},$$

where the first inequality use Hölder's inequality. Finally, plugging the above into (49) and recalling that $\|\boldsymbol{\pi} - \boldsymbol{\mu}\|_{\infty} < \frac{\epsilon}{24D^3 \|\boldsymbol{\mu}\|_{\infty}}$, we have

$$|\psi_{\omega,\mathbf{x},\eta}(\boldsymbol{\pi}) - \psi_{\omega,\mathbf{x},\eta}(\boldsymbol{\mu})| < \epsilon.$$

This concludes the proof. \square

G.1 An application of the maximum theorem

Recall that $\psi_{\omega,\mathbf{x},\eta}(\boldsymbol{\pi}) = \langle \nabla \bar{F}_{\boldsymbol{\pi},\eta}(\boldsymbol{\omega}), \mathbf{x} - \boldsymbol{\omega} \rangle$.

Lemma 17. *For any $\epsilon > 0$, there exists a constant $\xi_{\epsilon} > 0$ such that if $\|\boldsymbol{\pi} - \boldsymbol{\mu}\|_{\infty} < \xi_{\epsilon}$, then*

$$|\psi_{\omega,\mathbf{x},\eta}(\boldsymbol{\pi}) - \psi_{\omega,\mathbf{x},\eta}(\boldsymbol{\mu})| < \epsilon, \quad \forall (\boldsymbol{\omega}, \mathbf{x}) \in \Sigma_+ \times \mathcal{X}, \quad \forall \eta \in (0, \min_{k \in [K]} \omega_k). \quad (50)$$

The proof of Lemma 17 replies on the celebrated maximum theorem [FKV14], which is introduced below. After that, we then show its proof.

Maximum Theorem: Here we briefly introduce the maximum theorem and Lemma 17 will be proved at the end of this section. The definitions and results are taken from [FKV14] (see also Appendix K.1 of [WTP21]).

Definition 1. *Let $U \neq \emptyset$ be a subset of a topological space and $h : U \mapsto \mathbb{R}$ be a function. Define the level sets of h for $y \in \mathbb{R}$ as*

$$L_h(y, U) = \{x \in U : h(x) \leq y\} \quad \text{and} \quad L_h^<(y, U) = \{x \in U : h(x) < y\}.$$

The function h is said to be lower semi-continuous (resp. upper semi-continuous) on U if $L_h(y, U)$ are closed (resp. $L_h^<(y, U)$ are open) for all $y \in \mathbb{R}$; h is said to be inf-compact on U if $L_h(y, U)$ and $L_h^<(y, U)$ are compact for all $y \in \mathbb{R}$.

Definition 2. *Let \mathbb{X} and \mathbb{Y} be Hausdorff topological spaces and $\Phi : \mathbb{X} \rightrightarrows \mathbb{S}(\mathbb{Y})$ be a set-valued function, where $\mathbb{S}(\mathbb{Y})$ is the set of non-empty subsets of \mathbb{Y} . Define*

$$Gr_U(\Phi) = \{(x, y) \in U \times \mathbb{Y} : y \in \Phi(x)\}$$

as the graph of Φ restricted to U . The function $u : \mathbb{X} \times \mathbb{Y} \mapsto \mathbb{R}$ is said to be \mathbb{K} -inf-compact on $Gr_{\mathbb{X}}(\Phi)$ if for all non-empty compact subset C of \mathbb{X} , u is inf-compact on $Gr_C(\Phi)$.

Theorem 6 (Maximum theorem). *Suppose \mathbb{X} is compactly generated, $\Phi : \mathbb{X} \rightrightarrows \mathbb{S}(\mathbb{Y})$ is lower hemicontinuous, and $u : \mathbb{X} \times \mathbb{Y} \mapsto \mathbb{R}$ is \mathbb{K} -inf-compact and upper semi-continuous on $Gr_{\mathbb{X}}(\Phi)$. Then, the function $v(x) = \inf_{y \in \Phi(x)} u(x, y)$ is continuous and the set of its optimal solutions $\Phi^*(x) = \{y \in \Phi(x) : u(x, y) = v(x)\}$ is upper hemicontinuous and compact-valued.*

Proof of Lemma 17: Fix any $\boldsymbol{\mu} \in \Lambda$ and let $\boldsymbol{i}^* = \boldsymbol{i}^*(\boldsymbol{\mu})$. The goal is to show that for any $\epsilon > 0$, there exists a constant $\xi_\epsilon > 0$ such that if $\|\boldsymbol{\pi} - \boldsymbol{\mu}\|_\infty < \xi_\epsilon$, then

$$|\psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}(\boldsymbol{\pi}) - \psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}(\boldsymbol{\mu})| < \epsilon, \quad \forall (\boldsymbol{\omega}, \boldsymbol{x}) \in \Sigma_+ \times \mathcal{X}, \forall \eta \in (0, \min_{k \in [K]} \omega_k), \quad (50)$$

where $\psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}(\boldsymbol{\pi}) = \langle \nabla \bar{F}_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega}), \boldsymbol{x} - \boldsymbol{\omega} \rangle$. In what follows, we will use p to denote the probability distribution of $\text{Uniform}(B_2)$. We will first show that $\psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}$ is continuous for each fixed $(\boldsymbol{\omega}, \boldsymbol{x}, \eta) \in \Sigma_+ \times \mathcal{X} \times (0, 1)$, and then use Theorem 6 to show (50).

Continuity of $\psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}$: Fix $(\boldsymbol{\omega}, \boldsymbol{x}, \eta) \in \Sigma_+ \times \mathcal{X} \times (0, 1)$. Let $U_\eta = \{z \in B_2 : |\partial F_\mu(\boldsymbol{\omega} + \eta z)| > 1\}$ which is a measure-zero set under p (Lemma 20 in Appendix H). For its complement set $B_2 \setminus U_\eta$, we split $B_2 \setminus U_\eta = \cup_{\boldsymbol{y} \neq \boldsymbol{i}^*} B_\eta(\boldsymbol{y})$ into possibly overlapping sets $B_\eta(\boldsymbol{y}) = \{z \in B_2 \setminus U_\eta : \nabla F_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega} + \eta z) = \nabla_\omega f_{\boldsymbol{y}}(\boldsymbol{\omega} + \eta z, \boldsymbol{\pi})\}$, and define $\psi_{\boldsymbol{\omega}, \boldsymbol{y}, \eta}(\boldsymbol{y}, \cdot) = \int_{z \in B_\eta(\boldsymbol{y})} \langle \nabla_\omega f_{\boldsymbol{y}}(\boldsymbol{\omega} + \eta z, \cdot), \boldsymbol{x} - \boldsymbol{\omega} \rangle dp(z)$ on each of these sets $B_\eta(\boldsymbol{y})$. Observe that for any $\boldsymbol{\pi} \in \mathbb{R}^K$, we have

$$\psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}(\boldsymbol{\pi}) = \int_{z \in B_2 \setminus U_\eta} \langle \nabla F_{\boldsymbol{\pi}, \eta}(\boldsymbol{\omega} + \eta z), \boldsymbol{x} - \boldsymbol{\omega} \rangle dp(z) = \sum_{\boldsymbol{y} \neq \boldsymbol{i}^*} \psi_{\boldsymbol{\omega}, \boldsymbol{y}, \eta}(\boldsymbol{y}, \boldsymbol{\pi}).$$

To show the continuity of $\psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}(\boldsymbol{\pi})$, it suffices to show that each $\psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}(\boldsymbol{y}, \cdot)$ is continuous. Fix $\boldsymbol{y} \in \mathcal{X} \setminus \{\boldsymbol{i}^*\}$ and any sequence $\{\boldsymbol{\pi}_n\}_{n=1}^\infty$ converging to $\boldsymbol{\mu}$. Then, for any $\forall z \in B_2$, we have

- (i) $|\langle \nabla_\omega f_{\boldsymbol{y}}(\boldsymbol{\omega} + \eta z, \boldsymbol{\pi}_n), \boldsymbol{x} - \boldsymbol{\omega} \rangle| \leq \|\nabla_\omega f_{\boldsymbol{y}}(\boldsymbol{\omega} + \eta z, \boldsymbol{\pi}_n)\|_\infty \|\boldsymbol{x} - \boldsymbol{\omega}\|_1 \leq 2D\ell$
- (ii) $\lim_{n \rightarrow \infty} \langle \nabla_\omega f_{\boldsymbol{y}}(\boldsymbol{\omega} + \eta z, \boldsymbol{\pi}_n), \boldsymbol{x} - \boldsymbol{\omega} \rangle = \langle \nabla_\omega f_{\boldsymbol{y}}(\boldsymbol{\omega} + \eta z, \boldsymbol{\mu}), \boldsymbol{x} - \boldsymbol{\omega} \rangle$. This is because $\nabla_\omega f_{\boldsymbol{y}}(\boldsymbol{\omega} + \eta z, \cdot) = \frac{\langle \boldsymbol{i}^* - \boldsymbol{y}, \cdot \rangle^2 (\boldsymbol{i}^* \oplus \boldsymbol{y}) \odot (\boldsymbol{\omega} + \eta z)^{-2}}{2(\boldsymbol{i}^* \oplus \boldsymbol{y}, (\boldsymbol{\omega} + \eta z)^{-1})^2}$ by Lemma 19 and Proposition 1 (Appendix C.1) is obviously continuous and that function composition preserves continuity.

From (i) and (ii), the dominated convergence theorem implies that

$$\psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}(\boldsymbol{y}, \boldsymbol{\mu}) = \lim_{n \rightarrow \infty} \int_{z \in B_\eta(\boldsymbol{y})} \langle \nabla_\omega f_{\boldsymbol{y}}(\boldsymbol{\omega} + \eta z, \boldsymbol{\pi}_n), \boldsymbol{x} - \boldsymbol{\omega} \rangle dp(z).$$

This shows the continuity of $\psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}(\boldsymbol{y}, \cdot)$ for each $\boldsymbol{y} \neq \boldsymbol{i}^*$, and thus $\psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}$ is continuous.

Application of the maximum theorem (Theorem 6): For this part, we take the approach similar to Lemma 6 in [WTP21]. Define

$$\phi(\boldsymbol{\pi}) = \min \{ -|\psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}(\boldsymbol{\pi}) - \psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}(\boldsymbol{\mu})| : (\boldsymbol{\omega}, \boldsymbol{x}, \eta) \in \Sigma_+ \times \mathcal{X} \times (0, 1) \}.$$

We prove the continuity of ϕ on $\mathcal{S} = \mathbb{R}^K \setminus \text{cl}(\text{Alt}(\boldsymbol{\mu}))$ by invoking Theorem 6 with the following substitutions:

- $\mathbb{X} = \mathcal{S}$,
- $\Phi = \Sigma_+ \times \mathcal{X} \times (0, 1)$,
- $\mathbb{Y} = \Sigma_+ \times \mathcal{X} \times (0, 1)$,
- $u(\boldsymbol{\pi}, \boldsymbol{\omega}, \boldsymbol{x}, \eta) = -|\psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}(\boldsymbol{\pi}) - \psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}(\boldsymbol{\mu})|$.

Here we verify that the assumptions of Theorem 6 are satisfied. \mathbb{X} is compactly generated as \mathcal{S} is a metric space; Φ is continuous as it is a constant map; u is continuous due to the continuity of $\psi_{\boldsymbol{\omega}, \boldsymbol{x}, \eta}$. To show that u is \mathbb{K} -inf compact, consider any compact set $C \subset \mathcal{S}$ and any $y \in \mathbb{R}$. We see that $L_u(y, C \times \Sigma_+ \times \mathcal{X} \times (0, 1))$ is compact because it is bounded (as $\Sigma_+ \times \mathcal{X} \times (0, 1)$ is bounded and C is compact) and closed (as u is continuous and the preimage of $[0, y]$ is closed). Hence, ϕ is continuous on \mathcal{S} by Theorem 6. Finally, by $\phi(\boldsymbol{\mu}) = 0$ and the continuity of ϕ , there exists $\xi_\epsilon > 0$ such that $\phi(\boldsymbol{\pi}) > -\epsilon$ for any $\|\boldsymbol{\pi} - \boldsymbol{\mu}\|_\infty < \xi_\epsilon$. This completes the proof of (50). \square

G.2 The length of gradients

Throughout this subsection, we fix $\boldsymbol{\mu} \in \Lambda$ and denote $\boldsymbol{i}^* = \boldsymbol{i}^*(\boldsymbol{\mu})$, $\Delta_{\boldsymbol{x}} = \Delta_{\boldsymbol{x}}(\boldsymbol{\mu})$, and $\Delta_{\min}(\boldsymbol{\mu}) = \Delta_{\min}$ for short. Here we aim to present Proposition 4, in which (i) quantifies how close an estimate $\boldsymbol{\pi}$ of $\boldsymbol{\mu}$ should be such that $\boldsymbol{i}^*(\boldsymbol{\pi}) = \boldsymbol{i}^*$, and (ii) asserts the continuity of any component of $\nabla_\omega f_{\boldsymbol{x}}(\boldsymbol{\omega}, \boldsymbol{\pi})$ in $\boldsymbol{\pi}$, and that its gradient with respect to $\boldsymbol{\pi}$ is bounded.

Proposition 4. Any $\boldsymbol{\pi} \in \mathbb{R}^K$ such that $\|\boldsymbol{\pi} - \boldsymbol{\mu}\|_\infty < \frac{\Delta_{\min}}{\sqrt{2KD}}$ satisfies

$$(i) \mathbf{i}^*(\boldsymbol{\pi}) = \mathbf{i}^*,$$

(ii) $\forall \mathbf{x} \in \mathcal{X} \setminus \{\mathbf{i}^*\}$ and all $k \in [K]$, $\frac{\partial f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\pi})}{\partial \omega_k}$ is continuous in $\boldsymbol{\pi}$ and

$$\left\| \nabla_{\boldsymbol{\pi}} \left(\frac{\partial f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\pi})}{\partial \omega_k} \right) \right\|_1 \leq 12D^2 \|\boldsymbol{\mu}\|_{\infty}.$$

Proof Proof of (i): Lemma 18 is equivalent to that: any $\boldsymbol{\pi} \in \mathbb{R}^K$ satisfying $\|\boldsymbol{\mu} - \boldsymbol{\pi}\|_{\infty} < \frac{\Delta_{\min}}{\sqrt{2KD}}$ implies that $\boldsymbol{\pi} \notin \text{cl}(\text{Alt}(\boldsymbol{\mu}))$. As closure of finite union equals union of closures,

$$\begin{aligned} \mathbb{R}^K \setminus \text{cl}(\text{Alt}(\boldsymbol{\mu})) &= \mathbb{R}^K \setminus \left(\bigcup_{\mathbf{x} \neq \mathbf{i}^*} \text{cl}(\{\boldsymbol{\lambda} \in \mathbb{R}^K : \langle \mathbf{i}^* - \mathbf{x}, \boldsymbol{\lambda} \rangle < 0\}) \right) \\ &= \mathbb{R}^K \setminus \left(\bigcup_{\mathbf{x} \neq \mathbf{i}^*} \{\boldsymbol{\lambda} \in \mathbb{R}^K : \langle \mathbf{i}^* - \mathbf{x}, \boldsymbol{\lambda} \rangle \leq 0\} \right) \\ &= \{\boldsymbol{\lambda} \in \mathbb{R}^K : \mathbf{i}^*(\boldsymbol{\lambda}) = \mathbf{i}^*\}. \end{aligned}$$

Thus, $\boldsymbol{\pi} \notin \text{cl}(\text{Alt}(\boldsymbol{\mu}))$ is equivalent to $\mathbf{i}^*(\boldsymbol{\pi}) = \mathbf{i}^*$. This concludes the proof of (i).

Proof of (ii): Fix any $\boldsymbol{\pi} \in \mathbb{R}^K$ satisfying $\|\boldsymbol{\mu} - \boldsymbol{\pi}\|_{\infty} < \frac{\Delta_{\min}}{\sqrt{2KD}}$. By Lemma 19 and $\mathbf{i}^*(\boldsymbol{\pi}) = \mathbf{i}^*$,

$$\forall k \in [K], \quad \frac{\partial f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\pi})}{\partial \omega_k} = \frac{\langle \mathbf{i}^* - \mathbf{x}, \boldsymbol{\pi} \rangle^2 (x_k \oplus i_k^*)}{2 \langle \mathbf{x} \oplus \mathbf{i}^*, \boldsymbol{\omega}^{-1} \rangle^2 \omega_k^2}.$$

Fix $k \in [K]$. Note that the function $\boldsymbol{\pi} \mapsto \frac{\partial f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\pi})}{\partial \omega_k}$ is continuous and differentiable since it consists of inner products, element-wise products, and since its denominator is always positive. For its derivative,

$$\begin{aligned} \left\| \nabla_{\boldsymbol{\pi}} \left(\frac{\partial f_{\mathbf{x}}(\boldsymbol{\omega}, \boldsymbol{\pi})}{\partial \omega_k} \right) \right\|_1 &= \left\| \frac{(\mathbf{i}^* - \mathbf{x}) \langle \mathbf{i}^* - \mathbf{x}, \boldsymbol{\pi} \rangle (x_k \oplus i_k^*)}{\langle \mathbf{x} \oplus \mathbf{i}^*, \boldsymbol{\omega}^{-1} \rangle^2 \omega_k^2} \right\|_1 \\ &\leq \|(\mathbf{i}^* - \mathbf{x}) \langle \mathbf{i}^* - \mathbf{x}, \boldsymbol{\pi} \rangle (x_k \oplus i_k^*)\|_1 \\ &\leq \|\mathbf{i}^* - \mathbf{x}\|_1 |\langle \mathbf{i}^* - \mathbf{x}, \boldsymbol{\pi} \rangle| \leq 4D^2 \|\boldsymbol{\pi}\|_{\infty} \leq 12D^2 \|\boldsymbol{\mu}\|_{\infty}, \end{aligned}$$

where the first inequality is because $\langle \mathbf{x} \oplus \mathbf{i}^*, \boldsymbol{\omega}^{-1} \rangle \omega_k \geq 1$ if $(x_k \oplus i_k^*) = 1$; the second is because $x_k \oplus i_k^* \leq 1$; the third uses $\|\mathbf{i}^* - \mathbf{x}\|_1 \leq 2D$ and $|\langle \mathbf{i}^* - \mathbf{x}, \boldsymbol{\pi} \rangle| \leq \|\mathbf{i}^* - \mathbf{x}\|_1 \|\boldsymbol{\pi}\|_{\infty}$; the last uses the triangle inequality:

$$\|\boldsymbol{\pi}\|_{\infty} \leq \|\boldsymbol{\mu}\|_{\infty} + \|\boldsymbol{\mu} - \boldsymbol{\pi}\|_{\infty} \leq \|\boldsymbol{\mu}\|_{\infty} + \frac{\Delta_{\min}}{\sqrt{2KD}} \leq 3\|\boldsymbol{\mu}\|_{\infty},$$

where the last inequality is due to an application of Hölder's inequality to

$$\Delta_{\min} \leq \min_{\mathbf{x} \neq \mathbf{i}^*} \|\mathbf{i}^* - \mathbf{x}\|_1 \|\boldsymbol{\mu}\|_{\infty} \leq 2D \|\boldsymbol{\mu}\|_{\infty}.$$

□

Lemma 18. $\inf_{\boldsymbol{\lambda} \in \text{Alt}(\boldsymbol{\mu})} \|\boldsymbol{\mu} - \boldsymbol{\lambda}\|_{\infty} \geq \frac{\Delta_{\min}}{\sqrt{2KD}}$.

Proof We claim that

$$\inf_{\boldsymbol{\lambda} \in \Lambda : \langle \boldsymbol{\lambda}, \mathbf{i}^* - \mathbf{x} \rangle < 0} \|\boldsymbol{\mu} - \boldsymbol{\lambda}\|_2^2 = \frac{\Delta_{\mathbf{x}}}{\|\mathbf{i}^* \oplus \mathbf{x}\|_2}, \quad \forall \mathbf{x} \neq \mathbf{i}^*. \quad (51)$$

Observe that the proof immediately follows from (51) because the facts that $\|\mathbf{y}\|_2 \leq \sqrt{K} \|\mathbf{y}\|_{\infty}$ for any $\mathbf{y} \in \mathbb{R}^K$, $\text{Alt}(\boldsymbol{\mu}) = \bigcup_{\mathbf{x} \neq \mathbf{i}^*} \{\boldsymbol{\lambda} \in \Lambda : \langle \boldsymbol{\lambda}, \mathbf{i}^* - \mathbf{x} \rangle < 0\}$, and $\|\mathbf{i}^* \oplus \mathbf{x}\|_2 \leq \sqrt{2D}$.

Proof of (51): By solving the stationary conditions, i.e., $\nabla_{\boldsymbol{\lambda}} \mathcal{L}_{\mathbf{x}}(\boldsymbol{\lambda}_{\mathbf{x}}^*, \alpha^*) = 2(\boldsymbol{\mu} - \boldsymbol{\lambda}_{\mathbf{x}}^*) + \alpha^*(\mathbf{i}^* - \mathbf{x}) = 0$ and $\nabla_{\alpha} \mathcal{L}_{\mathbf{x}}(\boldsymbol{\lambda}_{\mathbf{x}}^*, \alpha^*) = \langle \boldsymbol{\lambda}_{\mathbf{x}}^*, \mathbf{i}^* - \mathbf{x} \rangle = 0$, we find

$$\boldsymbol{\lambda}_{\mathbf{x}}^* = \boldsymbol{\mu} - \frac{\Delta_{\mathbf{x}}(\boldsymbol{\mu}) \odot (\mathbf{i}^* - \mathbf{x})}{\|\mathbf{i}^* \oplus \mathbf{x}\|_2^2}$$

is a minimizer for $\inf_{\boldsymbol{\lambda} \in \Lambda : \langle \boldsymbol{\lambda}, \mathbf{i}^* - \mathbf{x} \rangle < 0} \|\boldsymbol{\mu} - \boldsymbol{\lambda}\|_2$. (51) follows by plugging $\boldsymbol{\lambda}_{\mathbf{x}}^*$ into $\|\boldsymbol{\mu} - \boldsymbol{\lambda}\|_2^2$. □

Remind that $\nabla_{\boldsymbol{\omega}} f_{\mathbf{x}}$ can be evaluated by the following Lemma 19.

Lemma 19 (Envelope theorem). *Let $(\omega, \mu) \in \Sigma_+ \times \Lambda$ and $\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{i}^*\}$. Define $\lambda_{\omega, \mu}^*(\mathbf{x}) \in \operatorname{argmin}_{\lambda \in \operatorname{cl}(\mathcal{C}_{\mathbf{x}})} \left\langle \omega, \frac{(\mu - \lambda)^2}{2} \right\rangle$. Then,*

$$\nabla_{\omega} f_{\mathbf{x}}(\omega, \mu) = \frac{(\mu - \lambda_{\omega, \mu}^*(\mathbf{x}))^2}{2} = \frac{\Delta_{\mathbf{x}}(\mu)^2 (\mathbf{x} \oplus \mathbf{i}^*) \odot \omega^{-2}}{2 \langle \mathbf{x} \oplus \mathbf{i}^*, \omega^{-1} \rangle^2}.$$

Proof The first equality is an application of Lemma 6 and Proposition 1 of [WTP21] with $\mathcal{I} = \mathcal{X}$, $\mathcal{J}_{\mathbf{x}} = \{\mathbf{x}\}$, $\Sigma = \Sigma_K$, $\mathcal{S}_{\mathbf{x}} = \{\lambda \in \Lambda : \mathbf{i}^*(\lambda) = \mathbf{x}\}$ (see Appendix K.2 and Appendix K.4 in [WTP21] for more details). The second equality substitutes $\lambda_{\omega, \mu}^*(\mathbf{x}) = \mu + \frac{\Delta_{\mathbf{x}}(\mu)(\mathbf{x} - \mathbf{i}^*) \odot \omega^{-1}}{\langle \mathbf{x} \oplus \mathbf{i}^*, \omega^{-1} \rangle}$ by using (11)-(14). \square

H Stochastic smoothing

This section is devoted to present Proposition 2 and verify the assumptions required for applying Proposition 2 to our objective F_μ .

Stochastic smoothing [FKM05, DBW12] is a well-studied technique and has been widely applied to online convex nonsmooth optimization [HK12, H⁺16]. Proposition 2 is a restatement of existing results. In particular, Proposition 2 (i), (ii) and (iii) directly follow from Lemma E.2 in [DBW12] with $(L_0, u) = (\ell, \eta)$, $f = -\Phi$ and $f_u = -\bar{\Phi}_\eta(\cdot)$, and Proposition 2 (iv) can be established by Jensen's inequality as done in the proof of Theorem 2.1 [DBW12].

Proposition 2. *Assume that $\Phi : \mathbb{R}_{>0}^K \mapsto \mathbb{R}$ is concave, ℓ -Lipschitz, and differentiable almost everywhere. Let $B_2 = \{\mathbf{v} \in \mathbb{R}^K : \|\mathbf{v}\|_2 \leq 1\}$. For any $\boldsymbol{\omega} \in \Sigma_+$ and $\eta \in (0, \min_{k \in [K]} \omega_k)$, define*

$$\bar{\Phi}_\eta(\boldsymbol{\omega}) = \mathbb{E}_{\mathbf{Z} \sim \text{Uniform}(B_2)}[\Phi(\boldsymbol{\omega} + \eta\mathbf{Z})]. \quad (6)$$

Then, $\bar{\Phi}_\eta(\boldsymbol{\omega})$ satisfies that:

- (i) $\Phi(\boldsymbol{\omega}) - \eta\ell \leq \bar{\Phi}_\eta(\boldsymbol{\omega}) \leq \Phi(\boldsymbol{\omega})$
- (ii) $\nabla \bar{\Phi}_{\mu, \eta}(\boldsymbol{\omega}) = \mathbb{E}_{\mathbf{Z} \sim \text{Uniform}(B_2)}[\nabla \Phi_\mu(\boldsymbol{\omega} + \eta\mathbf{Z})]$
- (iii) $\bar{\Phi}_\eta$ is $\frac{\ell K}{\eta}$ -smooth
- (iv) if $\eta > \eta' > 0$, then $\bar{\Phi}_{\eta'}(\boldsymbol{\omega}) \geq \bar{\Phi}_\eta(\boldsymbol{\omega})$

Now, we validate assumptions of Proposition 2 on F_μ . The concavity of F_μ , which is shown by [WTP21], follows from the facts that each $f_x(\cdot, \boldsymbol{\mu})$ is concave and that F_μ is a minimum of these functions $f_x(\cdot, \boldsymbol{\mu})$ over all possible \mathbf{x} . The Lipschitzness of F_μ is shown in Lemma 21 in Appendix I). Hence, it remains to show the almost-everywhere differentiability of F_μ . To show that the set of non-differentiable points of F_μ , i.e.,

$$\bigcup_{\mathbf{x}, \mathbf{x}' \in \mathcal{X} \setminus \{\mathbf{i}^*(\boldsymbol{\mu})\}, \mathbf{x} \neq \mathbf{x}'}$$

is measure-zero under $\text{Uniform}(B_2)$, it suffices to show the following lemma.

Lemma 20. *Let $\boldsymbol{\mu} \in \Lambda$ and $\mathbf{x}_1, \mathbf{x}_2$ be distinct actions in $\mathcal{X} \setminus \{\mathbf{i}^*(\boldsymbol{\mu})\}$. Then under the probability measure of $\text{Uniform}(B_2)$,*

$$\{\mathbf{z} \in B_2 : f_{\mathbf{x}_1}(\boldsymbol{\omega} + \eta\mathbf{z}, \boldsymbol{\mu}) = f_{\mathbf{x}_2}(\boldsymbol{\omega} + \eta\mathbf{z}, \boldsymbol{\mu})\}$$

is a measure-zero set.

Proof To simplify the notation, let $\mathbf{i}^* = \mathbf{i}^*(\boldsymbol{\mu})$ and $\Delta_{\mathbf{x}} = \Delta_{\mathbf{x}}(\boldsymbol{\mu})$. Thanks to the close-form expressions of $f_{\mathbf{x}_1}$ and $f_{\mathbf{x}_2}$, $\mathbf{z} \in B_2$ such that $f_{\mathbf{x}_1}(\boldsymbol{\omega} + \eta\mathbf{z}, \boldsymbol{\mu}) = f_{\mathbf{x}_2}(\boldsymbol{\omega} + \eta\mathbf{z}, \boldsymbol{\mu})$ are the points satisfying that:

$$\frac{\Delta_{\mathbf{x}_1}^2}{2 \langle \mathbf{x}_1 \oplus \mathbf{i}^*, (\boldsymbol{\omega} + \eta\mathbf{z})^{-1} \rangle} = \frac{\Delta_{\mathbf{x}_2}^2}{2 \langle \mathbf{x}_2 \oplus \mathbf{i}^*, (\boldsymbol{\omega} + \eta\mathbf{z})^{-1} \rangle}.$$

In other words, the set of interests is

$$\left\{ \mathbf{z} \in B_2 : \sum_{k=1}^K a_k \prod_{k' \neq k} (\omega_{k'} + \eta z_{k'}) = 0 \right\}, \quad (52)$$

where $a_k = (\mathbf{x}_2 \oplus \mathbf{i}^*)_k \Delta_{\mathbf{x}_1}^2 - (\mathbf{x}_1 \oplus \mathbf{i}^*)_k \Delta_{\mathbf{x}_2}^2$ for all $k \in [K]$. We claim that \mathbf{a} is a non-zero vector. Otherwise, $a_k = 0, \forall k \in [K]$, which together with the fact that $\Delta_{\mathbf{x}_1}^2, \Delta_{\mathbf{x}_2}^2 > 0$ directly imply $(\mathbf{x}_2 \oplus \mathbf{i}^*)_k = 0$ if and only if $(\mathbf{x}_1 \oplus \mathbf{i}^*)_k = 0$. That means $\mathbf{x}_1 = \mathbf{x}_2$, but this becomes a contradiction. Therefore, the set in (52) are the roots of a non-zero polynomial inside B_2 , and hence it is a measure-zero set (see e.g. Lemma in [Oka73]). \square

I Lipschitzness of F_μ and boundness of F_μ on $\Sigma_K \cap \mathbb{R}_{>0}^K$

In this section, we show the Lipschitzness of $F_\mu(\mathbf{v}) = \min_{\mathbf{x} \neq \mathbf{i}^*} f_{\mathbf{x}}(\mathbf{v}, \boldsymbol{\mu})$ for $\mathbf{v} \in \mathbb{R}_{>0}^K$. Let \mathbf{x}_e be an equilibrium action such that $F_\mu(\mathbf{v}) = f_{\mathbf{x}_e}(\mathbf{v}, \boldsymbol{\mu})$. We will use the envelope theorem (Lemma 19 in Appendix G) to evaluate $\nabla_\omega f_{\mathbf{x}_e}(\mathbf{v}, \boldsymbol{\mu})$ in closed-form, and then bound its length. We will also derive an upper bound of $F_\mu(\mathbf{v})$ valid for any positive vector \mathbf{v} in the $(K-1)$ -dimensional simplex Σ_K . In what below, we denote $\mathbf{i}^* = \mathbf{i}^*(\boldsymbol{\mu})$ and $\Delta_{\mathbf{x}} = \Delta_{\mathbf{x}}(\boldsymbol{\mu})$ for any $\mathbf{x} \neq \mathbf{i}^*$ for short.

Lemma 21. *Let $\boldsymbol{\mu} \in \Lambda$ and $\ell = 2D^2 \|\boldsymbol{\mu}\|_\infty^2$. Then, F_μ is ℓ -Lipschitz with respect to $\|\cdot\|_\infty$ on $\mathbb{R}_{>0}^K$,*

Proof Let $\mathbf{v} \in \mathbb{R}_{>0}^K$. Recall that $F_\mu(\mathbf{v}) = \min_{\mathbf{x} \neq \mathbf{i}^*} f_{\mathbf{x}}(\mathbf{v}, \boldsymbol{\mu})$, and each $f_{\mathbf{x}}(\mathbf{v}, \boldsymbol{\mu})$ is differentiable (proven in Lemma 19 in Appendix G.2). Hence if \mathbf{x} is the action such that $F_\mu(\mathbf{v}) = f_{\mathbf{x}}(\mathbf{v}, \boldsymbol{\mu})$, the concavity of $F_\mu(\mathbf{v})$ and the fact that $\nabla_\omega f_{\mathbf{x}}(\mathbf{v}, \boldsymbol{\mu})$ is the subdifferential of F_μ on \mathbf{v} yield that

$$\forall \mathbf{v}' \in \mathbb{R}_{>0}^K, |F_\mu(\mathbf{v}) - F_\mu(\mathbf{v}')| \leq |\langle \nabla_\omega f_{\mathbf{x}}(\mathbf{v}, \boldsymbol{\mu}), \mathbf{v} - \mathbf{v}' \rangle| \leq \|\nabla_\omega f_{\mathbf{x}}(\mathbf{v}, \boldsymbol{\mu})\|_1 \|\mathbf{v} - \mathbf{v}'\|_\infty,$$

where the last inequality stems from Hölder's inequality. From the above, the ℓ -Lipschitz can be derived by upper bounding $\|\nabla_\omega f_{\mathbf{x}}(\mathbf{v}, \boldsymbol{\mu})\|_1$ by ℓ . Now applying Lemma 19 in Appendix G.2 yields

$$\|\nabla_\omega f_{\mathbf{x}}(\mathbf{v}, \boldsymbol{\mu})\|_1 = \left\| \frac{(\boldsymbol{\mu} - \boldsymbol{\lambda}_{\mathbf{v}, \boldsymbol{\mu}}^*(\mathbf{x}, \alpha_{\mathbf{x}}^*))^2}{2} \right\|_1 = \frac{\|\mathbf{v}^{-2} \odot (\mathbf{x} \oplus \mathbf{i}^*)\|_1 \Delta_{\mathbf{x}}^2}{2 \langle \mathbf{x} \oplus \mathbf{i}^*, \mathbf{v}^{-1} \rangle^2}. \quad (53)$$

To simplify the above, we observe that

$$\langle \mathbf{x} \oplus \mathbf{i}^*, \mathbf{v}^{-1} \rangle^2 = \left(\sum_{k=1}^K v_k^{-1} \mathbb{1}\{x_k \neq \mathbf{i}_k^*\} \right)^2 \geq \sum_{k=1}^K v_k^{-2} \mathbb{1}\{x_k \neq \mathbf{i}_k^*\} = \|\mathbf{v}^{-2} \odot (\mathbf{x} \oplus \mathbf{i}^*)\|_1, \quad (54)$$

where the inequality uses the fact that $v_k > 0$ for all $k \in [K]$. Also,

$$\Delta_{\mathbf{x}} = \langle \mathbf{i}^* - \mathbf{x}, \boldsymbol{\mu} \rangle \leq \|\mathbf{i}^* - \mathbf{x}\|_1 \|\boldsymbol{\mu}\|_\infty \leq 2D \|\boldsymbol{\mu}\|_\infty. \quad (55)$$

Thus, (53)-(54)-(55) yields that $\|\nabla_\omega f_{\mathbf{x}}(\mathbf{v}, \boldsymbol{\mu})\|_1 \leq 2D^2 \|\boldsymbol{\mu}\|_\infty^2$. \square

Lemma 22. *Let $\boldsymbol{\mu} \in \Lambda$ and $\ell = 2D^2 \|\boldsymbol{\mu}\|_\infty^2$. Then, $\max_{\boldsymbol{\omega} \in \Sigma_K \cap \mathbb{R}_{>0}^K} F_\mu(\boldsymbol{\omega}) \leq \ell$.*

Proof Observe that $f_{\mathbf{x}}(\mathbf{v}, \boldsymbol{\mu}) = \langle \boldsymbol{\omega}, \nabla_\omega f_{\mathbf{x}}(\mathbf{v}, \boldsymbol{\mu}) \rangle$ for any $\mathbf{x} \neq \mathbf{i}^*$. Combining this observation with the fact that $\Delta_{\mathbf{x}} \leq 2D \|\boldsymbol{\mu}\|_\infty$ (as argued in (54) in proof of Lemma 21) implies:

$$F_\mu(\mathbf{v}) = \min_{\mathbf{x} \neq \mathbf{i}^*} \frac{\Delta_{\mathbf{x}}^2}{2 \langle \mathbf{x} \oplus \mathbf{i}^*, \mathbf{v}^{-1} \rangle} \leq \frac{(2D \|\boldsymbol{\mu}\|_\infty)^2}{2} = \ell,$$

where the first inequality is because $\langle \mathbf{x} \oplus \mathbf{i}^*, \mathbf{v}^{-1} \rangle \geq \min_{k \in [K]} v_k^{-1} \geq 1$ (as $\mathbf{v} \in \Sigma_K$ and $v_k > 0$ for all $k \in [K]$). The proof is completed since \mathbf{v} is taken arbitrarily. \square

J Proofs related to combinatorial sets

Assumption 1. (i) There exists a polynomial-time algorithm identifying $\mathbf{i}^*(\mathbf{v})$ for any $\mathbf{v} \in \mathbb{R}^K$; (ii) \mathcal{X} is inclusion-wise maximal, i.e., there is no $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ s.t. $\mathbf{x} < \mathbf{x}'$; (iii) for each $k \in [K]$, there exists $\mathbf{x} \in \mathcal{X}$ such that $x_k = 1$; (iv) $|\mathcal{X}| \geq 2$.

As claimed in §2.2, Assumption 1 holds for the following combinatorial sets:

- *m-sets*: $\mathcal{X} = \{\mathbf{x} \in \{0, 1\}^K : \|\mathbf{x}\|_1 = m\}$
- *spanning forests*: \mathcal{X} is a set of all spanning forests in a given graph
- *bipartite matchings*: \mathcal{X} is a set of all maximal matchings in a given bipartite graph
- *s-t paths*: \mathcal{X} is the set of all source-destination paths in a directed acyclic graph

In what below, we present a simple proof for the above examples.

Proof Suppose (iii) (iv) hold (as we can always achieve (iii) by removing arms not covered by \mathcal{X} and (iv) holds for non-trivial sets). For (i), it is well-known that a polynomial-time LM Oracle, i.e., $\mathbf{i}^*(\cdot)$, exists for each of the discussed combinatorial structures. For example, see Chapter 39 in [S⁺03] for the greedy algorithm for matroids (applicable to m -set and spanning forests), Chapter 41 in [S⁺03] for the augmentation-based algorithm for 2-matroid intersection (applicable to bipartite matchings), and algorithms such as Dijkstra's algorithm for s - t paths.

It remains to verify (ii) the inclusion-wise maximal property of \mathcal{X} . For \mathcal{X} as m -sets, the inclusion-wise maximal property clearly holds because any binary vector $\mathbf{x}' > \mathbf{x}$ (resp. $\mathbf{x}' < \mathbf{x}$) for some $\mathbf{x} \in \mathcal{X}$ must have $\sum_{k \in [K]} x'_k > m$ (resp. $< m$) and thus $\mathbf{x}' \notin \mathcal{X}$. The case is similar for \mathcal{X} as spanning forests since the number of edges of any spanning forests of a graph is the same. For \mathcal{X} as maximal matchings in which the term 'maximal' exactly refers to being inclusion-wise maximal, (ii) directly follows from the definition. For \mathcal{X} as the set of all source-destination paths in an acyclic graph, if there exists any source-destination path $\mathbf{x}' > \mathbf{x}$ for some $\mathbf{x} \in \mathcal{X}$ then \mathbf{x}' must contain a cycle, so inclusion-wise maximal property holds. \square

Lemma 2. Let $\mathbf{v} \in \mathbb{R}^K$ and $\mathbf{x} \in \mathcal{X}$. Under Assumption 1, there exists an algorithm that solves $\max_{\mathbf{x}' \in \mathcal{X}: \mathbf{x}' \neq \mathbf{x}} \langle \mathbf{v}, \mathbf{x}' \rangle$ by only making at most D queries to the LM Oracle.

Proof Fix $\mathbf{x} \in \mathcal{X}$. Assume $\mathbf{v} \neq \mathbf{0}_K$ (as otherwise, any $\mathbf{x}' \neq \mathbf{x}$ is a second-best action). Inspired by Lawler-Murty's m -best algorithm [Law72], we will prove this lemma by considering the algorithm described as follows. It first computes $\mathbf{i}^*(\mathbf{v})$ by the LM Oracle, and returns it as the output if $\mathbf{i}^*(\mathbf{v}) \neq \mathbf{x}$. Otherwise, we identify the second-best action by the program below:

$$\max_{k \in [K]: x_k = 1} \langle \mathbf{v}, \mathbf{i}^*(\mathbf{v}^{(k)}) \rangle, \quad \text{where} \quad v_i^{(k)} = \begin{cases} -3 \|\mathbf{v}\|_1 & \text{if } i = k \\ v_i & \text{otherwise.} \end{cases} \quad (56)$$

Intuitively, for each arm k of \mathbf{x} , the action $\mathbf{i}^*(\mathbf{v}^{(k)})$ represents the best one among all actions without k (we have a strong negative weight on the k -th component of $\mathbf{v}^{(k)}$). In the following, we will show that at least one of $\{\mathbf{i}^*(\mathbf{v}^{(k)}) : k \in [K], x_k = 1\}$ is the second-best action.

More precisely, we will show that for any maximizer $a \in [K]$ to (56), $\mathbf{i}^*(\mathbf{v}^{(a)})$ is a second-best action. Consider if $(\mathbf{i}^*(\mathbf{v}^{(a)}))_a = 0$, then the claim follows from the fact that $\mathbf{i}^*(\mathbf{v}^{(a)})$ is the best among all actions without a and also the best in $\{\mathbf{i}^*(\mathbf{v}^{(k)}) : k \in [K], x_k = 1\}$. It suffices to show that $(\mathbf{i}^*(\mathbf{v}^{(a)}))_a = 1$ cannot happen. If $(\mathbf{i}^*(\mathbf{v}^{(a)}))_a = 1$, then it follows from Assumption 1 (iv) $|\mathcal{X}| \geq 2$ and (ii) the inclusion-wise maximality of \mathcal{X} that there is another action \mathbf{x}' such that $x'_k = 0$ but $x_k = 1$ for some $k \in [K]$. So, by $\mathbf{i}^*(\mathbf{v}^{(a)})_a = 1$, $\mathbf{v} \neq \mathbf{0}_K$ and the definition of $\mathbf{v}^{(a)}$, we get

$$\langle \mathbf{v}, \mathbf{i}^*(\mathbf{v}^{(a)}) \rangle = \sum_{j \in [K]: \mathbf{i}^*(\mathbf{v}^{(a)})_j = 1, j \neq a} v_j - 3 \|\mathbf{v}\|_1 \leq -2 \|\mathbf{v}\|_1 < \langle \mathbf{v}, \mathbf{x}' \rangle \leq \langle \mathbf{v}, \mathbf{i}^*(\mathbf{v}^{(k)}) \rangle,$$

which contradicts the optimality of a (as it would imply that $\mathbf{i}^*(\mathbf{v}^{(k)})$ is better).

Finally, as $\|\mathbf{x}\|_1 \leq D$, the number of LM Oracle calls required for solving (56) is at most D . \square

Finally, we present the property of \mathcal{X}_0 briefly argued in § 4.2.

Lemma 23. *Let \mathbf{e}_k is the k -th column of an identity matrix. Under Assumption 1, \mathcal{X}_0 is a $[K]$ -covering set and $|\mathcal{X}_0| \geq 2$.*

Proof Showing that \mathcal{X}_0 covers $[K]$: Assumption 1 (iii) ensures $\{\mathbf{x} \in \mathcal{X} : x_k = 1\} \neq \emptyset$, and it follows that $\max_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x}, \mathbf{e}_k \rangle = 1$, i.e., $(\mathbf{i}^*(\mathbf{e}_k))_k = 1$. As $(\mathbf{i}^*(\mathbf{e}_k))_k = 1$ holds for all k , the proof is completed.

Showing that $|\mathcal{X}_0| \geq 2$: Suppose on the contrary, $|\mathcal{X}_0| = 1$. Thanks to Assumption 1 (iv) $|\mathcal{X}| \geq 2$, there exists $\mathbf{x} \in \mathcal{X}$ such that $x_k = \langle \mathbf{e}_k, \mathbf{x} \rangle \geq \langle \mathbf{e}_k, \mathbf{x}' \rangle = x'_k$ for all $k \in [K]$, $\mathbf{x}' \neq \mathbf{x}$. Together with Assumption 1 (iii), one can easily deduce that $x_k = 1$ for all $k \in [K]$. However, this implies $\mathbf{x}' < \mathbf{x}$ for any $\mathbf{x}' \neq \mathbf{x}$ and hence contradicts to Assumption 1 (ii) that \mathcal{X} is inclusion-wise maximal. \square

K Sample complexity lower bound

In this section, we assume $\boldsymbol{\mu} \in \Lambda$ and $\delta \in (0, 1)$ is fixed, and show Theorem 7 by adapting Lemma 19 in [KCG16].

Lemma 24 ([KCG16]). *Any δ -PAC algorithm satisfies*

$$\forall \boldsymbol{\lambda} \in \text{Alt}(\boldsymbol{\mu}), \sum_{k \in [K]} \sum_{\mathbf{x} \in \mathcal{X}: x_k=1} \mathbb{E}_{\boldsymbol{\mu}}[N_{\mathbf{x}}(\tau)] \frac{(\mu_k - \lambda_k)^2}{2} \geq \text{kl}(\delta, 1 - \delta). \quad (57)$$

Theorem 7. *Any δ -PAC strategy satisfies*

$$\mathbb{E}_{\boldsymbol{\mu}}[\tau] \geq T^*(\boldsymbol{\mu}) \text{kl}(\delta, 1 - \delta) \quad \text{with} \quad T^*(\boldsymbol{\mu})^{-1} = \sup_{\boldsymbol{\omega} \in \Sigma} \inf_{\boldsymbol{\lambda} \in \text{Alt}(\boldsymbol{\mu})} \left\langle \boldsymbol{\omega}, \frac{(\boldsymbol{\mu} - \boldsymbol{\lambda})^2}{2} \right\rangle, \quad (1)$$

where $\Sigma = \{\sum_{\mathbf{x} \in \mathcal{X}} w_{\mathbf{x}} : \mathbf{w} \in \Sigma_{|\mathcal{X}|}\}$ and $\text{Alt}(\boldsymbol{\mu}) = \{\boldsymbol{\lambda} \in \Lambda : \mathbf{i}^*(\boldsymbol{\lambda}) \neq \mathbf{i}^*(\boldsymbol{\mu})\}$.

Proof We have: under any algorithm,

$$\sup_{\boldsymbol{\omega} \in \Sigma} \inf_{\boldsymbol{\lambda} \in \text{Alt}(\boldsymbol{\mu})} \sum_{k \in [K]} \omega_k \frac{(\mu_k - \lambda_k)^2}{2} \geq \inf_{\boldsymbol{\lambda} \in \text{Alt}(\boldsymbol{\mu})} \sum_{k \in [K]} \sum_{\mathbf{x} \in \mathcal{X}: x_k=1} \frac{\mathbb{E}_{\boldsymbol{\mu}}[N_{\mathbf{x}}(\tau)]}{\mathbb{E}_{\boldsymbol{\mu}}[\tau]} \frac{(\mu_k - \lambda_k)^2}{2},$$

Hence if the algorithm is δ -PAC, by Lemma 24,

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\mu}}[\tau] \sup_{\boldsymbol{\omega} \in \Sigma} \inf_{\boldsymbol{\lambda} \in \text{Alt}(\boldsymbol{\mu})} \sum_{k \in [K]} \omega_k \frac{(\mu_k - \lambda_k)^2}{2} &\geq \inf_{\boldsymbol{\lambda} \in \text{Alt}(\boldsymbol{\mu})} \sum_{k \in [K]} \sum_{\mathbf{x} \in \mathcal{X}: x_k=1} \mathbb{E}_{\boldsymbol{\mu}}[N_{\mathbf{x}}(\tau)] \frac{(\mu_k - \lambda_k)^2}{2} \\ &\geq \text{kl}(\delta, 1 - \delta). \end{aligned}$$

□

Lemma 1. *For any $\boldsymbol{\mu} \in \Lambda$, $T^*(\boldsymbol{\mu}) \leq 4KD \Delta_{\min}(\boldsymbol{\mu})^{-2}$.*

Proof Take $\boldsymbol{\omega}_0 = \sum_{\mathbf{x} \in \mathcal{X}_0} \mathbf{x} / |\mathcal{X}_0| \in \Sigma$, where $\mathcal{X}_0 = \{\mathbf{i}^*(\mathbf{e}_k) : k \in [K]\}$. Observe that $\boldsymbol{\omega}_0 \geq \mathbf{1}_K / K$ by Lemma 23 (which leads to $\sum_{\mathbf{x} \in \mathcal{X}_0} \mathbf{x} \geq \mathbf{1}_K$ and $1/|\mathcal{X}_0| \geq 1/K$). Thus,

$$F_{\boldsymbol{\mu}}(\boldsymbol{\omega}_0) = \min_{\mathbf{x} \neq \mathbf{i}^*(\boldsymbol{\mu})} \frac{\Delta_{\mathbf{x}}(\boldsymbol{\mu})^2}{2 \langle \mathbf{x} \oplus \mathbf{i}^*(\boldsymbol{\mu}), \boldsymbol{\omega}_0^{-1} \rangle} \geq \frac{\Delta_{\min}(\boldsymbol{\mu})^2}{4KD},$$

where we used Proposition 1 in §3.1 to obtain the equality, and the last inequality is because

$$\langle \mathbf{x} \oplus \mathbf{i}^*(\boldsymbol{\mu}), \boldsymbol{\omega}_0^{-1} \rangle \leq \|\mathbf{x} \oplus \mathbf{i}^*(\boldsymbol{\mu})\|_1 \|\boldsymbol{\omega}_0^{-1}\|_{\infty} \leq \frac{2D}{\min_{k \in [K]} (\boldsymbol{\omega}_0)_k} \leq 2KD.$$

As $T^*(\boldsymbol{\mu})^{-1} = \max_{\boldsymbol{\omega} \in \Sigma} F_{\boldsymbol{\mu}}(\boldsymbol{\omega}) \geq F_{\boldsymbol{\mu}}(\boldsymbol{\omega}_0)$, we then have $T^*(\boldsymbol{\mu}) \leq \frac{4KD}{\Delta_{\min}(\boldsymbol{\mu})^2}$. □

L Extension to the transductive setting

In this section, we extend our results to the transductive combinatorial semi-bandits. In transductive best-arm identification with fixed confidence with semi-bandit feedback [JMKK21], the decision maker is given an exploration set $\mathcal{A} \subseteq \{0, 1\}^K$ and a decision set $\mathcal{X} \subseteq \{0, 1\}^K$ (\mathcal{A} might differ from \mathcal{X}), and at each round, she selects an action in \mathcal{A} to receive a semi-bandit feedback. Her goal is to identify the best action in \mathcal{X} using as few samples as possible.

Notation. Let $\mathcal{M} \subseteq \{0, 1\}^K$ be any set of actions. We use $\mathbf{i}_{\mathcal{M}}^*(\boldsymbol{\mu})$ to denote any maximizer in \mathcal{M} of the linear maximization $\max_{\mathbf{x} \in \mathcal{M}} \langle \mathbf{x}, \boldsymbol{\mu} \rangle$. We also use $\Sigma_{\mathcal{M}} = \{\sum_{\mathbf{x} \in \mathcal{M}} w_{\mathbf{x}} : \mathbf{w} \in \Sigma_{|\mathcal{M}|}\}$.

Sample complexity lower bound. The generalization of Theorem 7 to the transductive setting has been made in [JMKK21]: any δ -PAC algorithm satisfies

$$\mathbb{E}_{\boldsymbol{\mu}}[\tau] \geq T^*(\boldsymbol{\mu}) \text{kl}(\delta, 1 - \delta) \quad \text{with} \quad T^*(\boldsymbol{\mu})^{-1} = \sup_{\boldsymbol{\omega} \in \Sigma_{\mathcal{A}}} \inf_{\boldsymbol{\lambda} \in \text{Alt}(\boldsymbol{\mu})} \left\langle \boldsymbol{\omega}, \frac{(\boldsymbol{\mu} - \boldsymbol{\lambda})^2}{2} \right\rangle. \quad (58)$$

The inner optimization is still with respect to \mathcal{X} while the outer optimization is with respect to the exploration set \mathcal{A} . Refer to Appendix C in [JMKK21] for the proof.

Transductive P-FWS algorithm. Assumption 1 has to be extended. It now needs to ensure that $\mathbf{i}_{\mathcal{A}}^*(\mathbf{v})$ for any $\mathbf{v} \in \mathbb{R}^K$ can be computed in polynomial-time. The P-FWS algorithm also needs to be adapted to the transductive setting. This is done by the following two modifications:

- [K]-covering set: $\mathcal{X}_0 \leftarrow \{\mathbf{i}_{\mathcal{A}}^*(\mathbf{e}_k) : k \in [K]\}$
- FW update: $\mathbf{x}(t) \leftarrow \mathbf{i}_{\mathcal{A}}^* \left(\nabla \tilde{F}_{\hat{\boldsymbol{\mu}}(t-1), \eta_t, n_t}(\hat{\boldsymbol{\omega}}(t-1)) \right)$

Analysis of P-FWS. Let $D_{\mathcal{A}} = \max_{\mathbf{x} \in \mathcal{A}} \|\mathbf{x}\|_1$. The analysis is easily extended by replacing (D, \mathcal{X}) with $(D_{\mathcal{A}}, \mathcal{A})$ in Appendix D, Appendix E, Appendix F and Appendix G whenever the context is subject to the exploration set rather than the decision set.