

Exploration via Feature Perturbation in Contextual Bandits

Seouh-won Yi

Seoul National University
uniqeseouh@snu.ac.kr

Min-hwan Oh

Seoul National University
minoh@snu.ac.kr

Abstract

We propose *feature perturbation*, a simple yet effective exploration strategy for contextual bandits that injects randomness directly into feature inputs, instead of randomizing unknown parameters or adding noise to rewards. Remarkably, this algorithm achieves $\tilde{O}(d\sqrt{T})$ worst-case regret bound for generalized linear contextual bandits, while avoiding the $\tilde{O}(d^{3/2}\sqrt{T})$ regret typical of existing randomized bandit algorithms. Because our algorithm eschews parameter sampling, it is both computationally efficient and naturally extends to non-parametric or neural network models. We verify these advantages through empirical evaluations, demonstrating that feature perturbation not only surpasses existing methods but also unifies strong practical performance with the near-optimal regret guarantees.

1 Introduction

Multi-armed bandits (MABs) provide the canonical model for sequential decision-making under uncertainty: at each round a decision-making agent selects one of several arms to maximize cumulative reward while balancing exploration and exploitation. However, classical MABs ignore side information that often accompanies decisions in practice. Contextual bandits address this limitation by allowing the agent to first observe contextual information and then choose an action tailored to that context—e.g., features of users and/or items inform which arm to pull. This contextualized formulation has become a pivotal framework in online learning and sequential decision-making, with a rich literature on algorithms and guarantees [2, 7, 38, 9, 34].

A widely studied formulation of contextual bandits is the (*generalized*) *linear contextual bandit*, where the expected reward is modeled by a linear function [2, 7, 38, 11, 1] or, more generally, by a generalized linear model (GLM) [17, 39, 26, 37]. In both linear and GLM settings, deterministic methods based on *optimism in the face of uncertainty* (OFU) [7, 38, 1, 39] and randomized approaches such as Thompson Sampling (TS) [10, 6, 3] or Perturbed History Exploration (PHE) [31, 32, 35] have been extensively studied. Notably, OFU-type algorithms achieve near-optimal regret of $\tilde{O}(d\sqrt{T})$ in linear contextual bandits (and likewise in GLM bandits [39]), yet often underperform compared to TS and PHE in practice. In contrast, randomized exploration methods typically exhibit superior empirical performance but suffer from sub-optimal theoretical guarantees: standard analyses confirm a regret bound of $\tilde{O}(d^{3/2}\sqrt{T})$ [3, 6] in the frequentist (worst-case) setting.¹ Crucially, this gap is not merely an artifact of analysis: Hamidi and Bayati [21] show it reflects an inherent cost of randomization in (generalized) linear Thompson sampling. This result highlights a fundamental mismatch between the

¹LinTS [3, 6] achieves a regret of $\tilde{O}(\min(d^{3/2}\sqrt{T}, d\sqrt{T\log K}))$. While Kveton et al. [31] originally showed that LinPHE has a regret of $\tilde{O}(d\sqrt{T\log K})$, where K is the number of arms, a recent work [35] proves that LinPHE also satisfies $\tilde{O}(d^{3/2}\sqrt{T})$ regret. For further discussion on trading off factors of $\mathcal{O}(\sqrt{d})$ vs. $\mathcal{O}(\sqrt{\log K})$, see Agrawal and Goyal [6]. In this work, we consider even a large action space with $K > e^d$, so the $\tilde{O}(d^{3/2}\sqrt{T})$ regret of randomized (generalized) linear bandit algorithms is the main focus of discussion.

existing randomized exploration and the tighter optimism mechanism in OFU-based approaches. This dichotomy prompts a natural question: *is it possible to close the gap between randomized exploration and $\tilde{O}(d\sqrt{T})$ worst-case regret?* If one adheres strictly to randomly sampling the parameters (as in TS) or perturbing the observed rewards (as in PHE), there appears a fundamental barrier [21] preventing regret from achieving $\tilde{O}(d\sqrt{T})$.

In this work, we propose a simple yet effective alternative: instead of sampling parameters or perturbing rewards, we randomly perturb the observed *features* (or contexts). By shifting the focus of exploration from parameter space to feature space, we circumvent the limitations that impose higher regret on existing randomized algorithms. Remarkably, our analysis shows that this new approach not only retains the empirical advantages of randomization but also achieves $\tilde{O}(d\sqrt{T})$ worst-case regret in (generalized) linear bandit settings, with no additional dependence on the number of arms. Furthermore, our method avoids the computational overhead of sampling parameters, making it attractive for a wide range of real-world applications.

Beyond theoretical efficiency, feature perturbation can seamlessly extend to more flexible or non-parametric reward models, including neural networks. We demonstrate this empirically, showing that feature-based randomization can drive effective exploration even when specific parametric model assumptions may not hold. Hence, our proposed approach unifies strong theoretical guarantees with practical efficacy in (generalized) linear contextual bandits, and extends practically to more complex models. Our main contributions are summarized as follows:

- **Feature perturbation for contextual bandits.** We introduce a *new class of algorithms* for randomized exploration, termed *feature perturbation*, which focuses on perturbing feature inputs rather than parameters or rewards. This approach is straightforward to implement and conceptually distinct from existing randomized exploration strategies.
- **Tight regret bounds.** To the best of our knowledge, our work is the first *randomized algorithm* for generalized linear contextual bandits that achieves: (i) a regret bound of $\tilde{O}(d\sqrt{T})$, matching the best-known guarantees of deterministic (OFU-based) methods; and simultaneously (ii) benefiting from an instance-dependent constant κ . Notably, our algorithm’s regret does not increase (even logarithmically) with the number of arms.
- **Empirical validation.** Through extensive experiments on both synthetic and real-world data, we show that feature perturbation not only performs competitively against existing randomized methods but also generalizes beyond parametric models (e.g., deep neural networks), demonstrating robustness even when linear assumptions do not hold.

2 Related works

Contextual bandits have been extensively investigated under various modeling assumptions. In the *linear* bandit setting, deterministic methods based on OFU [1, 7] achieve near-optimal $\tilde{O}(d\sqrt{T})$ regret, but often exhibit conservative exploration in practice. By contrast, *randomized* algorithms such as TS [3, 6, 10] and PHE [31, 32] typically show better empirical performance yet suffer from a higher $\tilde{O}(d^{3/2}\sqrt{T})$ regret bound. Notably, Hamidi and Bayati [21] demonstrated that the extra \sqrt{d} inflation in TS-type algorithms is unavoidable in worst-case scenarios: eliminating this factor would lead to a linear dependence on T . Consequently, parameter-based randomization cannot, in general, achieve $\tilde{O}(d\sqrt{T})$ regret without further modifications.

Generalized linear bandits (GLB; [17, 39]) extend linear bandits to settings where rewards follow a nonlinear link function. UCB- and TS-based approaches [3, 17, 32, 46] have also been applied here, displaying the same contrast between deterministic and randomized exploration. While UCB-type methods reach $\tilde{O}(d\sqrt{T})$ regret, they tend to over-explore in practice; randomized strategies mitigate this over-exploration but retain an additional \sqrt{d} penalty in the worst case. Like their linear counterparts, these methods rely on sampling the unknown parameter or perturbing rewards rather than altering the feature representation.

By contrast, our work introduces a new class of *feature-perturbation* (FP) algorithms designed to circumvent the dimensional penalty inherent in standard randomized approaches. Instead of randomizing parameters or rewards, we propose to perturb the features directly. This perspective not

only preserves the empirical robustness associated with randomized strategies but also achieves a tight regret bound in both linear and generalized linear settings—thereby reconciling the theoretical and practical advantages of contextual bandit exploration.

3 Preliminaries

Notations. For vectors $x, y \in \mathbb{R}^d$, let $\|x\|$ denote the 2-norm and $\|x\|_A = \sqrt{x^\top A x}$ the weighted norm for a positive definite matrix $A \in \mathbb{R}^{d \times d}$. The inner product is $x^\top y = \langle x, y \rangle$, and the weighted version is $x^\top A y = \langle x, y \rangle_A$. The notation \tilde{O} hides logarithmic factors in big-O notation, retaining instance-dependent constants. For a real-valued function f , we write \dot{f} and \ddot{f} to denote its first and second derivatives. The set $\{1, \dots, K\}$ is abbreviated as $[K]$.

Generalized linear contextual bandit. A *generalized linear model* (GLM; [41]) describes a response $r \in \mathbb{R}$ drawn from an exponential-family distribution with mean $\mu(x^\top \theta^*)$, where $x \in \mathbb{R}^d$ is a feature vector and $\theta^* \in \mathbb{R}^d$ is an unknown parameter. Given differentiable functions g and h , and a base measure ν , the conditional density of r given x takes the form:

$$dp(r \mid x; \theta^*) = \exp(r x^\top \theta^* - g(x^\top \theta^*) + h(r)) d\nu, \quad (1)$$

where the derivative of g defines the *link function* μ .² Let $\mathcal{H}_{t-1} := \sigma(\{(x_\tau, r_\tau)\}_{\tau=1}^{t-1})$ denote the filtration up to round $t-1$. We define $\mathbb{P}_t(\cdot) := \mathbb{P}(\cdot \mid \mathcal{H}_{t-1})$ and $\mathbb{E}_t[\cdot] := \mathbb{E}[\cdot \mid \mathcal{H}_{t-1}]$. The negative log-likelihood and the *maximum likelihood estimate* (MLE) at round t are then given by:

$$L_t(\theta) = \sum_{\tau=1}^{t-1} (g(x_\tau^\top \theta) - r_\tau x_\tau^\top \theta), \quad \hat{\theta}_t := \operatorname{argmin}_{\theta \in \Theta} L_t(\theta).$$

In the generalized linear contextual bandit (GLB) setting, the agent observes a context $c_t \in \mathcal{C}$ and a corresponding set of feature vectors $\mathcal{X}_t \subset \mathbb{R}^d$ representing each allowable arm $a \in \mathcal{A}(c_t)$ at each round. Upon selecting $x_t \in \mathcal{X}_t$, the agent receives a stochastic reward $r_t \sim p(\cdot \mid x_t; \theta^*)$. The learner aims to minimize the regret: $R(T) = \sum_{t=1}^T (\mu(x_{t^*}^\top \theta^*) - \mu(x_t^\top \theta^*))$, where $x_{t^*} := \operatorname{argmax}_{x \in \mathcal{X}_t} \mu(x^\top \theta^*)$ is the optimal arm at round t , which depends on the context c_t .

4 Algorithm: GLM-FP

Algorithm 1 GLM-FP: Feature Perturbation in Generalized Linear Bandits

- 1: **Input:** Regularization parameter $\lambda > 0$, tuning parameter $\{c_t\}$
 - 2: **for** $t = 1, 2, \dots, T$ **do**
 - 3: Compute $\hat{\theta}_t = \operatorname{argmin}_{\theta \in \mathbb{R}^d} L_t(\theta; \{x_\tau, r_\tau\}_{\tau=1}^{t-1})$
 - 4: Sample $\zeta_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 5: Compute $\tilde{x}_{ti} = x_{ti} + c_t \cdot \frac{\|x_{ti}\|_{\hat{H}_t^{-1}}}{\|\hat{\theta}_t\|} \cdot \zeta_t$ for all i
 - 6: Choose $i_t = \operatorname{argmax}_{i \in [\mathcal{X}_t]} \mu(\tilde{x}_{ti}^\top \hat{\theta}_t)$ and observe reward r_t ▷ Let $x_t := x_{t, i_t}$
 - 7: **end for**
-

At each step t , given the filtration \mathcal{H}_{t-1} , the algorithm computes the MLE $\hat{\theta}_t$ (line 3). Since MLE lacks a closed-form solution, we employ *Sequential Quadratic Programming* (SQP; [14]) or *Iteratively Reweighted Least Squares* (IRLS; [47]). Instead of perturbing rewards or parameters, the algorithm injects controlled randomness into feature vectors using a perturbing distribution \mathcal{D} . By default, we use a multivariate normal distribution (line 4), but any distribution satisfying concentration and anti-concentration properties can be employed, as analyzed in Section 5.2.

A fundamental property of the GLB is that the strictly increasing link function simplifies the problem structure, making it resemble linear bandits. While prior works [3, 39, 46] adopt near-identical

²We normalize the reward model so that $\operatorname{Var}[r_t \mid x_t] = \mu(x_t^\top \theta^*)$. Scaling the variance by σ^2 (by inserting σ^2 in the denominator of the exponential term of Eq. (1)) accordingly yields a σ -inflated regret bound.

methods to linear bandits—typically using the *vanilla Gram matrix* $V_t = \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} x_\tau x_\tau^\top$ —our approach takes a more refined direction. To precisely control the perturbation magnitude for each feature representation, our approach utilizes a weighted Gram matrix, $\hat{H}_t := \lambda \mathbf{I} + \nabla^2 L_t(\hat{\theta}_t)$.

This adaptive structure, modulated by a tunable parameter c_t , enables a more targeted perturbation strategy. Using the resulting scaling factors, the algorithm perturbs each feature vector to construct a set of perturbed vectors, $\{\tilde{x}_{ti}\}$ for reachable arms at round t (line 5). Importantly, the perturbation noise ζ_t is shared across all arms, coupling each feature vector with the same random variable. This design eliminates the explicit dependence on K in the regret bound, thereby yielding stronger theoretical guarantees and improved empirical performance. Finally, the algorithm selects the arm that maximizes $\mu(\tilde{x}_{ti}^\top \hat{\theta}_t)$ and updates its history upon observing the reward r_t (lines 6).

4.1 Extension to general function class

The algorithm extends naturally to more flexible function classes, as described in Algorithm C.1. Under the realizability assumption (i.e., $f^* \in \mathcal{F}$), the estimate $\hat{f}_t \in \mathcal{F}$ is obtained via a least squares oracle on the history \mathcal{H}_{t-1} [19, 45, 49]. Given the structural of the bandit reward model, a sampling distribution is then defined for each arm and used to construct a set of perturbed contexts $\{\tilde{x}_{ti}\}$. The algorithm selects the arm that maximizes $\hat{f}_t(\tilde{x}_{ti})$, thereby balancing exploration and exploitation.

The proposed GLM-FP algorithm is an instance of this framework, characterized by two specific design choices: (i) a Gaussian sampling distribution $\mathcal{D}(x_{ti}, \Sigma_{ti}) \triangleq \mathcal{N}(x_{ti}, \Sigma_{ti})$ centered at x_{ti} with elliptical covariance scaling with the uncertainty in the direction of x_{ti} and normalized by the estimated parameter norm; and (ii) a *coupled* perturbation scheme, where a single shared random vector ζ_t perturbs all arms simultaneously, in contrast to perturbing each arm independently.

4.2 Intuition behind the algorithm

In contextual bandit problems, randomized algorithms can be broadly categorized into two types: those that introduce randomness into the underlying model and those that inject randomness directly into the estimated expected rewards. The former category includes methods such as PHE [31] and TS [3, 6], which compute perturbed model parameters to induce exploration. However, as noted by Hamidi and Bayati [21], this approach can be suboptimal even in linear settings. The latter category includes algorithms such as RandUCB [46], where the model is trained deterministically and randomness is introduced by adding stochastic bonuses to the estimated rewards of each arm. While effective in inducing exploration, such reward perturbation can violate the inductive bias of the function class, since the resulting scores for arms may not be realized simultaneously by any single model $f \in \mathcal{F}$ —even in simple linear cases.

To address these limitations, we propose an alternative strategy that retains the estimated model \hat{f} from past data and introduces randomness through input perturbations at decision time. Exploring directly in the feature space preserves the structural assumptions of the function class and remains effective even in overparameterized settings where $p \gg d$, such as neural bandits. It also aligns naturally with real-world scenarios in which contextual features contain inherent noise [8, 28, 29].

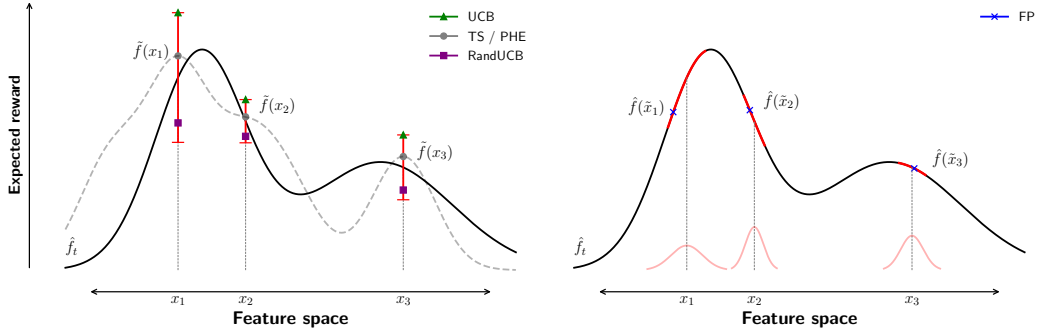


Figure 1: (Left) Model perturbation methods randomize rewards via off-estimated models \tilde{f}_t . (Right) Feature Perturbation (FP) perturbs inputs and evaluates them with a fixed model \hat{f}_t .

As illustrated in Figure 1, whereas UCB, TS, PHE, and RandUCB randomize parameters or reward estimates via modified models \tilde{f}_t , GLM-FP explores through feature-space perturbations evaluated under a fixed \hat{f} . This structural shift preserves the inductive bias of the model class and decouples exploration from parameter uncertainty, a property that becomes crucial in high-dimensional regimes.

5 Regret analysis of GLM-FP

In this section, we establish the regret guarantee for the proposed algorithm, GLM-FP, and outline the key steps in its proof. We introduce fundamental concepts and lemmas that form the backbone of our analysis. We begin by presenting the standard assumptions commonly adopted in the analysis of generalized linear bandit algorithms [1, 3, 4, 6, 15, 17, 32, 37, 39, 52].

Assumption 1 (Boundedness). *The feature space \mathcal{X} and the parameter space Θ are compact subsets of \mathbb{R}^d . For any $x \in \mathcal{X}_{[T]}$ and $\theta^* \in \Theta$, we have $\|x\| \leq 1$ and $\|\theta^*\| \leq 1$.*

Assumption 2 (Self-concordance). *The function g is three times differentiable, and its derivative $\dot{g} = \mu$ is strictly increasing and \mathcal{L}_μ -Lipschitz continuous. Furthermore, g is self-concordant, characterized by the constant $M_\mu := \sup_{x \in \mathcal{X}, \theta \in \Theta} |\ddot{\mu}(\langle x, \theta \rangle)| / \dot{\mu}(\langle x, \theta \rangle)$.*

Many widely studied GLB instances naturally satisfy these assumptions [43]. For example, the triple $(\mu(z), \mathcal{L}_\mu, M_\mu)$ takes the form $(z, 1, 0)$ in linear, $(\frac{1}{1+e^{-z}}, \frac{1}{4}, 1)$ in logistic, and $(e^z, e, 1)$ in Poisson.

5.1 Confidence bound for the true parameter

An important step in analyzing the regret bound of the algorithm is to establish a confidence set for the underlying parameter θ^* . This involves constructing a region that reliably contains θ^* throughout the learning process. To obtain a practical and tighter bound, we adopt confidence sets derived from the log-likelihood function using an *ellipsoidal relaxation*. The confidence width is provided by recent work [37], which can be substituted with alternative, potentially tighter bounds.

Lemma 1 (Adapted from Theorem 3.2. in Lee et al. [37]). *Let $\mathcal{L}_t := \max_{\theta \in \Theta} \|\nabla L_t(\theta)\|$ denote the Lipschitz constant of the loss function $L_t(\cdot)$.³ For any $\lambda > 0$, define the regularized Hessian matrix at $\hat{\theta}_t$ as $\hat{H}_t := \nabla^2 L_t(\hat{\theta}_t) + \lambda \mathbf{I}$. Then, with probability at least $1 - \delta$, for all $t \geq 1$, it holds that $\theta^* \in \Theta_t(\delta, \lambda) := \{\theta \in \mathbb{R}^d \mid \|\theta - \hat{\theta}_t\|_{\hat{H}_t} \leq \beta_t(\delta)\}$, where*

$$\beta_t(\delta) = \sqrt{4\lambda + 2(1 + M_\mu) \left(\log \frac{1}{\delta} + d \log \left(\frac{2e\mathcal{L}_t}{d} \right) \right)}.$$

5.2 Concentration and anti-concentration

In addition to ensuring that $\hat{\theta}_t$ remains close to the true parameter θ^* , it is crucial to balance the degree of randomization in the algorithm. Our regret analysis relies on showing that, with an appropriate choice of the tuning parameter c_t , the perturbed feature \tilde{x}_t is stochastically optimistic, while concentrated around its estimated value $\mu(x_t^\top \hat{\theta}_t)$. These properties are fundamental to the analysis, and we introduce the relevant components in this section.

Definition 1. *Let $t \in [T]$. We define the following events:*

- (i) \hat{E}_t : $\hat{\theta}_\tau$ remains close to θ^* for all steps $\tau \leq t$.
- (ii) \tilde{E}_t : all perturbed vectors $\tilde{x}_{\tau i}$ concentrated around their corresponding $x_{\tau i}$ for all steps $\tau \leq t$.

For a given confidence level $\delta \in (0, 1)$, define $\delta' = \delta/(4T)$ and $\gamma_t(\delta) = \beta_t(\delta') \sqrt{c \log(c'/\delta)}$, where c and c' are constants consistent with Eq. (2). The events are formally defined as:

$$\hat{E}_t := \left\{ \forall \tau \leq t; \|\hat{\theta}_\tau - \theta^*\|_{\hat{H}_\tau} \leq \beta_t(\delta') \right\} \quad \text{and} \quad \tilde{E}_t := \left\{ \forall \tau \leq t, x_{\tau i} \in \mathcal{X}_\tau; \tilde{x}_{\tau i} \in \mathcal{E}_\tau(x_{\tau i}) \right\},$$

where $\mathcal{E}_t(x) := \{\tilde{x} \in \mathbb{R}^d \mid |\langle \tilde{x} - x, \hat{\theta}_t \rangle| \leq \gamma_t(\delta') \|x\|_{\hat{H}_t^{-1}}\}$ represents a high-probability region for the perturbed feature vector associated with each arm x .

³It has been shown by Lee et al. [37] that $\mathcal{L}_t = \mathcal{O}(t)$ for linear, logistic, and Poisson bandit instances.

A key requirement for the perturbation distribution \mathcal{D} is the following concentration property: for some constants $c, c' > 0$ and any unit vector u , we have

$$\mathbb{P}_{\zeta \sim \mathcal{D}} \left(|u^\top \zeta| \leq \sqrt{c \log(c'/\delta)} \right) \geq 1 - \delta. \quad (2)$$

Remark. The perturbing distribution of TS on θ is described in Appendix B. Unlike in TS, the concentration here is evaluated along a fixed direction u , not over all coordinates. This avoids the need for a union bound over d dimensions, and the resulting bound is independent of d . This dimensionality reduction is a key advantage of perturbing features instead of parameters.

By construction, the events satisfy the nested structure $\hat{E}_T \subset \dots \subset \hat{E}_1$ and $\tilde{E}_T \subset \dots \subset \tilde{E}_1$. Building on these definitions, we show that the proposed perturbation distribution induces an appropriate balance between exploration and exploitation, as formalized in the following lemmas.

Lemma 2 (Concentration). *Under Assumptions 1 and 2, with $c_t = \beta_t(\delta')$, $\mathbb{P}(\hat{E}_T \cap \tilde{E}_T) \geq 1 - \frac{\delta}{2}$.*

Lemma 3 (Stochastic optimism). *For each round t , given that the events \hat{E}_t and \tilde{E}_t occur, the probability of anti-concentration, conditioned on the filtration \mathcal{H}_{t-1} and under Assumptions 1 and 2, is lower bounded as*

$$\mathbb{P}_t(\mu(\tilde{x}_t^\top \hat{\theta}_t) \geq \mu(x_{t*}^\top \theta^*) \mid \hat{E}_t, \tilde{E}_t) \geq \frac{1}{4\sqrt{e\pi}}.$$

In our construction, the concentration term in Eq. (2) induces only a constant-order inflation, yielding both β_t and γ_t scaling as $\mathcal{O}(\sqrt{d})$. In contrast, TS incurs an extra \sqrt{d} factor, resulting in $\gamma_t = \mathcal{O}(d)$. This leads to a tighter exploration term in our algorithm and improved regret performance.

5.3 Regret bound of GLM-FP

The complexity of the GLB problem is fundamentally determined by the following quantities, which captures the degree of nonlinearity in the reward function:

$$\kappa_* := \frac{\sum_{t=1}^T \dot{\mu}(x_{t*}^\top \theta^*)}{T}, \quad \kappa := \min_{x \in \mathcal{X}_{[T]}, \theta \in \Theta} \dot{\mu}(x^\top \theta), \quad \text{where } \mathcal{X}_{[T]} := \bigcup_{t=1}^T \mathcal{X}_t. \quad (3)$$

These may scale exponentially small, particularly in the case of logistic bandits [15]. The following theorem presents the regret guarantee for our algorithm.

Theorem 1. *For all $\delta \in (0, 1)$, define $\delta' = \delta/(4T)$. Under Assumptions 1 and 2, with $c_t = \beta_t(\delta')$ and $\lambda = \mathcal{O}(d)$, the cumulative regret $R(T)$ is bounded with probability at least $1 - \delta$ as follows:*

$$R(T) = \tilde{\mathcal{O}} \left(d\sqrt{\kappa_* T} + d^2/\kappa \right).$$

Discussion of Theorem 1. The leading term of the regret guarantee is $\tilde{\mathcal{O}}(d\sqrt{\kappa_* T})$, which matches the minimax optimal regret bound in terms of the dimensionality d , the horizon T , and the instance-dependent constant κ_* [4, 37]. While RandUCB [46], another randomized algorithm, also achieves a regret bound of $\tilde{\mathcal{O}}(d\sqrt{T})$, it is penalized by its inverse dependence on κ , lacking adaptation to instance-dependent complexity. In contrast, to the best of our knowledge, our result is the first to show that a randomized algorithm achieves a regret bound with linear d -dependency, without additional dependence on the number of arms, benefiting from κ_* in GLB problems.

5.4 Proof sketch of Theorem 1

Our proof begins by decomposing the instantaneous regret into two components: Reg_{FP} , the regret which arises from the randomization through FP, and Reg_{EST} , the regret accounting for the estimation error. By bounding each component separately, we derive the overall regret bound:

$$R(T) = \underbrace{\sum_{t=1}^T \left(\overbrace{\mu(x_{t*}^\top \theta^*) - \mu(\tilde{x}_t^\top \hat{\theta}_t)}^{\text{Optimism}} + \overbrace{\mu(\tilde{x}_t^\top \hat{\theta}_t) - \mu(x_t^\top \hat{\theta}_t)}^{\text{Perturbation Concentration}} \right)}_{\text{Reg}_{\text{FP}}} + \underbrace{\sum_{t=1}^T \left(\mu(x_t^\top \hat{\theta}_t) - \mu(x_t^\top \theta^*) \right)}_{\text{Reg}_{\text{EST}}}. \quad (4)$$

To bound Reg_{FP} , we rely on two key properties: (i) optimism induced by selecting the best estimated arm, and (ii) concentration of the perturbation around the original context. Under the high-probability events \hat{E}_T and \tilde{E}_T , both properties are well controlled, leading to the following bound:

$$\text{Reg}_{\text{FP}} \leq (8\sqrt{e\pi} + 1) \sum_{t=1}^T \left| \max_{x \in \mathcal{E}_t(x_t)} \mu(x^\top \hat{\theta}_t) - \mu(x_t^\top \hat{\theta}_t) \right| + \tilde{\mathcal{O}} \left(\sqrt{\frac{d^2 T}{\lambda^2}} \right),$$

where the $\tilde{\mathcal{O}}(\cdot)$ term results from an Azuma–Hoeffding concentration and becomes negligible under the choice $\lambda = \mathcal{O}(d)$. Note that since the stochastic optimistic probability in Lemma 3 is lower bounded by a constant, Reg_{FP} can be effectively bounded by the sum of per-round concentration widths, multiplied by a constant that is independent of d and T .

Unlike previous randomized algorithms [3, 32, 46] that linearize the reward function and thereby suffer regret bounds inversely proportional to κ , our analysis avoids such linearization. Instead, we directly utilize the gradient of the link function μ to characterize the shape of the elliptical confidence region, enabling more efficient exploration tailored to the reward model. This is reflected in the weighted Gram matrix \hat{H}_t , where the curvature of μ enters via the term $\dot{\mu}(x_\tau^\top \hat{\theta}_t)$.

However, this construction introduces a dependency on t (rather than solely on τ) in the weighted Gram matrix, which precludes a direct application of the standard *Elliptical Potential Lemma* (EPL). To address this, we build upon recent analytical developments [4, 15, 35], and introduce a lower envelope of derivatives by defining $\bar{\theta}_t$ as the minimizer of $\dot{\mu}(x_t^\top \theta)$ over the union $\cup_{\tau \in [t, T]} \Theta_\tau(\delta, \lambda)$.

This yields a matrix $\bar{H}_t := \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} \dot{\mu}(x_\tau^\top \bar{\theta}_\tau) x_\tau x_\tau^\top$ which satisfies $\hat{H}_t \succeq \bar{H}_t$, thereby enabling the application of EPL. We bound Reg_{EST} in a similar fashion and show that both Reg_{FP} and Reg_{EST} admit the same upper bound. This reduces the analysis to solving a quadratic inequality of the form:

$$\text{Reg}_{\text{max}} \leq A\sqrt{B + C\text{Reg}_{\text{max}}} + D, \quad \text{where} \quad \text{Reg}_{\text{max}} = \max\{\text{Reg}_{\text{FP}}, \text{Reg}_{\text{EST}}\}.$$

The constants A , B , C , and D are explicitly analyzed in the full proof, deferred to Appendix E.

6 Carving off the \sqrt{d} factor compared to TS

The proposed algorithm, GLM-FP, adopts a novel exploration strategy by perturbing the input feature vectors, in contrast to conventional randomized algorithms such as Thompson Sampling (TS), which introduce randomness into the model parameter θ . This design yields a regret bound with linear dependence on d , whereas TS incurs a higher $\mathcal{O}(d^{3/2})$ dependence. We examine the origin of this discrepancy by comparing the linear variants of both algorithms (see Appendix C.2), highlighting how each introduces randomness to facilitate exploration.

The randomized evaluation score $\tilde{f}_t(x_i)$ (either $x_{ti}^\top \hat{\theta}_t$ for TS or $\tilde{x}_{ti}^\top \hat{\theta}_t$ for FP) for each arm used for action selection in both algorithms is straightforward to compute as follows:

$$(\text{TS}) \quad \tilde{f}_t(x_i) = x_{ti}^\top \tilde{\theta}_t = x_{ti}^\top \hat{\theta}_t + c_t x_{ti}^\top V_t^{-1/2} \zeta_t, \quad (\text{FP}) \quad \tilde{f}_t(x_i) = \tilde{x}_{ti}^\top \hat{\theta}_t = x_{ti}^\top \hat{\theta}_t + c_t z_t \|x_{ti}\|_{V_t^{-1}},$$

where $\zeta_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $z_t \sim \mathcal{N}(0, 1)$. Thus, for each arm individually, both methods induce the same Gaussian distribution,⁴ $\tilde{f}_t(x_i) \sim \mathcal{N}(x_{ti}^\top \hat{\theta}_t, c_t^2 \|x_{ti}\|_{V_t^{-1}}^2)$. However, the *object of perturbation*—parameter in TS versus feature in FP—fundamentally alters how exploration bonuses are assigned and how arm comparisons are coupled at each timestep.

In TS, the bonus $\langle x_{ti}, \zeta_t \rangle_{V_t^{-1/2}}$ projects each arm onto a shared random direction ζ_t in the $V_t^{-1/2}$ -transformed space. As conceptually illustrated in Figure 2a, this shared dependence can produce counterintuitive effects: well-explored arms may occasionally receive large bonuses simply due to alignment with ζ_t , while under-explored ones may be neglected. Because the same ζ_t governs all arms, the analysis must ensure uniform reliability of exploration across directions, requiring high-probability control of the d -dimensional Gaussian vector ζ_t . Applying a union bound over d coordinates introduces an additional \sqrt{d} factor into the regret bound.

⁴In the linear bandit setting, this distribution also matches RandUCB [46], though its derivation is conceptually distinct and diverges beyond the GLB case; see Section 4.2 and Appendix A.

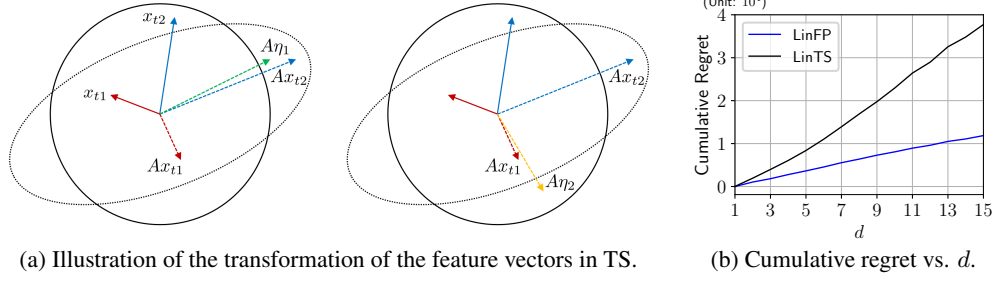


Figure 2: **(a)** Transformation of the well-explored arm x_{t1} and under-explored arm x_{t2} using $A = D^{1/2}P^\top$, where $V_t^{-1/2} = PDP^\top$. Left: ζ_1 induces a proper bonus. Right: ζ_2 reverses the effect. **(b)** Average terminal regret $R(T)$ over 100 runs with $T = 200,000$, $K = 50$, and varying d .

FP, in contrast, decouples exploration from directional uncertainty. Its bonus $z_t \cdot \|x_{ti}\|_{V_t^{-1}}$ scales directly with per-arm uncertainty, ensuring under-explored arms systematically receive larger bonuses. Since randomness enters only through the scalar z_t , the analysis reduces to bounding a one-dimensional Gaussian projection $u^\top \zeta_t$ for some fixed unit vector u .

From an equation-level viewpoint, the regret of TS also admits the decomposition in Eq. (4) [3]. For both algorithms, the estimation component Reg_{EST} is bounded in the same manner, and the anti-concentration event occurs with the same probability p . Consequently, the optimism-driven term in the regret scales with the concentration width divided by p . The essential difference therefore lies in the concentration width—or, equivalently, in how each algorithm controls the perturbation magnitude that also determines the second part of the decomposition. We express the upper bound for each perturbation term as

$$\begin{aligned}
 \text{(TS)} \quad & |x_t^\top (\tilde{\theta}_t - \hat{\theta}_t)| = c_t |x_{ti}^\top V_t^{-1/2} \zeta_t| \leq c_t \|x_{ti}\|_{V_t^{-1}} \cdot \|V_t^{-1/2} \zeta_t\|_{V_t} = c_t \|x_{ti}\|_{V_t^{-1}} \cdot |\zeta_t|, \\
 \text{(FP)} \quad & |(\tilde{x}_t - x_t)^\top \hat{\theta}_t| = c_t \left| \|x_t\|_{V_t^{-1}} \frac{\hat{\theta}_t^\top \zeta_t}{|\hat{\theta}_t|_2} \right| = c_t \|x_{ti}\|_{V_t^{-1}} \cdot |u^\top \zeta_t|.
 \end{aligned}$$

The concentration width in TS depends on $\|\zeta_t\|_2$, the norm of a d -dimensional Gaussian vector, whereas that of FP scales with the one-dimensional projection $|u^\top \zeta_t|$. To obtain a uniform high-probability guarantee across all coordinates, a union bound over the d -dimensional perturbation space introduces an additional $\mathcal{O}(\sqrt{d})$ factor for TS. Thus, despite having identical marginal distributions for individual arms, the two algorithms differ fundamentally in how their perturbations couple across arms: TS requires concentration over all directions in \mathbb{R}^d , whereas FP relies on a single scalar randomization. This structural decoupling eliminates the extraneous \sqrt{d} factor, yielding linear $\mathcal{O}(d)$ dependence in the regret bound and clarifying the geometric origin of FP’s improvement over TS.

7 Experiments

We conduct experiments in two contextual bandit settings: (i) generalized linear bandits (GLBs), including linear and logistic models, and (ii) nonlinear contextual bandits based on neural networks. In each setting, we compare the proposed method with state-of-the-art baselines across varying feature dimensions and datasets. All results are averaged over 100 independent runs to ensure robustness. Detailed experimental setups are provided in Appendix H.

7.1 Generalized linear bandits

We evaluate GLM-FP in both linear and logistic contextual bandit settings, where the expected reward follows a generalized linear model. In the linear bandit case, the reward is generated as $r_t = x_t^\top \theta^* + \varepsilon_t$ with $\varepsilon_t \sim \mathcal{N}(0, 1)$, while in the logistic bandit case, $r_t \sim \text{Bernoulli}(\mu(x_t^\top \theta^*))$ with the logistic function μ . We compare against widely used baselines including ε -greedy, UCB, TS, PHE, and RandUCB. Parameter estimation is performed via regularized weighted least squares (WLS) in linear bandit setting or IRLS in logistic bandit setting.

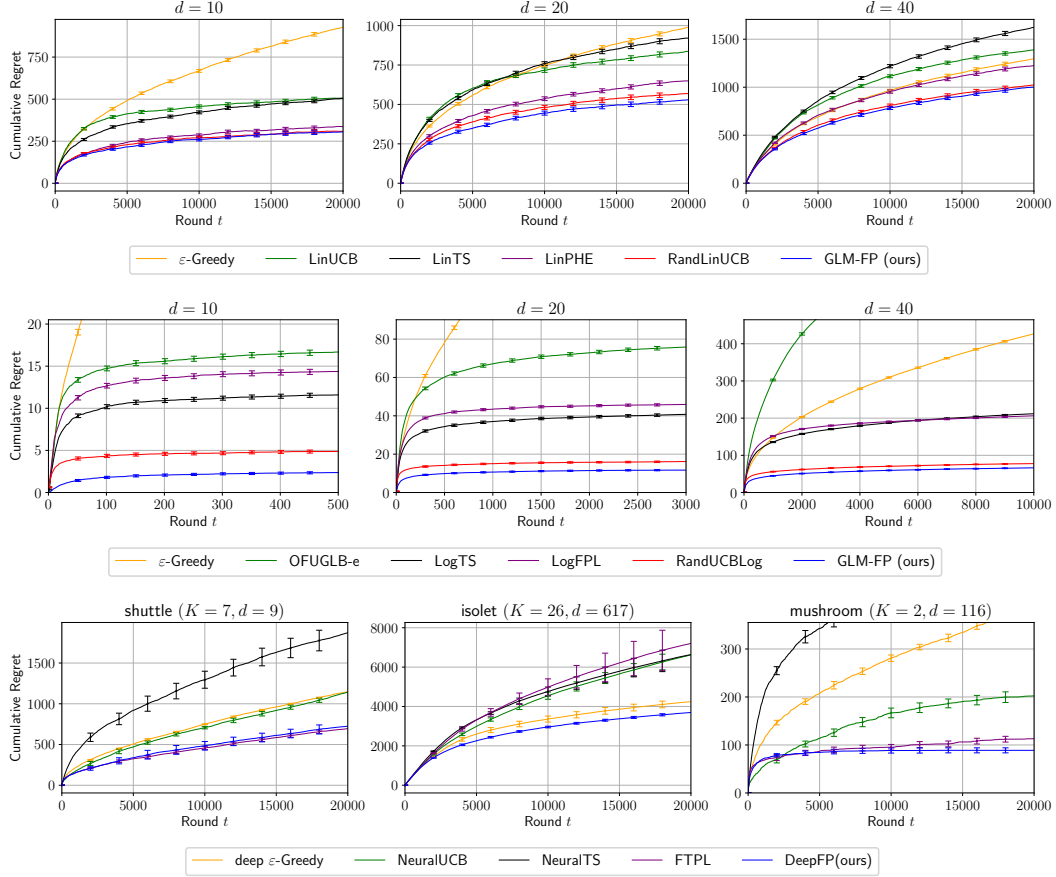


Figure 3: Comparison of cumulative regret across contextual bandit algorithms: linear (top), logistic (middle), and neural (bottom).

As shown in Figure 3 (top and middle), GLM-FP consistently achieves the lowest regret across all tested dimensions. While RandLinUCB performs competitively in the linear case, GLM-FP exhibits superior robustness, particularly in fixed-arm settings or environments with non-stationary arm sets (see Appendix H for more details). In the logistic setting, GLM-FP outperforms all baselines across all tested configurations, demonstrating its effectiveness even when the utility scale (i.e., $x^\top \theta^*$, the inner product term inside $\mu(\cdot)$) or the reward noise variance is varied. These results highlight the reliability and adaptability of our approach across diverse GLB scenarios.

7.2 Neural bandits

We extend the feature perturbation framework to the neural contextual bandit setting through a simple and scalable algorithm, DeepFP. At each round, DeepFP injects independent, arm-wise Gaussian noise into the input features before prediction: $\tilde{x}_{ti} = x_{ti} + \zeta_{ti}$, where $\zeta_{ti} \sim \mathcal{N}(0, I/t)$. The variance decay $1/t$ reflects the diminishing feature uncertainty over time, analogous to the confidence scaling $\|x\|_{\hat{H}_t^{-1}}$ in GLBs. This approach enables exploration without requiring access to model parameters or gradients. We evaluate DeepFP on three UCI benchmark datasets: *shuttle* (7 classes, 9 features), *isolet* (26 classes, 617 features), and *mushroom* (binary, 112 features). Baselines include ϵ -greedy, NeuralUCB, NeuralTS, and Follow-the-Perturbed-Leader (FTPL).

Unlike NeuralTS, which adds randomness to predicted rewards (essentially functioning as a randomized UCB variant), or FTPL, which perturbs historical rewards, DeepFP directly perturbs the *input features*. This structural distinction simplifies implementation and preserves the exploration intent in a more principled way. Furthermore, unlike parameter-sampling-based approaches whose application becomes increasingly unstable in high-capacity models where $p \gg d$, DeepFP perturbs a much lower-dimensional space, making it scalable and numerically stable.

As shown in Figure 3 (bottom), DeepFP outperforms all baselines across the three datasets. Notably, it achieves strong performance without relying on posterior approximations or gradient-based confidence intervals, demonstrating that simple feature perturbation can remain effective even in complex, high-dimensional settings.

8 Conclusion

We introduced a new paradigm for randomized exploration in contextual bandits. By shifting stochasticity to feature space, FP bridges the gap between randomized and optimistic methods while achieving optimal regret. This perspective offers a broadly applicable exploration principle and invites future work on leveraging structured randomness for efficient decision making.

Acknowledgements

This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (No. RS-2022-NR071853 and RS-2023-00222663), by Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No. RS-2025-02263754), and by AI-Bio Research Grant through Seoul National University.

References

- [1] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24:2312–2320, 2011.
- [2] Naoki Abe and Philip M Long. Associative reinforcement learning using linear probabilistic concepts. In *International conference on machine learning*, pages 3–11, 1999.
- [3] Marc Abeille and Alessandro Lazaric. Linear Thompson Sampling Revisited. In Aarti Singh and Jerry Zhu, editors, *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pages 176–184. PMLR, PMLR, 20–22 Apr 2017.
- [4] Marc Abeille, Louis Faury, and Clément Calauzènes. Instance-wise minimax-optimal algorithms for logistic bandits. In *International Conference on Artificial Intelligence and Statistics*, 2020.
- [5] Marc Abeille, David Janz, and Ciara Pike-Burke. When and why randomised exploration works (in linear bandits). In Gautam Kamath and Po-Ling Loh, editors, *Proceedings of The 36th International Conference on Algorithmic Learning Theory*, volume 272 of *Proceedings of Machine Learning Research*, pages 4–22. PMLR, 24–27 Feb 2025.
- [6] Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear payoffs. In *International conference on machine learning*, pages 127–135. PMLR, 2013.
- [7] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2):235–256, 2002.
- [8] Hamsa Bastani and Mohsen Bayati. Online decision making with high-dimensional covariates. *Oper. Res.*, 68(1):276–294, January 2020. ISSN 0030-364X.
- [9] Sébastien Bubeck and Nicolò Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012. ISSN 1935-8237.
- [10] Olivier Chapelle and Lihong Li. An empirical evaluation of thompson sampling. *Advances in neural information processing systems*, 24, 2011.
- [11] Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages 208–214. JMLR Workshop and Conference Proceedings, 2011.

- [12] Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics*, 2019.
- [14] Jean Bosco Etoa Etoa. A globally convergent sequential linear programming algorithm for mathematical programs with linear complementarity constraints. *Journal of Information & Optimization Sciences*, 31, 09 2010.
- [15] Louis Faury, Marc Abeille, Clement Calauzenes, and Olivier Fercoq. Improved optimistic algorithms for logistic bandits. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 3052–3060. PMLR, 13–18 Jul 2020.
- [16] Louis Faury, Marc Abeille, Kwang-Sung Jun, and Clement Calauzenes. Jointly efficient and optimal algorithms for logistic bandits. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera, editors, *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pages 546–580. PMLR, 28–30 Mar 2022.
- [17] Sarah Filippi, Olivier Cappé, Aurélien Garivier, and Csaba Szepesvári. Parametric bandits: The generalized linear case. In *Proceedings of the 23rd International Conference on Neural Information Processing Systems - Volume 1*, NIPS’10, page 586–594, Red Hook, NY, USA, 2010. Curran Associates Inc.
- [18] Meire Fortunato, Mohammad Gheshlaghi Azar, Bilal Piot, Jacob Menick, Matteo Hessel, Ian Osband, Alex Graves, Volodymyr Mnih, Remi Munos, Demis Hassabis, Olivier Pietquin, Charles Blundell, and Shane Legg. Noisy networks for exploration. In *International Conference on Learning Representations*, 2018.
- [19] Dylan Foster and Alexander Rakhlin. Beyond UCB: Optimal and efficient contextual bandits with regression oracles. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 3199–3210. PMLR, 13–18 Jul 2020.
- [20] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [21] Nima Hamidi and Mohsen Bayati. On frequentist regret of linear thompson sampling. *arXiv preprint arXiv:2006.06790*, 2020.
- [22] Osama A. Hanna, Lin F. Yang, and Christina Fragouli. Efficient batched algorithm for contextual linear bandits with large action space via soft elimination. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, Red Hook, NY, USA, 2023. Curran Associates Inc.
- [23] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2015.
- [24] Ian Osband, Van Roy Benjamin, and Russo Daniel. (more) efficient reinforcement learning via posterior sampling. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, NIPS’13, page 3003–3011, Red Hook, NY, USA, 2013. Curran Associates Inc.
- [25] David Janz, Shuai Liu, Alex Ayoub, and Csaba Szepesvári. Exploration via linearly perturbed loss minimisation. In Sanjoy Dasgupta, Stephan Mandt, and Yingzhen Li, editors, *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pages 721–729. PMLR, 02–04 May 2024.

- [26] Kwang-Sung Jun, Aniruddha Bhargava, Robert Nowak, and Rebecca Willett. Scalable generalized linear bandits: Online computation and hashing. *Advances in Neural Information Processing Systems*, 30, 2017.
- [27] Kwang-Sung Jun, Lalit Jain, Blake Mason, and Houssam Nassif. Improved confidence bounds for the linear logistic model and applications to bandits. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 5148–5157. PMLR, 18–24 Jul 2021.
- [28] Sampath Kannan, Jamie Morgenstern, Aaron Roth, Bo Waggoner, and Zhiwei Steven Wu. A smoothed analysis of the greedy algorithm for the linear contextual bandit problem. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, NIPS’18, page 2231–2241, Red Hook, NY, USA, 2018. Curran Associates Inc.
- [29] Jung-Hun Kim, Se-Young Yun, Minchan Jeong, Junhyun Nam, Jinwoo Shin, and Richard Combes. Contextual linear bandits under noisy features: Towards bayesian oracles. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent, editors, *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pages 1624–1645. PMLR, 25–27 Apr 2023.
- [30] Alexander Kolesnikov, Alexey Dosovitskiy, Dirk Weissenborn, Georg Heigold, Jakob Uszkoreit, Lucas Beyer, Matthias Minderer, Mostafa Dehghani, Neil Houlsby, Sylvain Gelly, Thomas Unterthiner, and Xiaohua Zhai. An image is worth 16x16 words: Transformers for image recognition at scale. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- [31] Branislav Kveton, Csaba Szepesvári, Mohammad Ghavamzadeh, and Craig Boutilier. Perturbed-history exploration in stochastic linear bandits. In *Uncertainty in Artificial Intelligence*, pages 530–540. PMLR, 2020.
- [32] Branislav Kveton, Manzil Zaheer, Csaba Szepesvari, Lihong Li, Mohammad Ghavamzadeh, and Craig Boutilier. Randomized exploration in generalized linear bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 2066–2076. PMLR, 2020.
- [33] John Langford and Tong Zhang. The epoch-greedy algorithm for multi-armed bandits with side information. In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc., 2007.
- [34] Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- [35] Harin Lee and Min-hwan Oh. Improved regret of linear ensemble sampling. In *Advances in Neural Information Processing Systems*, volume 38, 2024.
- [36] Junghyun Lee, Se-Young Yun, and Kwang-Sung Jun. Improved regret bounds of (multinomial) logistic bandits via regret-to-confidence-set conversion. In Sanjoy Dasgupta, Stephan Mandt, and Yingzhen Li, editors, *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pages 4474–4482. PMLR, 02–04 May 2024.
- [37] Junghyun Lee, Se-Young Yun, and Kwang-Sung Jun. A unified confidence sequence for generalized linear models, with applications to bandits. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [38] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670, 2010.
- [39] Lihong Li, Yu Lu, and Dengyong Zhou. Provably optimal algorithms for generalized linear contextual bandits. In *International Conference on Machine Learning*, pages 2071–2080. PMLR, 2017.
- [40] Kolby Nottingham Markelle Kelly, Rachel Longjohn. The uci machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.

- [41] P. McCullagh and J.A. Nelder. *Generalized Linear Models*. Monographs on statistics and applied probability. Chapman and Hall/CRC, 2001.
- [42] Ian Osband, John Aslanides, and Albin Cassirer. Randomized prior functions for deep reinforcement learning. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, NIPS’18, page 8626–8638, Red Hook, NY, USA, 2018. Curran Associates Inc.
- [43] Ayush Sawarni, Nirjhar Das, Siddharth Barman, and Gaurav Sinha. Generalized linear bandits with limited adaptivity. In *Proceedings of the 38th International Conference on Neural Information Processing Systems*, NIPS ’24, Red Hook, NY, USA, 2025. Curran Associates Inc. ISBN 9798331314385.
- [44] Connor Shorten and Taghi Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of Big Data*, 6, 07 2019.
- [45] David Simchi-Levi and Yunzong Xu. Bypassing the monster: A faster and simpler optimal algorithm for contextual bandits under realizability. *Math. Oper. Res.*, 47(3):1904–1931, August 2022. ISSN 0364-765X.
- [46] Sharan Vaswani, Abbas Mehrabian, Audrey Durand, and Branislav Kveton. Old dog learns new tricks: Randomized ucb for bandit problems. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 1988–1998. PMLR, 26–28 Aug 2020.
- [47] R. Wolke and H. Schwetlick. Iteratively reweighted least squares: Algorithms, convergence analysis, and numerical comparisons. *SIAM Journal on Scientific and Statistical Computing*, 9(5):907–921, 1988.
- [48] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *ArXiv*, abs/1708.07747, 2017.
- [49] Yunbei Xu and Assaf J. Zeevi. Upper counterfactual confidence bounds: a new optimism principle for contextual bandits. *ArXiv*, abs/2007.07876, 2020.
- [50] Andrea Zanette, David Brandfonbrener, Emma Brunskill, Matteo Pirodda, and Alessandro Lazaric. Frequentist regret bounds for randomized least-squares value iteration. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 1954–1964. PMLR, 26–28 Aug 2020.
- [51] Weitong Zhang, Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural thompson sampling. In *International Conference on Learning Representation (ICLR)*, 2021.
- [52] Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural contextual bandits with UCB-based exploration. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 11492–11502. PMLR, 13–18 Jul 2020.

Appendix

Table of Contents

A Further related wokrs	14
B Properties of FP distributions	16
C Application of FP algorithm	17
C.1 Generalized version of GLM-FP	17
C.2 Application to the linear bandit problem	17
D Table of notations	18
E Proof of main theorem	19
E.1 Proof of Lemma 1	19
E.2 Proof of Lemma 3	19
E.3 Proof of Lemma 2	20
E.4 Proof of Theorem 1	21
F Proof supplement and Lemmas	29
F.1 Supporting Lemmas for main proof	29
F.2 Auxiliary bounding inequalities	31
G Limitations	31
H Experimental settings and additional results	32
H.1 Experimental details	32
H.2 Additional experiments	33

A Further related wokrs

Landscape of GLB Algorithms. Generalized linear bandit (GLB) algorithms can be broadly divided into *deterministic OFU-type* and *randomized exploration*-based methods. OFU approaches such as GLM-UCB [17] and Logistic-UCB [26] achieve the tight $\tilde{O}(d\sqrt{T})$ or $\tilde{O}(d\sqrt{T}/\kappa)$ regret bound, and refinements further improve confidence construction [27], achieving $\tilde{O}(\sqrt{dT \log K})$ in the finite- K arm setting. Abeille et al. [4] provided an instance-dependent analysis showing that the regret can benefit from the curvature constant κ_* , achieving $\tilde{O}(d\sqrt{\kappa_* T})$. Subsequent works [36, 37] relaxed the dependence on S or improved computational efficiency [16] while preserving the same order. Randomized methods such as Thompson Sampling (TS) and Perturbed History Exploration (PHE) [3, 6, 31, 32], as well as more recent algorithms like EVILL [25] and RandUCB [46], typically achieve superior empirical performance but incur an additional \sqrt{d} penalty in the worst case, yielding $\tilde{O}(d^{3/2}\sqrt{T})$ for infinite arms or $\tilde{O}(d\sqrt{T \log K}/\kappa)$ for finite arms. Our proposed feature perturbation (FP) departs from parameter- or reward-perturbation by randomizing the *input features*, thereby closing this gap: as summarized in Table A.1, FP is the first randomized algorithm for GLBs to provably achieve $\tilde{O}(d\sqrt{\kappa_* T} + d^2/\kappa)$ regret with no dependence on K , unifying the tight guarantees of OFU-type methods with the empirical robustness of randomized exploration.

Comparison to RandUCB Algorithm. Our linear variant LinFP (see C.2) and RandUCB [46] coincide in the linear bandit setting. Both sample randomized scores $\tilde{f}_t(x_i) \sim \mathcal{N}(x_{ti}^\top \hat{\theta}_t, \beta_t^2 \|x_{ti}\|_{V_t^{-1}}^2)$ and couple the arms identically, yielding matching $\tilde{O}(d\sqrt{T})$ regret bounds. The difference lies only

Table A.1: Representative GLB algorithms: regret bounds and source of stochasticity.

Type	Algorithm	Regret Upper Bound	Stochasticity
Deterministic	GLM-UCB [17]	$\tilde{\mathcal{O}}(d\sqrt{T}/\kappa)$	–
	Logistic-UCB-2 [26]	$\tilde{\mathcal{O}}(d\sqrt{T} + d^2/\kappa)$	–
	OFUGLB [37]	$\tilde{\mathcal{O}}(d\sqrt{\kappa_*T} + d^2/\kappa)$	–
Randomized	LinTS [3]	$\tilde{\mathcal{O}}(d^{3/2}\sqrt{T}/\kappa)$	Parameter (θ)
	GLM-TSL [32]	$\tilde{\mathcal{O}}(d^{3/2}\sqrt{T}/\kappa)$	Parameter (θ)
	GLM-FPL [32]	$\tilde{\mathcal{O}}(d^{3/2}\sqrt{T}/\kappa)$	Reward (r)
	RandUCB [46]	$\tilde{\mathcal{O}}(d\sqrt{T}/\kappa)$	Linear utility ($x^\top \theta$)
	GLM-FP (Ours)	$\tilde{\mathcal{O}}(d\sqrt{\kappa_*T} + d^2/\kappa)$	Feature vector (x)

in the *source* of randomness: FP perturbs the input features, whereas RandUCB perturbs the reward estimate itself. These distinct mechanisms collapse to the same Gaussian rule under linear models, though they have been analyzed through different theoretical perspectives. Equivalently, one may view RandUCB as a special instance of FP with an identical perturbation distribution, differing only in interpretation and analytical framework.

In generalized linear bandits (GLBs), however, the two algorithms diverge fundamentally. RandUCB extends its linear recipe by linearizing the link function, introducing a multiplicative κ^{-1} penalty and yielding a regret bound of $\tilde{\mathcal{O}}(d\sqrt{T}/\kappa)$. In contrast, GLM-FP perturbs the inputs directly using a curvature-aware Gram matrix that weights past features by $\hat{\mu}(x^\top \hat{\theta}_t)$, enabling both anti-concentration (for exploration) and concentration (for confidence). This yields a tighter regret of $\tilde{\mathcal{O}}(d\sqrt{\kappa_*T})$, linear in d and directly in κ_* . Conceptually, FP injects stochasticity *before* inference, so that each sampled reward $\hat{f}(\tilde{x})$ remains within the hypothesis class—preserving inductive bias and reflecting epistemic uncertainty. RandUCB, in contrast, adds randomness *after* inference, producing post hoc scores that may not correspond to any $f \in \mathcal{F}$. As a result, in expressive models FP remains aligned with the model structure, while RandUCB may misalign exploration incentives, leading to divergent empirical and theoretical behaviors.

Geometry and Scalability in Randomized Exploration. Recent studies have examined when randomized exploration can match the $\tilde{\mathcal{O}}(d\sqrt{T})$ guarantees of optimistic approaches. Abeille et al. [5] identified a class of geometric conditions—absorbing, strongly convex, and smooth action sets—under which Thompson Sampling (TS) achieves optimal dependence on d . While these conditions provide valuable insights into the role of geometry, they often fail to hold in high-dimensional or unstructured settings. Subsequent works [21, 32] further clarified that posterior variance inflation can inherently introduce the extra \sqrt{d} factor observed in randomized methods. In contrast, feature-level perturbation achieves similar statistical optimality under the standard boundedness assumption, bridging geometric optimality with more general feature-level regularity. From a computational standpoint, randomized exploration in large or continuous action spaces raises significant scalability challenges. Several strategies have been proposed to mitigate this issue, including lazy or delayed updates of the Gram matrix [1], two-stage candidate selection using approximate nearest neighbors, and optimization-oracle-based methods such as batched soft elimination [22]. These approaches highlight a trade-off between statistical tightness and computational efficiency: while algorithms like FP prioritize theoretical optimality in the online setting, batched or oracle-based techniques offer scalable alternatives for large-scale practical applications.

Connections to Feature Perturbation in Broader ML Feature perturbation is a common idea in other areas of machine learning, most notably in computer vision and natural language processing, where it is employed for robustness or regularization during training. Examples include data augmentation [44], adversarial training [20], or NoisyNets for exploration in deep reinforcement learning [18]. In these contexts, perturbations are introduced at training time to improve model generalization or robustness. In contrast, our FP algorithm introduces perturbations at *decision time* as a principled mechanism for exploration in online learning. This distinction highlights the novelty

of FP: rather than making a static predictor robust, we leverage feature perturbations dynamically to induce stochasticity in action selection, enabling efficient exploration. The same principle applies naturally when contextual information comes from high-dimensional embeddings, such as ResNet or ViT features for images [23, 30] or BERT embeddings for language [13], where FP can induce semantically meaningful exploration by perturbing compact representations. Thus, FP not only closes a theoretical gap in contextual bandits but also suggests a unifying exploration paradigm that resonates with broader trends in modern ML. Finally, this perspective also provides a bridge to reinforcement learning (RL), where perturbing the state–action feature representation can serve as an efficient alternative to parameter-space randomization used in posterior sampling or Noisy Networks [24, 42, 50]. Extending feature perturbation to structured settings such as Linear MDPs or value-function approximation is a promising direction for future work, potentially unifying exploration principles across bandit and reinforcement learning paradigms.

B Properties of FP distributions

Perturbation in Thompson Sampling The perturbation distribution utilized in the TS algorithm to bring randomness to the parameter, as described by Abeille and Lazaric [3] is as followed:

Definition B.1 (Definition 1. in Abeille and Lazaric [3]). \mathcal{D}^{TS} is a multivariate distribution on \mathbb{R}^d , absolutely continuous with respect to the Lebesgue measure, and satisfies the following properties:

1. (Anti-concentration) There exists a positive probability $p > 0$ such that for any unit vector $u \in \mathbb{R}^d$,

$$\mathbb{P}_{\zeta \sim \mathcal{D}^{TS}}(u^\top \zeta \geq 1) \geq p,$$

2. (Concentration) There exist positive constants c and c' such that for all $\delta \in (0, 1)$,

$$\mathbb{P}_{\zeta \sim \mathcal{D}^{TS}}(\|\zeta\| \leq \sqrt{cd \log(c'd/\delta)}) \geq 1 - \delta.$$

Below, we provide examples of distributions satisfying these *anti-concentration* and *concentration* properties, with the latter condition restated in Eq. (2) in Section 5.

Example 1: Gaussian distribution $\zeta \sim \mathcal{N}(0, I)$ The concentration property comes directly from Lemma F.6, as the inner product of a standard multivariate normal random variable ζ and an arbitrary unit vector u follows a standard normal distribution. In the same manner, for a unit vector u ,

$$\mathbb{P}_{\zeta \sim \mathcal{N}(0, I)}(u^\top \zeta \geq 1) = \mathbb{P}_{z \sim \mathcal{N}(0, 1)}(z \geq 1) = \frac{1}{2} \text{erfc}\left(\frac{1}{\sqrt{2}}\right) \geq \frac{1}{4\sqrt{e\pi}}. \quad (\text{B.1})$$

Thus, the standard Gaussian distribution satisfies the concentration property with $c = c' = 2$ and anti-concentration property with $p = \frac{1}{4\sqrt{e\pi}}$. Adjusting the scale of the covariance matrix, we can easily prove other variants satisfy the conditions.

Example 2: Uniform distribution $\zeta \sim \mathcal{U}_{B_d(0, \sqrt{d})}$ Let the random variable $\zeta = rv$, where $r = \|\zeta\| \in [0, \sqrt{d}]$ and $v = \zeta/\|\zeta\|$ is a unit vector. Then, $u^\top \zeta$ can be expressed as the product of two independent random variables, r and $u^\top v$ ($\sim \text{Beta}(\frac{1}{2}, \frac{d-1}{2})$), as $r \cdot (u^\top v)$. These two random variables follow the distributions:

$$f_r(r) = \frac{dr^{d-1}}{\sqrt{d}^d}, \quad r \in [0, \sqrt{d}], \quad \text{and} \quad f_{u^\top v}(x) = \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{1}{2})\Gamma(\frac{d-1}{2})}(1-x^2)^{\frac{d-3}{2}}, \quad x \in [-1, 1].$$

Based on these random variables, we can write:

$$f_{u^\top \zeta}(z) = \int_0^{\sqrt{d}} \int_{-1}^1 \delta(z - rx) f_r(r) f_{u^\top v}(x) dx dr.$$

Using Monte Carlo simulations, we observe that $u^\top \zeta$ has a lighter tail distribution compared a standard normal distribution. Accordingly, $c = c' = 2$ satisfies the concentration property. By proposition 9 and 10 in Abeille and Lazaric [3],

$$\mathbb{P}(u^\top \zeta \geq 1) = \frac{1}{2} I_{1-\frac{1}{d}}\left(\frac{d+1}{2}, \frac{1}{2}\right) \geq \frac{1}{16\sqrt{3\pi}},$$

where $I_x(a, b)$ is the incomplete regularized beta function. This suggests that the Uniform distribution satisfies the anti-concentration property with $p = \frac{1}{16\sqrt{3\pi}}$.

C Application of FP algorithm

C.1 Generalized version of GLM-FP

In Section 4, we introduced how FP can be applied to the contextual bandit settings in which the reward model extends beyond generalized linear models. We provide the general algorithmic framework below.

Algorithm C.1 Feature Perturbation in Bandit Problems

- 1: **Input:** Regularization parameter $\lambda > 0$, tuning parameter $\{c_t\}$
 - 2: **for** $t = 1, 2, \dots, T$ **do**
 - 3: Compute $\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{\tau=1}^{t-1} (f(x_{\tau, i_\tau}) - r_\tau)^2$ via a least squares oracle
 - 4: Sample $\tilde{x}_{ti} \sim \mathcal{D}(x_{ti}, \Sigma_{ti})$ for all i ▷ e.g., $\mathcal{D}(x_{ti}, \Sigma_{ti}) = \mathcal{N}(x_{ti}, \mathbf{I}/t)$
 - 5: Select arm $i_t = \arg \max_{i \in [|\mathcal{X}_t|]} \hat{f}(\tilde{x}_{ti})$ ▷ either i.i.d. or via shared perturbation
 - 6: Observe reward $r_t = f^*(x_{t, i_t}) + \xi_t$
 - 7: **end for**
-

C.2 Application to the linear bandit problem

While line 4 in Algorithm C.1 merely defines the sampling distribution for each arm, in practice—mirroring the design of GLM-FP—one may introduce a shared perturbing factor that is first sampled and then applied to all arms. This construction induces dependencies across the perturbed arms and can serve as the basis for the arm selection mechanism.

Algorithm C.2 LinFP: Feature Perturbation in Linear bandits

- 1: **Input:** Regularization parameter $\lambda > 0$, tuning parameter $\{c_t\}$
 - 2: **Initialize:** $V_1 \leftarrow \lambda \mathbf{I}$, $b_1 \leftarrow \mathbf{0}_d$
 - 3: **for** $t = 1, 2, \dots, T$ **do**
 - 4: Compute $\hat{\theta}_t = V_t^{-1} b_t$
 - 5: Sample $\zeta_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ▷ Shared perturbing factor
 - 6: Compute $\tilde{x}_{ti} = x_{ti} + c_t \cdot \frac{\|x_{ti}\|_{\hat{V}_t^{-1}}}{\|\hat{\theta}_t\|} \cdot \zeta_t$ for all i
 - 7: Choose $i_t = \arg \max_{i \in [|\mathcal{X}_t|]} \tilde{x}_{ti}^\top \hat{\theta}_t$ ▷ $x_{ti}^\top \tilde{\theta}_t \sim \mathcal{N}(x_{ti}^\top \hat{\theta}_t, c_t^2 \|x_{ti}\|_{V_t^{-1}}^2)$ for all i .
 - 8: Observe reward $r_t = x_{t, i_t}^\top \theta^* + \xi_t$
 - 9: Update $V_{t+1} = V_t + x_{t, i_t} x_{t, i_t}^\top$, $b_{t+1} = b_t + x_{t, i_t} r_t$
 - 10: **end for**
-

In Section 6, we compare LinFP, the linear variant of our approach, to LinTS [3, 6]. The primary algorithmic difference lies in lines 6–7 of Algorithm C.2. In our method, the shared random vector ζ_t is used to perturb each feature vector. In contrast, LinTS perturbs the model parameter as $\tilde{\theta}_t = \hat{\theta}_t + c_t \cdot V_t^{-1/2} \zeta_t$, and selects the arm maximizing $x_{ti}^\top \tilde{\theta}_t$. For LinTS, we have:

$$\begin{aligned} \mathbb{E}[x_{ti}^\top \tilde{\theta}_t] &= x_{ti}^\top \hat{\theta}_t + c_t \cdot x_{ti}^\top V_t^{-1/2} \mathbb{E}[\zeta_t] = x_{ti}^\top \hat{\theta}_t \\ \operatorname{Var}[x_{ti}^\top \tilde{\theta}_t] &= c_t^2 \cdot \operatorname{Var}[x_{ti}^\top V_t^{-1/2} \zeta_t] = c_t^2 \cdot x_{ti}^\top V_t^{-1} x_{ti} = c_t^2 \|x_{ti}\|_{V_t^{-1}}^2 \end{aligned}$$

While the marginal distribution of $x_{ti}^\top \tilde{\theta}_t$ under both methods is identical, the use of a shared perturbation ζ_t in LinFP induces algorithmic coupling across arms. This distinction is further discussed in Section 6.

D Table of notations

This section introduces additional notations and concepts essential for the analysis. For ease of reference, Table D.1 summarizes the key notations defined in this paper, along with newly introduced notations. Conventional concepts such as $d, T, \mathcal{A}, \mathcal{C}, \mathcal{X}, K$ or r are omitted here. The concepts will be reintroduced as needed in subsequent sections.

Table D.1: Notations and concepts for the analysis of the main theorem

Notation	Definition
M_μ	Self-concordance constant
\mathcal{L}_μ	Lipschitz constant of the link function μ
\mathcal{L}_t	Lipschitz constant of the negative log-likelihood function
$\Theta_t(\delta, \lambda)$	$1 - \delta$ probability ellipsoidal relaxed confidence set with regularization λ for the true parameter θ^*
$\beta_t(\delta)$	$\sqrt{4S^2\lambda + 2(1 + SM_\mu)(\log(1/\delta) + d\log(2e\mathcal{L}_t/d))} = \tilde{O}(\sqrt{d})$
$\gamma_t(\delta)$	$\beta_t(\delta/(4T))\sqrt{c\log(4c'T/\delta)}$ (c, c' : constant satisfying concentration property)
$\mathcal{E}_t(x)$	$\left\{ \tilde{x} \in \mathbb{R}^d \mid \ \tilde{x} - x, \hat{\theta}_t\ \leq \gamma_t(\delta/(4T))\ x\ _{\hat{H}_t^{-1}} \right\}$
κ_*	Average derivative of link function at the true optimal arm over T rounds
κ	Minimum reachable derivative of link function
Warm-up stage	
I_T	$\left\{ t \in [T] : \left(\ \sqrt{\dot{\mu}(x_t^\top \hat{\theta}_t)}x_t\ _{\hat{H}_t^{-1}} \geq 1 \right) \vee \left(\ x_t\ _{V_t^{-1}} \geq 1 \right) \right\}$
Taylor remainder term	
$\bar{\alpha}_t(x)$	$\int_0^1 (1-u)\dot{\mu}(x_t^\top \hat{\theta}_t + u(x^\top \hat{\theta}_t - x_t^\top \hat{\theta}_t))du$
$\bar{\alpha}_t(\theta, \nu)$	$\int_0^1 (1-u)\dot{\mu}(x_t^\top \theta + u(x_t^\top \nu - x_t^\top \theta))du$
Matrices	
V_t	$\sum_{\tau=1}^{t-1} x_\tau x_\tau^\top + \lambda \mathbf{I}$
\bar{V}_t	$\sum_{\tau=1}^{t-1} x_\tau x_\tau^\top + \lambda \mathbf{I}/\kappa$
\hat{H}_t	$\sum_{\tau=1}^{t-1} \dot{\mu}(x_\tau^\top \hat{\theta}_t)x_\tau x_\tau^\top + \lambda \mathbf{I}$
\bar{H}_t	$\sum_{\tau=1}^{t-1} \dot{\mu}(x_\tau^\top \bar{\theta}_\tau)x_\tau x_\tau^\top + \lambda \mathbf{I}$
$\tilde{H}_t(\theta, \nu)$	$\sum_{\tau=1}^{t-1} \bar{\alpha}_\tau(\theta, \nu)x_\tau x_\tau^\top$
Other notations	
R^*	$\max_{x \in \mathcal{X}} \mu(x^\top \theta^*) $
$\bar{\theta}_t$	$\operatorname{argmin}_{\theta \in \cup_{\tau \in [t, T]} \Theta_\tau(\delta, \lambda)} \dot{\mu}(x_t^\top \theta)$
$(\tau(t), \omega_t)$	$\operatorname{argmax}_{\tau \in [t, T], \theta \in \Theta_\tau(\delta, \lambda)} \left \mu(x_t^\top \theta) - \mu(x_t^\top \hat{\theta}_\tau) \right $
$(p, c_{\delta'})$	constants related to standard normal distribution ($p = 1/(4\sqrt{e\pi})$, $c_{\delta'} = \sqrt{2\log(2/\delta')}$)

E Proof of main theorem

The first step in our proof is to derive the high-probability confidence bound for the estimate $\hat{\theta}_t$. Using this bound, we ensure that the events defined in Definition 1 occur with high probability. Finally, we compute the regret bound of our algorithm under these events.

E.1 Proof of Lemma 1

We rederived Theorem 3.2 of Lee et al. [37] to obtain a tighter bound with improved dependence on the regularization parameter λ , such that the λ term no longer scales with M_μ .

By Theorem 3.1. in [37], with probability at least $1 - \delta$, for all $t \geq 1$, the following inequality holds:

$$L_t(\theta^*) - L_t(\hat{\theta}_t) \leq \log \frac{1}{\delta} + d \log \left(\frac{2eS\mathcal{L}_t}{d} \right) := \mathcal{W}_t(\delta)^2$$

Then we observe:

$$\begin{aligned} \int_0^1 (1-u) \nabla^2 L_t(\hat{\theta}_t + u(\theta^* - \hat{\theta}_t)) du &= \int_0^1 (1-u) \sum_{\tau=1}^{t-1} \dot{\mu}(x_\tau^\top (\hat{\theta}_t + u(\theta^* - \hat{\theta}_t))) x_\tau x_\tau^\top du \\ &= \sum_{\tau=1}^{t-1} \underbrace{\left(\int_0^1 (1-u) \dot{\mu}(x_\tau^\top (\hat{\theta}_t + u(\theta^* - \hat{\theta}_t))) du \right)}_{\bar{\alpha}_\tau(\hat{\theta}_t, \theta^*)} x_\tau x_\tau^\top = \tilde{H}_t(\hat{\theta}_t, \theta^*), \end{aligned}$$

where the second equality follows from Fubini's theorem, where the order of the integral and the summation can be switched. Using Taylor's theorem with an integral remainder (Proposition F.2), we can further deduce that, with probability $1 - \delta$:

$$\begin{aligned} \mathcal{W}_t(\delta)^2 &\geq L_t(\theta^*) - L_t(\hat{\theta}_t) \\ &= \langle \nabla L_t(\hat{\theta}_t), \theta^* - \hat{\theta}_t \rangle + (\theta^* - \hat{\theta}_t)^\top \left(\int_0^1 (1-u) \nabla^2 L_t(\hat{\theta}_t + u(\theta^* - \hat{\theta}_t)) du \right) (\theta^* - \hat{\theta}_t). \end{aligned}$$

As the optimality condition at $\hat{\theta}_t$ infers that $\langle \nabla L_t(\hat{\theta}_t), \theta^* - \hat{\theta}_t \rangle \geq 0$ and by equation (Proposition F.1), we have that

$$\mathcal{W}_t(\delta)^2 \geq \|\theta^* - \hat{\theta}_t\|_{\tilde{H}_t(\hat{\theta}_t, \theta^*)}^2 \geq \frac{1}{2 + 2SM_\mu} \|\theta^* - \hat{\theta}_t\|_{\hat{H}_t - \lambda \mathbf{I}}^2, \quad (\text{Proposition F.1})$$

where the last inequality holds, since

$$\tilde{H}_t(\hat{\theta}_t, \theta^*) \succeq \sum_{\tau=1}^{t-1} \left(\frac{\dot{\mu}(x_\tau^\top \hat{\theta}_t)}{2 + 2SM_\mu} \right) x_\tau x_\tau^\top = \frac{1}{2 + 2SM_\mu} (\hat{H}_t - \lambda \mathbf{I}).$$

Accordingly,

$$\|\theta^* - \hat{\theta}_t\|_{\hat{H}_t}^2 \leq \|\theta^* - \hat{\theta}_t\|_{\hat{H}_t - \lambda \mathbf{I}}^2 + \lambda \|\theta^* - \hat{\theta}_t\|^2 \leq 4S^2\lambda + 2(1 + SM_\mu)\mathcal{W}_t(\delta)^2 = \beta_t(\delta)^2,$$

and by Assumption 1, we finish a proof.

E.2 Proof of Lemma 3

Let the event \tilde{E}_t be defined as $\tilde{E}_t := \left\{ \mu(\tilde{x}_t^\top \hat{\theta}_t) \geq \mu(x_{t^*}^\top \theta^*) \right\}$. This event corresponds to the case where the chosen perturbed feature vector yields an optimistic expected reward relative to the true optimal arm at step t . To bound the anti-concentration probability as required in Lemma 3, we aim to

lower bound $\mathbb{P}_t(\ddot{E}_t \mid \hat{E}_t, \tilde{E}_t)$, conditioned on the events \hat{E}_t and \tilde{E}_t . For any $t \in [T]$, we have:

$$\begin{aligned}
\mathbb{P}_t(\ddot{E}_t \mid \hat{E}_t, \tilde{E}_t) &= \mathbb{P}_t\left(\mu(\tilde{x}_t^\top \hat{\theta}_t) \geq \mu(x_{t*}^\top \theta^*) \mid \hat{E}_t, \tilde{E}_t\right) \\
&= \mathbb{P}_t\left(\tilde{x}_t^\top \hat{\theta}_t \geq x_{t*}^\top \theta^* \mid \hat{E}_t, \tilde{E}_t\right) \quad (\because \mu \text{ is strictly increasing}) \\
&\geq \mathbb{P}_t\left(\tilde{x}_{t*}^\top \hat{\theta}_t - x_{t*}^\top \hat{\theta}_t \geq x_{t*}^\top \theta^* - x_{t*}^\top \hat{\theta}_t \mid \hat{E}_t, \tilde{E}_t\right) \quad (\because x_t = \operatorname{argmax}_{i \in [\mathcal{K}]} x_{ti}^\top \hat{\theta}_t) \\
&\geq \mathbb{P}_t\left(\langle \tilde{x}_{t*} - x_{t*}, \hat{\theta}_t \rangle \geq \left| \langle x_{t*}, \theta^* - \hat{\theta}_t \rangle \right| \mid \hat{E}_t, \tilde{E}_t\right) \\
&\geq \mathbb{P}_t\left(\left(\beta_t(\delta') \frac{\|x_{t*}\|_{\hat{H}_t^{-1}}}{\|\hat{\theta}_t\|} \zeta_t\right)^\top \hat{\theta}_t \geq \|x_{t*}\|_{\hat{H}_t^{-1}} \|\hat{\theta}_t - \theta^*\|_{\hat{H}_t} \mid \hat{E}_t, \tilde{E}_t\right) \\
&\geq \mathbb{P}\left(\beta_t(\delta') \|x_{t*}\|_{\hat{H}_t^{-1}} \cdot \langle \zeta_t, u_t \rangle \geq \beta_t(\delta') \|x_{t*}\|_{\hat{H}_t^{-1}} \mid \hat{E}_t, \tilde{E}_t\right) \\
&= \mathbb{P}(\langle \zeta_t, u_t \rangle \geq 1) \geq \frac{1}{4\sqrt{e\pi}} := p,
\end{aligned}$$

where the third inequality follows from the Cauchy-Schwarz inequality, and the fourth from the assumption that under the event \hat{E}_t , we have $\|\theta^* - \hat{\theta}_t\|_{\hat{H}_t} \leq \beta_t(\delta')$. The final inequality follows from the anti-concentration property of the standard normal distribution, as detailed in (B.1). For simplicity, we henceforth fix $p := 1/(4\sqrt{e\pi})$ as the corresponding lower bound on this probability.

E.3 Proof of Lemma 2

We now proceed to establish how the probability of each event defined in Definition 1 can be ensured using the confidence bound $\beta_t(\delta)$ derived in Appendix E.1. Each event is analyzed and bounded individually, and the results are then combined to complete the proof of the lemma.

Bounding \hat{E}_T Let $\delta' = \delta/(4T)$. By the choice of $\beta_t(\delta)$ in Lemma 1, we have that

$$\begin{aligned}
\forall 1 \leq t \leq T, \quad \mathbb{P}\left(\|\hat{\theta}_t - \theta^*\|_{\hat{H}_t} \leq \beta_t(\delta')\right) &\geq 1 - \delta' \\
\text{from union bound, } \mathbb{P}\left(\bigcap_{t=1}^T \left\{\|\hat{\theta}_t - \theta^*\|_{\hat{H}_t} \leq \beta_t(\delta')\right\}\right) &\geq 1 - \sum_{t=1}^T \mathbb{P}\left(\|\hat{\theta}_t - \theta^*\|_{\hat{H}_t} \geq \beta_t(\delta')\right) \\
\implies \mathbb{P}\left(\bigcap_{t=1}^T \left\{\|\hat{\theta}_t - \theta^*\|_{\hat{H}_t} \leq \beta_t(\delta')\right\}\right) &\geq 1 - \sum_{t=1}^T \delta' \\
\implies \mathbb{P}(\hat{E}_T) &\geq 1 - T\delta' = 1 - \frac{\delta}{4}.
\end{aligned}$$

Bounding \tilde{E}_T The expression for the perturbed feature vector \tilde{x}_{ti} is given as the expression $\tilde{x}_{ti} = x_{ti} + \beta_t(\delta') \frac{\|x_{ti}\|_{\hat{H}_t^{-1}}}{\|\hat{\theta}_t\|} \zeta_t$ with the choice of $c_t = \beta_t(\delta')$, where ζ_t is drawn *i.i.d.* from $\mathcal{N}(\mathbf{0}, \mathbf{I})$. Note that the constants c and c' indicated in (2) are both 2, as proved in (B.1), from now on for simplicity, we let $c_{\delta'} := \sqrt{2 \log(2/\delta')}$. Since all arms are coupled⁵ with same ζ_t , we can write

$$\forall 1 \leq t \leq T, \quad \mathbb{P}\left(\forall x_{ti} \in \mathcal{X}_t; \quad \left|\langle \tilde{x}_{ti} - x_{ti}, \hat{\theta}_t \rangle\right| \leq \gamma_t(\delta') \|x_{ti}\|_{\hat{H}_t^{-1}}\right) \quad (\text{E.1})$$

$$= \mathbb{P}\left(\forall x_{ti} \in \mathcal{X}_t; \quad \beta_t(\delta') \|x_{ti}\|_{\hat{H}_t^{-1}} \left|\left\langle \zeta_t, \frac{\hat{\theta}_t}{\|\hat{\theta}_t\|} \right\rangle\right| \leq c_{\delta'} \cdot \beta_t(\delta') \|x_{ti}\|_{\hat{H}_t^{-1}}\right) \quad (\text{E.2})$$

$$= \mathbb{P}(|\langle \zeta_t, u_t \rangle| \leq c_{\delta'}) \geq 1 - \delta', \quad (\text{E.3})$$

⁵Uncoupled sampling means ζ_{ti} 's are sampled for each arm respectively and it results in extra $\log K$ term in the regret bound because of the union bound over the number of arms at each step.

where u_t is a unit vector. The first equality holds by the definition of $c_{\delta'}$ and Definition 1, and the inequality follows from the concentration property. The cancellation in the second equality plays a critical role in removing arm-wise dependence in GLB setting. A union bound over T rounds yields

$$\mathbb{P}(\tilde{E}_T) \geq 1 - T\delta' = 1 - \frac{\delta}{4}.$$

Finally, applying the union bound across the events \hat{E}_T and \tilde{E}_T , we have that

$$\mathbb{P}(\hat{E}_T \cap \tilde{E}_T) \geq 1 - \frac{\delta}{2}.$$

Remark. To guarantee the same probability level $1 - \delta'$ as in equations (E.1)–(E.3), which bound the randomness arising from perturbations to the feature vectors, Thompson Sampling requires the confidence parameters β and γ to be inflated by an additional factor of \sqrt{d} . This inflation arises due to the right-hand side of the concentration bound in Definition B.1, which scales with \sqrt{d} . Such adjustment is necessary to control the deviation in the perturbed estimated expected reward, which takes the form $x^\top(\tilde{\theta} - \hat{\theta})$. This is consistent with the reasoning discussed in Section 6.

E.4 Proof of Theorem 1

In this section, we establish the regret guarantee for our algorithm. Given the complexity of the analysis, we divide the proof into multiple steps. Supporting lemmas and their proofs are deferred to Appendix F.

Step 1 (Warm-up) We begin by partitioning the T rounds into a “warm-up” stage and the primary stage. The set of time steps corresponding to the primary stage is defined as:

$$I_T := \left\{ t \in [T] : \left(\left\| \sqrt{\dot{\mu}(x_t^\top \bar{\theta}_t)} x_t \right\|_{\bar{H}_t^{-1}} \leq 1 \right) \wedge \left(\|x_t\|_{\bar{V}_t^{-1}} \leq 1 \right) \right\},$$

where \bar{H}_t , \bar{V}_t , and $\bar{\theta}_t$ are given by:

$$\bar{H}_t := \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} \dot{\mu}(x_\tau^\top \bar{\theta}_\tau) x_\tau x_\tau^\top, \quad \bar{V}_t := \lambda \mathbf{I} / \kappa + \sum_{\tau=1}^{t-1} x_\tau x_\tau^\top, \quad \bar{\theta}_t := \underset{\theta \in \cup_{\tau \in [t, T]} \Theta_\tau(\delta, \lambda)}{\operatorname{argmin}} \dot{\mu}(x_t^\top \theta).$$

The introduction of \bar{H}_t is crucial because $\hat{H}_t = \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} \dot{\mu}(x_\tau^\top \hat{\theta}_\tau) x_\tau x_\tau^\top$ depends on t , which prevents direct application of the Elliptical Potential Lemma (EPL; Lemma F.1), as discussed in Section 5. To address this, we leverage \bar{H}_t , which incorporates the minimum derivative of μ within future confidence sets ($\bar{\theta}_\tau$), ensuring it serves as a smaller Gram matrix suitable for bounding the regret. Similarly, \bar{V}_t is introduced to directly apply EPL.

Next, we bound each weighted 2-norm using the Elliptical Potential Count Lemma (EPCL; Lemma F.2), which guarantees that the regret incurred during the warm-up phase remains manageable. Consequently, the cumulative regret over T rounds is decomposed as follows:

$$\begin{aligned} R(T) &= \underbrace{\sum_{t \in I_T} \{\mu(x_{t^*}^\top \theta^*) - \mu(x_t^\top \theta^*)\}}_{\operatorname{Reg}(T)} + \underbrace{\sum_{t \notin I_T} \{\mu(x_{t^*}^\top \theta^*) - \mu(x_t^\top \theta^*)\}}_{\text{warm-up regret}} \\ &\leq \operatorname{Reg}(T) + 2R^* \left(\sum_{t=1}^T \mathbb{1} \left\{ \left\| \sqrt{\dot{\mu}(x_t^\top \bar{\theta}_t)} x_t \right\|_{\bar{H}_t^{-1}} \geq 1 \right\} + \sum_{t=1}^T \mathbb{1} \left\{ \|x_t\|_{\bar{V}_t^{-1}} \geq 1 \right\} \right) \\ &\leq \operatorname{Reg}(T) + \frac{4dR^*}{\log 2} \left\{ \log \left(1 + \frac{\mathcal{L}_\mu}{\lambda \log 2} \right) + \log \left(1 + \frac{\kappa}{\lambda \log 2} \right) \right\}, \end{aligned}$$

where $R^* := \max_{x \in \mathcal{X}} |\mu(x^\top \theta^*)|$ is the maximum expected reward achievable under the underlying model. The first and the second inequalities hold from the definition of I_T and by EPCL (Lemma F.2), respectively. Note that the warm-up regret is $\tilde{\mathcal{O}}(d)$, which is independent of T .

Step 2-1 (Decomposition) We decompose the cumulative regret for the primary stage into three components:

$$\text{Reg}(T) = \sum_{t \in I_T} \left(\underbrace{\left\{ \mu(x_{t*}^\top \theta^*) - \mu(\tilde{x}_t^\top \hat{\theta}_t) \right\}}_{A_t} + \underbrace{\left\{ \mu(\tilde{x}_t^\top \hat{\theta}_t) - \mu(x_t^\top \hat{\theta}_t) \right\}}_{B_t} + \underbrace{\left\{ \mu(x_t^\top \hat{\theta}_t) - \mu(x_t^\top \theta^*) \right\}}_{C_t} \right).$$

Here, A_t and B_t relate to the perturbations' effect on the estimated reward, while C_t concerns the closeness of $\hat{\theta}_t$ to θ^* . We will bound each term under the events \hat{E}_t and \tilde{E}_t .

Bounding C_t Bounding C_t is straightforward. Using the confidence set $\Theta_t(\delta, \lambda)$, abbreviated as Θ_t , we define $(\tau(t), \omega_t)$ as the pair maximizing the confidence width computed on the selected action at round t , x_t , after round t : $\arg\max_{\tau \in [t, T], \theta \in \Theta_\tau} \left| \mu(x_t^\top \theta) - \mu(x_t^\top \hat{\theta}_\tau) \right|$. Under the event \hat{E}_t , we know that $\theta^* \in \Theta_\tau$ for all $t \leq \tau \leq T$. Thus,

$$C_t \mathbb{1}\{\hat{E}_t \cap \tilde{E}_t\} = \left(\mu(x_t^\top \hat{\theta}_t) - \mu(x_t^\top \theta^*) \right) \mathbb{1}\{\hat{E}_t \cap \tilde{E}_t\} \leq \left| \mu(x_t^\top \omega_t) - \mu(x_t^\top \hat{\theta}_{\tau(t)}) \right| \mathbb{1}\{\hat{E}_t \cap \tilde{E}_t\}.$$

Bounding B_t Let $\tilde{x}_t^* := \arg\max_{x \in \mathcal{E}_t(x_t)} \left| \mu(x^\top \hat{\theta}_t) - \mu(x_t^\top \hat{\theta}_t) \right|$. Under the event \tilde{E}_t , $\tilde{x}_t \in \mathcal{E}_t(x_t)$ holds and we can write:

$$B_t \mathbb{1}\{\hat{E}_t \cap \tilde{E}_t\} = \left(\mu(\tilde{x}_t^\top \hat{\theta}_t) - \mu(x_t^\top \hat{\theta}_t) \right) \mathbb{1}\{\hat{E}_t \cap \tilde{E}_t\} \leq \left| \mu(\tilde{x}_t^{\top*} \hat{\theta}_t) - \mu(x_t^\top \hat{\theta}_t) \right| \mathbb{1}\{\hat{E}_t \cap \tilde{E}_t\}.$$

Before proceeding, note that for $\tilde{x}_t \in \mathcal{E}_t(x_t)$, we can derive an upper bound using Taylor's theorem with an integral remainder (Proposition F.2). Define $\bar{\alpha}_t(x) = \int_0^1 (1-u) \dot{\mu}(x_t^\top \hat{\theta}_t + u(x^\top \hat{\theta}_t - x_t^\top \hat{\theta}_t)) du$, which accounts for higher-order terms based on the estimated parameter $\hat{\theta}_t$ and feature vectors x and x_t . The difference $\left| \mu(\tilde{x}_t^\top \hat{\theta}_t) - \mu(x_t^\top \hat{\theta}_t) \right|$ then can be bounded as:

$$\begin{aligned} \left| \mu(\tilde{x}_t^\top \hat{\theta}_t) - \mu(x_t^\top \hat{\theta}_t) \right| &= \left| \dot{\mu}(x_t^\top \hat{\theta}_t) \langle \tilde{x}_t - x_t, \hat{\theta}_t \rangle + \int_{x_t^\top \hat{\theta}_t}^{\tilde{x}_t^\top \hat{\theta}_t} (\mu(\tilde{x}_t^\top \hat{\theta}_t) - z) \ddot{\mu}(z) dz \right| \\ &\leq \dot{\mu}(x_t^\top \hat{\theta}_t) \left| \langle \tilde{x}_t - x_t, \hat{\theta}_t \rangle \right| + \langle \tilde{x}_t - x_t, \hat{\theta}_t \rangle^2 \int_0^1 (1-u) \left| \ddot{\mu}(x_t^\top \hat{\theta}_t + u(\tilde{x}_t^\top \hat{\theta}_t - x_t^\top \hat{\theta}_t)) \right| du \\ &\leq \dot{\mu}(x_t^\top \hat{\theta}_t) \left| \langle \tilde{x}_t - x_t, \hat{\theta}_t \rangle \right| + M_\mu \langle \tilde{x}_t - x_t, \hat{\theta}_t \rangle^2 \underbrace{\int_0^1 (1-u) \left| \ddot{\mu}(x_t^\top \hat{\theta}_t + u(\tilde{x}_t^\top \hat{\theta}_t - x_t^\top \hat{\theta}_t)) \right| du}_{=\bar{\alpha}_t(\tilde{x}_t)} \\ &\leq \dot{\mu}(x_t^\top \hat{\theta}_t) \left| \langle \tilde{x}_t - x_t, \hat{\theta}_t \rangle \right| + M_\mu \bar{\alpha}_t(\tilde{x}_t) \langle \tilde{x}_t - x_t, \hat{\theta}_t \rangle^2 \\ &\leq \dot{\mu}(x_t^\top \hat{\theta}_t) \gamma_t(\delta') \|x_t\|_{\hat{H}_t}^{-1} + M_\mu \bar{\alpha}_t(\tilde{x}_t) \gamma_t(\delta')^2 \|x_t\|_{\hat{H}_t}^{-1}, \end{aligned}$$

where the second and the last inequalities hold from Assumption 2 and the definition of $\mathcal{E}_t(x_t)$ in Definition 1. This bound captures both the linear and higher-order contributions to the regret from perturbations in the feature vectors.

Bounding A_t With $\ddot{E}_t := \left\{ \mu(\tilde{x}_t^\top \hat{\theta}_t) \geq \mu(x_{t*}^\top \theta^*) \right\}$ and \tilde{x}_t^* defined above, we write:

$$\begin{aligned}
A_t \mathbb{1}\{\hat{E}_t \cap \tilde{E}_t\} &= \left(\mu(x_{t*}^\top \theta^*) - \mu(\tilde{x}_t^\top \hat{\theta}_t) \right) \mathbb{1}\{\hat{E}_t \cap \tilde{E}_t\} \\
&\leq \left(\mu(x_{t*}^\top \theta^*) - \inf_{\tilde{x}_t \in \mathcal{E}_t(x_t)} \mu(\tilde{x}_t^\top \hat{\theta}_t) \right) \mathbb{1}\{\hat{E}_t \cap \tilde{E}_t\} \\
&\leq \mathbb{E}_t \left[\left(\mu(\tilde{x}_t^\top \hat{\theta}_t) - \inf_{\tilde{x}_t \in \mathcal{E}_t(x_t)} \mu(\tilde{x}_t^\top \hat{\theta}_t) \right) \mathbb{1}\{\hat{E}_t \cap \tilde{E}_t\} \middle| \ddot{E}_t \right] \\
&= \mathbb{E}_t \left[\left(\mu(\tilde{x}_t^\top \hat{\theta}_t) - \mu(x_t^\top \hat{\theta}_t) \right) + \left(\mu(x_t^\top \hat{\theta}_t) - \inf_{\tilde{x}_t \in \mathcal{E}_t(x_t)} \mu(\tilde{x}_t^\top \hat{\theta}_t) \right) \middle| \hat{E}_t, \tilde{E}_t, \ddot{E}_t \right] \mathbb{P}(\hat{E}_t \cap \tilde{E}_t) \\
&\leq 2\mathbb{E}_t \left[\left(\sup_{\tilde{x}_t \in \mathcal{E}_t(x_t)} \left| \mu(\tilde{x}_t^\top \hat{\theta}_t) - \mu(x_t^\top \hat{\theta}_t) \right| \right) \middle| \hat{E}_t, \tilde{E}_t, \ddot{E}_t \right] \mathbb{P}(\hat{E}_t \cap \tilde{E}_t) \\
&\leq \frac{2}{p} \mathbb{E}_t \left[\left| \mu(\tilde{x}_t^{*\top} \hat{\theta}_t) - \mu(x_t^\top \hat{\theta}_t) \right| \mathbb{1}\{\hat{E}_t \cap \tilde{E}_t\} \right].
\end{aligned}$$

We justify the second inequality under the specified event, and the final inequality follows from the following reasoning: we use the bound $C \leq \mathbb{E}[Z \mid Z \geq C]$, and compensate for introducing the conditional expectation by incorporating the inverse of the probability of the conditioning event. Specifically, define $C := (\mu(x_{t*}^\top \theta^*) - \inf_x \mu(x^\top \hat{\theta}_t)) \cdot \mathbb{1}\{\hat{E}_t \cap \tilde{E}_t\}$ and $Z := (\mu(\tilde{x}_t^\top \hat{\theta}_t) - \inf_x \mu(x^\top \hat{\theta}_t)) \cdot \mathbb{1}\{\hat{E}_t \cap \tilde{E}_t\}$. Then the second inequality holds. To compensate for conditioning on the favorable event, we use the following logic:

$$\begin{aligned}
\mathbb{E}_t \left\{ \cdot \middle| \hat{E}_t, \tilde{E}_t \right\} &\geq \mathbb{E}_t \left\{ \cdot \middle| \hat{E}_t, \tilde{E}_t, \ddot{E}_t \right\} \mathbb{P}_t(\ddot{E}_t \mid \hat{E}_t, \tilde{E}_t) \\
&\geq \mathbb{E}_t \left\{ \cdot \middle| \hat{E}_t, \tilde{E}_t, \ddot{E}_t \right\} \cdot p.
\end{aligned} \tag{Lemma 3}$$

The upper bound for this term is similar to the previous part (B_t), differing only by a constant and the inclusion of the expectation over the filtration. However, since our goal is to bound the cumulative sum over T rounds rather than the expectation itself, directly handling the expectation complicates the application of the Elliptical Potential Lemma (EPL). To address this, we eliminate the expectation at the cost of introducing a concentration error, which we control using Azuma-Hoeffding's inequality.

Step 2-2 (Azuma-Hoeffding's inequality) Unless otherwise specified, we now analyze the regret bound under the assumption that events \hat{E}_T and \tilde{E}_T hold.

$$\begin{aligned}
\frac{p}{2} \sum_{t \in I_T} A_t &\leq \sum_{t \in I_T} \mathbb{E}_t \left[\left| \mu(\tilde{x}_t^{*\top} \hat{\theta}_t) - \mu(x_t^\top \hat{\theta}_t) \right| \right] \\
&\leq \gamma_T(\delta') \sum_{t \in I_T} \underbrace{\left(\dot{\mu}(x_t^\top \hat{\theta}_t) \|x_t\|_{\hat{H}_t^{-1}} + \gamma_T(\delta') \sum_{t \in I_T} \left(\mathbb{E}_t \left[\dot{\mu}(x_t^\top \hat{\theta}_t) \|x_t\|_{\hat{H}_t^{-1}} \right] - \dot{\mu}(x_t^\top \hat{\theta}_t) \|x_t\|_{\hat{H}_t^{-1}} \right) \right)}_{R_1} \\
&\quad + M_\mu \gamma_T(\delta')^2 \sum_{t \in I_T} \bar{\alpha}_t(\tilde{x}_t) \|x_t\|_{\hat{H}_t^{-1}}^2 + M_\mu \gamma_T(\delta')^2 \sum_{t \in I_T} \underbrace{\left(\mathbb{E}_t \left[\bar{\alpha}_t(\tilde{x}_t) \|x_t\|_{\hat{H}_t^{-1}}^2 \right] - \bar{\alpha}_t(\tilde{x}_t) \|x_t\|_{\hat{H}_t^{-1}}^2 \right)}_{R_2}
\end{aligned}$$

Note that R_1 and R_2 are constructed as martingales. Since the norm of each feature vector satisfies $\|x_t\| \leq 1$, and given $\hat{H}_t^{-1} \preceq \hat{H}_0^{-1} = \mathbf{I}/\lambda$ and $\dot{\mu}(x_t^\top \hat{\theta}_t) \leq \mathcal{L}_\mu$, the following bounds hold:

$$0 \leq \dot{\mu}(x_t^\top \hat{\theta}_t) \|x_t\|_{\hat{H}_t^{-1}} \leq \mathcal{L}_\mu \sqrt{x_t^\top \hat{H}_t^{-1} x_t} \leq \mathcal{L}_\mu \sqrt{\frac{1}{\lambda} \|x_t\|^2} \leq \frac{\mathcal{L}_\mu}{\sqrt{\lambda}}.$$

This provides an upper bound for each instantaneous element of R_1 as $\mathcal{L}_\mu/\sqrt{\lambda}$. Applying Azuma-Hoeffding's inequality (Proposition F.3), with probability at least $1 - \delta/4$, we obtain:

$$R_1 \leq \sqrt{\frac{2T\mathcal{L}_\mu^2}{\lambda}} \log \frac{8}{\delta}.$$

Due to the convexity of Θ_t , the term $\dot{\mu}$ in $\bar{\alpha}_t(x_t)$ is bounded by \mathcal{L}_μ . Consequently, we obtain:

$$\bar{\alpha}_t(x) := \int_0^1 (1-u) \dot{\mu} \left(x_t^\top \hat{\theta}_t + u(x^\top \hat{\theta}_t - x_t^\top \hat{\theta}_t) \right) du \leq \int_0^1 (1-u) \mathcal{L}_\mu du = \mathcal{L}_\mu/2 \quad (\text{E.4})$$

Similarly, we can show that

$$0 \leq \bar{\alpha}_t(\tilde{x}_t) \|x_t\|_{\hat{H}_t^{-1}}^2 \leq \frac{\mathcal{L}_\mu}{2} \left(x_t^\top \hat{H}_t^{-1} x_t \right) \leq \frac{\mathcal{L}_\mu}{2\lambda} \|x_t\|^2 \leq \frac{\mathcal{L}_\mu}{2\lambda}.$$

This provides an upper bound for each instantaneous element of R_2 as $\mathcal{L}_\mu/(2\lambda)$. Applying Proposition F.3, with probability at least $1 - \delta/4$, we have:

$$R_2 \leq \sqrt{\frac{2T\mathcal{L}_\mu^2}{4\lambda^2} \log \frac{8}{\delta}}.$$

By applying a union bound, we conclude that with probability at least $1 - \delta/2$, both R_1 and R_2 are bounded. Therefore, the regret term $\sum_{t \in I_T} A_t$ can be bounded by:

$$\sum_{t \in I_T} A_t \leq \frac{2}{p} \left(\gamma_T(\delta') \sum_{t \in I_T} \dot{\mu}(x_t^\top \hat{\theta}_t) \|x_t\|_{\hat{H}_t^{-1}} + M_\mu \gamma_T(\delta')^2 \sum_{t \in I_T} \bar{\alpha}_t(\tilde{x}_t) \|x_t\|_{\hat{H}_t^{-1}}^2 \right) + \varepsilon,$$

where ε is a function of T , δ , and λ , defined as:

$$\varepsilon = \varepsilon(T, \delta, \lambda) := \frac{2}{p} \left(\gamma_T(\delta') \sqrt{\frac{2T\mathcal{L}_\mu^2}{\lambda} \log \frac{8}{\delta}} + M_\mu \gamma_T(\delta')^2 \sqrt{\frac{2T\mathcal{L}_\mu^2}{4\lambda^2} \log \frac{8}{\delta}} \right) = \tilde{\mathcal{O}} \left(d \sqrt{\frac{T}{\lambda^2}} \right).$$

In summary, combining this result with Lemma 3-2 and the bounds for B_t and C_t , the total regret $\text{Reg}(T)$ can be bounded with probability at least $1 - \delta$ as:

$$\begin{aligned} \text{Reg}(T) &\leq \left(\frac{2}{p} + 1 \right) \underbrace{\sum_{t \in I_T} \left(\gamma_T(\delta') \dot{\mu}(x_t^\top \hat{\theta}_t) \|x_t\|_{\hat{H}_t^{-1}} + M_\mu \gamma_T(\delta')^2 \bar{\alpha}_t(\tilde{x}_t) \|x_t\|_{\hat{H}_t^{-1}}^2 \right)}_{\text{Reg}_{\text{FP}}} \\ &\quad + \underbrace{\sum_{t \in I_T} \left| \mu(x_t^\top \omega_t) - \mu(x_t^\top \hat{\theta}_{\tau(t)}) \right|}_{\text{Reg}_{\text{EST}}} + \varepsilon \end{aligned}$$

Here, Reg_{FP} captures the regret arising from perturbing the feature vectors, while Reg_{EST} accounts for the estimation error of $\hat{\theta}_t$ compared to the true parameter θ^* .

Step 3 (Bounding Reg_{FP}) The presence of \hat{H}_t^{-1} in the weighted norm makes it challenging to directly apply the Elliptical Potential Lemma (EPL; Lemma F.1). To address this, we introduce \bar{H}_t , which allows us to leverage EPL by splitting $\dot{\mu}(x_t^\top \hat{\theta}_t)$ into two components: a leading term based on $\bar{\theta}_t$ (used in defining \bar{H}_t) and a transient term that accounts for deviations from $\hat{\theta}_t$.

$$\begin{aligned} \text{Reg}_{\text{FP}} &= \gamma_T(\delta') \sum_{t \in I_T} \dot{\mu}(x_t^\top \hat{\theta}_t) \|x_t\|_{\hat{H}_t^{-1}} + M_\mu \gamma_T(\delta')^2 \sum_{t \in I_T} \bar{\alpha}_t(\tilde{x}_t) \|x_t\|_{\hat{H}_t^{-1}}^2 \\ &\leq \gamma_T(\delta') \sum_{t \in I_T} \left(\dot{\mu}(x_t^\top \bar{\theta}_t) + \left| \dot{\mu}(x_t^\top \bar{\theta}_t) - \dot{\mu}(x_t^\top \hat{\theta}_t) \right| \right) \|x_t\|_{\hat{H}_t^{-1}} + M_\mu \gamma_T(\delta')^2 \sum_{t \in I_T} \bar{\alpha}_t(\tilde{x}_t) \|x_t\|_{\hat{H}_t^{-1}}^2 \\ &\leq \underbrace{\sum_{t \in I_T} \dot{\mu}(x_t^\top \bar{\theta}_t) \gamma_T(\delta') \|x_t\|_{\hat{H}_t^{-1}}}_{D_t^{\text{FP}}} + \underbrace{\sum_{t \in I_T} \left| \dot{\mu}(x_t^\top \bar{\theta}_t) - \dot{\mu}(x_t^\top \hat{\theta}_t) \right| \gamma_T(\delta') \|x_t\|_{\hat{H}_t^{-1}}}_{E_t^{\text{FP}}} + \underbrace{\sum_{t \in I_T} M_\mu \bar{\alpha}_t(\tilde{x}_t) \gamma_T(\delta')^2 \|x_t\|_{\hat{H}_t^{-1}}^2}_{F_t^{\text{FP}}}. \end{aligned}$$

Bounding $\sum_t D_t^{\text{FP}}$ We define $\bar{H}_t := \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} \dot{\mu}(x_\tau^\top \bar{\theta}_\tau) x_\tau x_\tau^\top$. Then for all $\tau \leq t$, as the equation $\dot{\mu}(x_\tau^\top \bar{\theta}_\tau) \leq \dot{\mu}(x_\tau^\top \hat{\theta}_t)$ holds, we write:

$$\hat{H}_t = \lambda \mathbf{I} + \nabla^2 \mathcal{L}_t(\hat{\theta}_t) = \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} \dot{\mu}(x_\tau^\top \hat{\theta}_t) x_\tau x_\tau^\top \succeq \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} \dot{\mu}(x_\tau^\top \bar{\theta}_\tau) x_\tau x_\tau^\top \succeq \bar{H}_t. \quad (\text{E.5})$$

Then, we can bound $\sum_t D_t^{\text{FP}}$ as following:

$$\begin{aligned}
\sum_{t \in I_T} D_t^{\text{FP}} &\leq \gamma_T(\delta') \sum_{t \in I_T} \dot{\mu}(x_t^\top \bar{\theta}_t) \|x_t\|_{\hat{H}_t^{-1}} \\
&\leq \gamma_T(\delta') \sqrt{\sum_{t \in I_T} \dot{\mu}(x_t^\top \bar{\theta}_t)} \sqrt{\sum_{t \in I_T} \dot{\mu}(x_t^\top \bar{\theta}_t) \|x_t\|_{\hat{H}_t^{-1}}^2} \quad (\text{Cauchy-Schwartz inequality}) \\
&\leq \gamma_T(\delta') \sqrt{\sum_{t \in I_T} \dot{\mu}(x_t^\top \bar{\theta}_t)} \sqrt{\sum_{t \in I_T} \dot{\mu}(x_t^\top \bar{\theta}_t) \|x_t\|_{\hat{H}_t^{-1}}^2} \quad (\text{By (E.5)}) \\
&\leq \gamma_T(\delta') \sqrt{\sum_{t \in I_T} \dot{\mu}(x_t^\top \bar{\theta}_t)} \sqrt{\sum_{t=1}^T \min \left\{ 1, \sqrt{\dot{\mu}(x_t^\top \bar{\theta}_t) x_t^\top x_t} \|x_t\|_{\hat{H}_t^{-1}}^2 \right\}} \quad (\text{Definition of } I_T)
\end{aligned}$$

Applying Lemma F.1, the second square-root term is bounded by

$$\sqrt{\sum_{t=1}^T \min \left\{ 1, \sqrt{\dot{\mu}(x_t^\top \bar{\theta}_t) x_t^\top x_t} \|x_t\|_{\hat{H}_t^{-1}}^2 \right\}} \leq \sqrt{2d \log \left(1 + \frac{\mathcal{L}_\mu T}{d\lambda} \right)}.$$

To handle the first term in the square root, we decompose it as:

$$\sqrt{\sum_{t \in I_T} \dot{\mu}(x_t^\top \bar{\theta}_t)} \leq \sqrt{\sum_{t=1}^T \dot{\mu}(x_{t*}^\top \theta^*) + \sum_{t \in I_T} \{\dot{\mu}(x_t^\top \bar{\theta}_t) - \dot{\mu}(x_{t*}^\top \theta^*)\}} = \sqrt{\kappa_* T + \sum_{t \in I_T} \{\dot{\mu}(x_t^\top \bar{\theta}_t) - \dot{\mu}(x_{t*}^\top \theta^*)\}}.$$

Note that $\bar{\theta}_t := \arg\min_{\theta \in \cup_{\tau \in [t, T]} \Theta_\tau(\delta, \lambda)} \dot{\mu}(x_t^\top \theta)$. Let τ' be an arbitrary τ whose $\Theta_\tau(\delta, \lambda)$ contains $\bar{\theta}_t$. Then, the latter term in the square root is then bounded as follows:

$$\begin{aligned}
\sum_{t \in I_T} \{\dot{\mu}(x_t^\top \bar{\theta}_t) - \dot{\mu}(x_{t*}^\top \theta^*)\} &= \sum_{t \in I_T} \{\dot{\mu}(x_t^\top \bar{\theta}_t) - \dot{\mu}(x_t^\top \theta^*)\} + \sum_{t \in I_T} \{\dot{\mu}(x_t^\top \theta^*) - \dot{\mu}(x_{t*}^\top \theta^*)\} \\
&\leq M_\mu \left\{ \sum_{t \in I_T} |\mu(x_t^\top \bar{\theta}_t) - \mu(x_t^\top \theta^*)| + \sum_{t \in I_T} |\mu(x_t^\top \theta^*) - \mu(x_{t*}^\top \theta^*)| \right\} \quad (\text{Lemma F.3}) \\
&\leq M_\mu \left\{ \sum_{t \in I_T} |\mu(x_t^\top \bar{\theta}_t) - \mu(x_t^\top \hat{\theta}_{\tau'})| + \sum_{t \in I_T} |\mu(x_t^\top \hat{\theta}_{\tau'}) - \mu(x_t^\top \theta^*)| + \sum_{t \in I_T} \{\mu(x_{t*}^\top \theta^*) - \mu(x_t^\top \theta^*)\} \right\} \\
&\leq M_\mu \left\{ 2 \sum_{t \in I_T} |\mu(x_t^\top \omega_t) - \mu(x_t^\top \hat{\theta}_{\tau(t)})| + \text{Reg}(T) \right\} \leq M_\mu \{2\text{Reg}_{\text{EST}} + \text{Reg}(T)\},
\end{aligned}$$

where the second inequality holds from triangular inequality. For the last inequality, we leverage the definition of $\bar{\theta}_t$ and the pair $(\tau(t), \omega_t)$, as shown below:

$$\begin{aligned}
|\mu(x_t^\top \bar{\theta}_t) - \mu(x_t^\top \hat{\theta}_{\tau'})| &\leq \max_{\theta \in \Theta_{\tau'}(\delta, \lambda)} |\mu(x_t^\top \theta) - \mu(x_t^\top \hat{\theta}_{\tau'})| \\
&\leq \max_{\tau \in [t, T], \theta \in \Theta_\tau(\delta, \lambda)} |\mu(x_t^\top \theta) - \mu(x_t^\top \hat{\theta}_\tau)| = |\mu(x_t^\top \omega_t) - \mu(x_t^\top \hat{\theta}_{\tau(t)})|,
\end{aligned}$$

and that under the event \hat{E}_T , as $\theta^* \in \cup_{\tau \in [t, T]} \Theta_\tau(\delta, \lambda)$, we can easily show that the second term $|\mu(x_t^\top \hat{\theta}_{\tau'}) - \mu(x_t^\top \theta^*)|$ can be upper bounded by $|\mu(x_t^\top \omega_t) - \mu(x_t^\top \hat{\theta}_{\tau(t)})|$. Note that after decomposing $\text{Reg}(T)$, we are currently in the process of bounding its individual components. However, during this process, $\text{Reg}(T)$ itself appears in the upper bound. This issue will be addressed in Step 5, where we will provide a strategy to effectively resolve this recursive dependency.

Bounding $\sum_t E_t^{\text{FP}}$ We define the Gram matrix $\bar{V}_t := \lambda \mathbf{I} / \kappa + \sum_{\tau=1}^{t-1} x_\tau x_\tau^\top$. The introduction of this matrix is to apply EPL as in the previous section. As the weight $\dot{\mu}(x_t^\top \bar{\theta}_t)$ on the weighted Gram matrix \hat{H}_t is greater than or equal to $\kappa := \min_{x \in \mathcal{X}_{[T]}, \theta \in \Theta} \dot{\mu}(x^\top \theta)$, we have:

$$\hat{H}_t = \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} \dot{\mu}(x_\tau^\top \hat{\theta}_\tau) x_\tau x_\tau^\top \succeq \kappa \left(\frac{\lambda}{\kappa} \mathbf{I} + \sum_{\tau=1}^{t-1} x_\tau x_\tau^\top \right) \succeq \kappa \bar{V}_t. \quad (\text{E.6})$$

Then, we can bound $\sum_t E_t^{\text{FP}}$ as following:

$$\begin{aligned}
\sum_{t \in I_T} E_t^{\text{FP}} &\leq \gamma_T(\delta') \sum_{t \in I_T} \left| \dot{\mu}(x_t^\top \bar{\theta}_t) - \dot{\mu}(x_t^\top \hat{\theta}_t) \right| \|x_t\|_{\hat{H}_t^{-1}} \\
&\leq M_\mu \gamma_T(\delta') \sum_{t \in I_T} \left| \mu(x_t^\top \bar{\theta}_t) - \mu(x_t^\top \hat{\theta}_t) \right| \|x_t\|_{\hat{H}_t^{-1}} \quad (\text{Lemma F.3}) \\
&\leq M_\mu \gamma_T(\delta') \sqrt{\frac{1}{\kappa}} \sum_{t \in I_T} \left| \mu(x_t^\top \omega_t) - \mu(x_t^\top \hat{\theta}_{\tau(t)}) \right| \|x_t\|_{\hat{V}_t^{-1}} \\
&\leq \frac{2}{\kappa} M_\mu \mathcal{L}_\mu \gamma_T(\delta') \beta_T(\delta') \sum_{t \in I_T} \|x_t\|_{\hat{V}_t^{-1}}^2, \quad (\text{Lemma F.4})
\end{aligned}$$

where the third inequality holds from the definition of the pair $(\tau(t), \omega_t)$, and by (E.6) Here, Lemma F.4 plays a crucial role by introducing the weighted norm $\|x_t\|_{\hat{V}_t^{-1}}$, which enables the application of Lemma F.1. Utilizing this lemma, we further proceed to bound as:

$$\begin{aligned}
\sum_{t \in I_T} E_t^{\text{FP}} &\leq \frac{2}{\kappa} M_\mu \mathcal{L}_\mu \gamma_T(\delta') \beta_T(\delta') \sum_{t=1}^T \min \left\{ 1, \|x_t\|_{\hat{V}_t^{-1}}^2 \right\} \quad (\text{Definition of } I_T) \\
&\leq \frac{4}{\kappa} d M_\mu \mathcal{L}_\mu \gamma_T(\delta')^2 \log \left(1 + \frac{\kappa T}{d\lambda} \right), \quad (\text{Lemma F.1})
\end{aligned}$$

where the final inequality use the fact that $\gamma_T(\delta') = c_{\delta'} \cdot \beta_T(\delta') \geq \beta_T(\delta')$, as $c_{\delta'} \geq 1$.

Bounding $\sum_t F_t^{\text{FP}}$ The process closely resembles that of bounding $\sum_t E_t^{\text{FP}}$. we have:

$$\begin{aligned}
\sum_{t \in I_T} F_t^{\text{FP}} &\leq M_\mu \gamma_T(\delta')^2 \sum_{t \in I_T} \bar{\alpha}_t(\tilde{x}_t) \|x_t\|_{\hat{H}_t^{-1}}^2 \\
&\leq \frac{\mathcal{L}_\mu}{2\kappa} M_\mu \gamma_T(\delta')^2 \sum_{t \in I_T} \|x_t\|_{\hat{V}_t^{-1}}^2 \quad (\text{By (E.4) and (E.6)}) \\
&\leq \frac{\mathcal{L}_\mu}{2\kappa} M_\mu \gamma_T(\delta')^2 \sum_{t=1}^T \min \left\{ 1, \|x_t\|_{\hat{V}_t^{-1}}^2 \right\} \quad (\text{Definition of } I_T) \\
&\leq \frac{d}{\kappa} M_\mu \mathcal{L}_\mu \gamma_T(\delta')^2 \log \left(1 + \frac{\kappa T}{d\lambda} \right). \quad (\text{Lemma F.1})
\end{aligned}$$

Combining all three terms, we derive an upper bound for $\text{Reg}_{\text{FP}} \leq \sum_t \{D_t^{\text{FP}} + E_t^{\text{FP}} + F_t^{\text{FP}}\}$ as follows:

$$\begin{aligned}
\text{Reg}_{\text{FP}} &\leq \gamma_T(\delta') \sqrt{2d \log \left(1 + \frac{\mathcal{L}_\mu T}{d\lambda} \right)} \sqrt{\kappa_* T + M_\mu \{2\text{Reg}_{\text{EST}} + \text{Reg}(T)\}} \\
&\quad + \frac{5d}{\kappa} M_\mu \mathcal{L}_\mu \gamma_T(\delta')^2 \log \left(1 + \frac{\kappa T}{d\lambda} \right). \quad (\text{E.7})
\end{aligned}$$

Step 4 (Bounding Reg_{EST}) Next, we proceed to bound Reg_{EST} . The overall process closely mirrors that of Reg_{FP} . To begin, we decompose each instantaneous regret term. Let $\bar{\alpha}_t(\theta_1, \theta_2) = \int_0^1 (1-u) \dot{\mu}(x_t^\top \theta_1 + u(x_t^\top \theta_2 - x_t^\top \theta_1)) du$, which is derived from Taylor's theorem when estimating the regret for the selected arm x_t using two parameter vectors, θ_1 and θ_2 . Using Proposition F.2, we

expand the instantaneous regret term as:

$$\begin{aligned}
|\mu(x_t^\top \omega_t) - \mu(x_t^\top \hat{\theta}_{\tau(t)})| &= \left| \dot{\mu}(x_t^\top \hat{\theta}_{\tau(t)}) \langle x_t, \omega_t - \hat{\theta}_{\tau(t)} \rangle + \int_{x_t^\top \hat{\theta}_{\tau(t)}}^{x_t^\top \omega_t} (\mu(x_t^\top \omega_t) - z) \ddot{\mu}(z) dz \right| \\
&\leq \dot{\mu}(x_t^\top \hat{\theta}_{\tau(t)}) \left| \langle x_t, \omega_t - \hat{\theta}_{\tau(t)} \rangle \right| + \langle x_t, \omega_t - \hat{\theta}_{\tau(t)} \rangle^2 \int_0^1 (1-u) \left| \ddot{\mu}(x_t^\top \hat{\theta}_{\tau(t)} + u(x_t^\top \omega_t - x_t^\top \hat{\theta}_{\tau(t)})) \right| du \\
&\leq \dot{\mu}(x_t^\top \hat{\theta}_t) \left| \langle x_t, \omega_t - \hat{\theta}_{\tau(t)} \rangle \right| + M_\mu \langle x_t, \omega_t - \hat{\theta}_{\tau(t)} \rangle^2 \underbrace{\int_0^1 (1-u) \dot{\mu}(x_t^\top \hat{\theta}_{\tau(t)} + u(x_t^\top \omega_t - x_t^\top \hat{\theta}_{\tau(t)})) du}_{=\bar{\alpha}_t(\hat{\theta}_{\tau(t)}, \omega_t)} \\
&\leq \left(\dot{\mu}(x_t^\top \bar{\theta}_t) + \left| \dot{\mu}(x_t^\top \bar{\theta}_t) - \dot{\mu}(x_t^\top \hat{\theta}_{\tau(t)}) \right| \right) \left| \langle x_t, \omega_t - \hat{\theta}_{\tau(t)} \rangle \right| + M_\mu \bar{\alpha}_t(\hat{\theta}_{\tau(t)}, \omega_t) \langle x_t, \omega_t - \hat{\theta}_{\tau(t)} \rangle^2 \\
&\leq \left(\dot{\mu}(x_t^\top \bar{\theta}_t) + \left| \dot{\mu}(x_t^\top \bar{\theta}_t) - \dot{\mu}(x_t^\top \hat{\theta}_{\tau(t)}) \right| \right) \left| \langle x_t, \omega_t - \hat{\theta}_{\tau(t)} \rangle \right| + \underbrace{M_\mu \bar{\alpha}_t(\hat{\theta}_{\tau(t)}, \omega_t) \beta_{\tau(t)}(\delta')^2 \|x_t\|_{\hat{H}_{\tau(t)}^{-1}}^2}_{F_t^{\text{EST}}},
\end{aligned}$$

where the second inequality follows from Assumption 2, and the last inequality results from the Cauchy-Schwartz inequality. The key distinction between bounding Reg_{FP} and Reg_{EST} lies is the use of $\beta(\delta')$ instead of $\gamma(\delta')$. The first term then can be upper bounded by the following terms, using triangular inequality:

$$\begin{aligned}
&\left(\dot{\mu}(x_t^\top \bar{\theta}_t) + \left| \dot{\mu}(x_t^\top \bar{\theta}_t) - \dot{\mu}(x_t^\top \hat{\theta}_{\tau(t)}) \right| \right) \|x_t\|_{\hat{H}_{\tau(t)}^{-1}} \|\omega_t - \hat{\theta}_{\tau(t)}\|_{\hat{H}_{\tau(t)}} \\
&\leq \underbrace{\dot{\mu}(x_t^\top \bar{\theta}_t) \beta_{\tau(t)}(\delta') \|x_t\|_{\hat{H}_{\tau(t)}^{-1}}}_{D_t^{\text{EST}}} + \underbrace{\left| \dot{\mu}(x_t^\top \bar{\theta}_t) - \dot{\mu}(x_t^\top \hat{\theta}_{\tau(t)}) \right| \beta_{\tau(t)}(\delta') \|x_t\|_{\hat{H}_{\tau(t)}^{-1}}}_{E_t^{\text{EST}}}.
\end{aligned}$$

Bounding $\sum_t D_t^{\text{EST}}$ By the equation(E.5), instead of t -dependent \hat{H}_t , we use τ -only dependent \bar{H}_t to upper bound each term. We have:

$$\begin{aligned}
\sum_{t \in I_T} D_t^{\text{EST}} &\leq \beta_T(\delta') \sum_{t \in I_T} \dot{\mu}(x_t^\top \bar{\theta}_t) \|x_t\|_{\bar{H}_t^{-1}} \quad (t \leq \tau(t) \leq T) \\
&\leq \beta_T(\delta') \sqrt{2d \log \left(1 + \frac{\mathcal{L}_\mu T}{d\lambda} \right)} \sqrt{\sum_{t=1}^T \dot{\mu}(x_{t*}^\top \theta^*) + \sum_{t \in I_T} \{ \dot{\mu}(x_t^\top \bar{\theta}_t) - \dot{\mu}(x_{t*}^\top \theta^*) \}} \\
&\leq \beta_T(\delta') \sqrt{2d \log \left(1 + \frac{\mathcal{L}_\mu T}{d\lambda} \right)} \sqrt{\kappa_* T + M_\mu \{ 2\text{Reg}_{\text{EST}} + \text{Reg}(T) \}}.
\end{aligned}$$

Since the derivation proceeds identically to the case of $\sum_t D_t^{\text{FP}}$, with the only difference being the use of β_T in place of γ_T , we omit the detailed derivation here for brevity.

Bounding $\sum_t E_t^{\text{EST}}$ Following the bounding process of $\sum_t E_t^{\text{FP}}$, and using the equation (E.6),

$$\begin{aligned}
\sum_{t \in I_T} E_t^{\text{EST}} &\leq \beta_T(\delta') \sum_{t \in I_T} \left| \dot{\mu}(x_t^\top \bar{\theta}_t) - \dot{\mu}(x_t^\top \hat{\theta}_t) \right| \|x_t\|_{\hat{H}_t^{-1}} \quad (t \leq \tau(t) \leq T) \\
&\leq \frac{4d}{\kappa} M_\mu \mathcal{L}_\mu \beta_T(\delta')^2 \log \left(1 + \frac{\kappa T}{d\lambda} \right).
\end{aligned}$$

Bounding $\sum_t F_t^{\text{EST}}$ For $\bar{\alpha}_t(\hat{\theta}_{\tau(t)}, \omega_t) := \int_0^1 (1-u) \dot{\mu}(x_t^\top \hat{\theta}_{\tau(t)} + u(x_t^\top \omega_t - x_t^\top \hat{\theta}_{\tau(t)})) du$, as shown in (E.4), it follows that $\bar{\alpha}_t(\hat{\theta}_{\tau(t)}, \omega_t) \leq \mathcal{L}_\mu/2$. Following the bounding process of $\sum_t F_t^{\text{FP}}$,

$$\begin{aligned}
\sum_{t \in I_T} F_t^{\text{EST}} &\leq M_\mu \beta_T(\delta')^2 \sum_{t \in I_T} \bar{\alpha}_t(\hat{\theta}_{\tau(t)}, \omega_t) \|x_t\|_{\hat{H}_t^{-1}}^2 \quad (t \leq \tau(t) \leq T) \\
&\leq \frac{d}{\kappa} M_\mu \mathcal{L}_\mu \beta_T(\delta')^2 \log \left(1 + \frac{\kappa T}{d\lambda} \right).
\end{aligned}$$

In summary, combining three regret components, we bound $\text{Reg}_{\text{EST}} \leq \sum_t \{D_t^{\text{EST}} + E_t^{\text{EST}} + F_t^{\text{EST}}\}$ as follows:

$$\begin{aligned} \text{Reg}_{\text{FP}} &\leq \beta_T(\delta') \sqrt{2d \log \left(1 + \frac{\mathcal{L}_\mu T}{d\lambda}\right)} \sqrt{\kappa_* T + M_\mu \{2\text{Reg}_{\text{EST}} + \text{Reg}(T)\}} \\ &\quad + \frac{5d}{\kappa} M_\mu \mathcal{L}_\mu \beta_T(\delta')^2 \log \left(1 + \frac{\kappa T}{d\lambda}\right). \end{aligned} \quad (\text{E.8})$$

Step 5 (Solving equation) The equation (E.8) is indeed $c_{\delta'} > 1$ times the equation (E.7). Therefore, both Reg_{FP} and Reg_{EST} can be bounded using the regret bound of Reg_{FP} . Let $\text{Reg}_{\text{max}} = \max\{\text{Reg}_{\text{FP}}, \text{Reg}_{\text{EST}}\}$, and we can bound Reg_{max} with the equation (E.7). Then we have,

$$\text{Reg}(T) \leq \left(\frac{2}{p} + 1\right) \text{Reg}_{\text{FP}} + \text{Reg}_{\text{EST}} + \varepsilon \leq \left(\frac{2}{p} + 2\right) \text{Reg}_{\text{max}} + \varepsilon,$$

and accordingly, we get the following inequality:

$$\begin{aligned} \text{Reg}_{\text{max}} &\leq \gamma_T(\delta') \sqrt{2d \log \left(1 + \frac{\mathcal{L}_\mu T}{d\lambda}\right)} \sqrt{\kappa_* T + M_\mu \left\{ \left(\frac{2}{p} + 4\right) \text{Reg}_{\text{max}} + \varepsilon \right\}} \\ &\quad + \frac{5d}{\kappa} M_\mu \mathcal{L}_\mu \gamma_T(\delta')^2 \log \left(1 + \frac{\kappa T}{d\lambda}\right) \end{aligned}$$

Now, The bound takes the form of $\text{Reg}_{\text{max}} \leq A\sqrt{B + C\text{Reg}_{\text{max}}} + D$, where the following quantities are defined:

$$\begin{aligned} A &:= \gamma_T(\delta') \sqrt{2d \log \left(1 + \frac{\mathcal{L}_\mu T}{d\lambda}\right)} = \tilde{\mathcal{O}}(d) \\ B &:= \kappa_* T + M_\mu \varepsilon = \tilde{\mathcal{O}}\left(\kappa_* T + d\sqrt{\frac{T}{\lambda^2}}\right) \\ C &:= M_\mu \left(\frac{2}{p} + 4\right) = \mathcal{O}(1) \\ D &:= \frac{5d}{\kappa} M_\mu \mathcal{L}_\mu \gamma_T(\delta')^2 \log \left(1 + \frac{\kappa T}{d\lambda}\right) = \tilde{\mathcal{O}}\left(\frac{d^2}{\kappa}\right), \end{aligned}$$

focusing on the terms involving d, T , and κ . With the choice of $\lambda = \mathcal{O}(d)$ and applying Lemma F.5, the upper bound for Reg_{max} simplifies to:

$$\text{Reg}_{\text{max}} = \tilde{\mathcal{O}}\left(d\sqrt{\kappa_* T} + d^2 + \frac{d^2}{\kappa}\right) = \tilde{\mathcal{O}}\left(d\sqrt{\kappa_* T} + \frac{d^2}{\kappa}\right).$$

Now, combining all terms, the cumulative regret $R(T)$ can be bounded as:

$$\begin{aligned} R(T) &\leq \text{Reg}(T) + \frac{4dR^*}{\log 2} \left\{ \log \left(1 + \frac{\mathcal{L}_\mu}{\lambda \log 2}\right) + \log \left(1 + \frac{\kappa}{\lambda \log 2}\right) \right\} \\ &\lesssim \left(\frac{2}{p} + 2\right) \text{Reg}_{\text{max}} + \varepsilon \\ &= \tilde{\mathcal{O}}\left(d\sqrt{\kappa_* T} + \frac{d^2}{\kappa} + d\sqrt{\frac{T}{\lambda^2}}\right) = \tilde{\mathcal{O}}\left(d\sqrt{\kappa_* T} + \frac{d^2}{\kappa}\right). \quad (\text{Choose } \lambda = \mathcal{O}(d)) \end{aligned}$$

Our derived regret bound, $\tilde{\mathcal{O}}(d\sqrt{\kappa_* T} + d^2/\kappa)$, aligns with the state-of-the-art regret guarantee [4, 15, 36, 37].

F Proof supplement and Lemmas

In this section, we provide key propositions, lemmas, and inequality bounds that are essential for the main proof. These results serve as the mathematical foundation for the regret analysis and other theoretical guarantees established in this work.

F.1 Supporting Lemmas for main proof

Proposition F.1 (Lemma D.1. of Lee et al. [37]). *Let μ be increasing and self-concordant with M_μ . Let $\mathcal{Z} \subseteq \mathcal{B}(S) := \{z \in \mathbb{R} \mid |z| \leq S\}$ in \mathbb{R} . Then for any $z_1, z_2 \in \mathcal{Z}$, the following holds:*

$$\int_0^1 (1-u) \dot{\mu}(z_1 + u(z_2 - z_1)) du \geq \frac{\dot{\mu}(z_1)}{2 + 2SM_\mu}$$

This proposition establishes a lower bound for the integral involving a self-concordant, particularly useful for controlling weighted norms and derivatives in regret analysis. The following proposition provides the well-known Taylor's expansion expression with integral remainder. We highlight specific cases that are frequently used throughout the main proof.

Proposition F.2 (Taylor's Theorem with Integral Remainder Form). *Let $n \geq 0$ be an integer and let the function $f : \mathbb{R} \rightarrow \mathbb{R}$ be $(n+1)$ times differentiable at the point $x_0 \in \mathbb{R}$. Let $f^{(n)}$ to denote its n -th derivatives, then $f(x)$ can be expressed as:*

$$f(x) = \sum_{i=0}^n \frac{f^{(i)}(x_0)}{i!} (x - x_0)^i + \frac{1}{n!} \int_{x_0}^x f^{(n+1)}(t) (x - t)^n dt$$

Epecially for $n = 0$, by letting $t = x_0 + u(x - x_0)$, $f(x)$ can be expressed as:

$$f(x) = f(x_0) + \int_{x_0}^x f'(t) dt = f(x_0) + (x - x_0) \int_0^1 f'(x_0 + u(x - x_0)) du \quad (\text{F.1})$$

Epecially for $n = 1$, by letting $t = x_0 + u(x - x_0)$, $f(x)$ can be expressed as:

$$\begin{aligned} f(x) &= f(x_0) + f'(x_0)(x - x_0) + \int_{x_0}^x f''(t)(x - t) dt \\ &= f(x_0) + f'(x_0)(x - x_0) + \int_0^1 f''(x_0 + u(x - x_0)) ((1-u)(x - x_0)) \cdot (x - x_0) du \\ &= f(x_0) + f'(x_0)(x - x_0) + (x - x_0)^2 \int_0^1 f''(x_0 + u(x - x_0)) (1-u) du \end{aligned} \quad (\text{F.2})$$

Similarly, with multivariate function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ and $n = 1$, we have that:

$$\begin{aligned} f(x) &= f(x_0) + \sum_{i=1}^d \frac{\partial}{\partial x_i} f(x_0) (x - x_0)_i + \sum_{i=1}^d (x - x_0)_i^2 \int_0^1 (1-u) \frac{\partial^2}{\partial x_i^2} f(x_0 + u(x - x_0)) du \\ &\quad + \sum_{i \neq j} 2(x - x_0)_i (x - x_0)_j \int_0^1 (1-u) \frac{\partial}{\partial x_i} \frac{\partial}{\partial x_j} f(x_0 + u(x - x_0)) du \\ &= f(x_0) + \nabla f(x_0)^\top (x - x_0) + (x - x_0)^\top \left(\int_0^1 (1-u) \nabla^2 f(x_0 + u(x - x_0)) du \right) (x - x_0). \end{aligned}$$

The following lemmas provide critical tools for analyzing regret bounds in the context of bandit algorithms. Specifically, they focus on bounding terms that involve weighted norms of feature vectors under Gram matrices. These results are pivotal in deriving efficient bounds on cumulative regret and exploration.

Lemma F.1 (Elliptical Potential Lemma). *let $\{x_t\}_{t=1}^T$ be a sequence in \mathbb{R}^d satisfying $\|x_t\| \leq R$ for all $t \leq T$. For a gram matrix $V_t := \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} x_\tau x_\tau^\top$, we have that*

$$\sum_{t=1}^T \min \left\{ 1, \|x_t\|_{V_t^{-1}}^2 \right\} \leq 2d \log \left(1 + \frac{R^2 T}{d\lambda} \right).$$

Lemma F.2 (Elliptical Potential Count Lemma). *For $R, c > 0$, let $\{x_t\}_{t=1}^T$ be a sequence in \mathbb{R}^d satisfying $\|x_t\| \leq R$ for all $t \leq T$. For a gram matrix $V_t := \lambda \mathbf{I} + \sum_{\tau=1}^{t-1} x_\tau x_\tau^\top$, the length of the sequence $\mathcal{N}_T := \left\{ t \in [T] \mid \|x_t\|_{V_t^{-1}} > c \right\}$ is bounded as:*

$$|\mathcal{N}_T| \leq \frac{2d}{\log(1+c^2)} \log \left(1 + \frac{R^2}{\lambda \log(1+c^2)} \right)$$

The following connects the smoothness of $\dot{\mu}$ to the difference in μ .

Lemma F.3. *For $x, y \in \mathbb{R}$, $|\dot{\mu}(x) - \dot{\mu}(y)| \leq M_\mu |\mu(x) - \mu(y)|$.*

Proof.

$$\begin{aligned} |\dot{\mu}(x) - \dot{\mu}(y)| &= \left| (x-y) \int_0^1 \ddot{\mu}(y+u(x-y)) du \right| && \text{(By (F.1) with } f = \dot{\mu}) \\ &\leq |x-y| \int_0^1 |\ddot{\mu}(y+u(x-y))| du \\ &\leq M_\mu |x-y| \int_0^1 \dot{\mu}(y+u(x-y)) du && \text{(Assumption 2)} \\ &= M_\mu \left| (x-y) \int_0^1 \dot{\mu}(y+u(x-y)) du \right| \\ &= M_\mu |\mu(x) - \mu(y)| && \text{(By (F.1) with } f = \mu) \end{aligned}$$

where the second equality holds since the integral term is positive as $\dot{\mu} \geq 0$. \square

The following lemma bounds the difference in μ values, using linearization as in the previous works.

Lemma F.4. *For any $t \geq 1$ and $\theta_1, \theta_2 \in \Theta_{\tau(t)}$, we have the following:*

$$|\mu(x_t^\top \theta_1) - \mu(x_t^\top \theta_2)| \leq 2\gamma_T(\delta) \mathcal{L}_\mu \sqrt{\frac{1}{\kappa}} \|x_t\|_{V_t^{-1}}$$

Proof.

$$\begin{aligned} |\mu(x_t^\top \theta_1) - \mu(x_t^\top \theta_2)| &= \left| \langle x_t, \theta_1 - \theta_2 \rangle \int_0^1 \dot{\mu}(x_t^\top \theta_2 + u(x_t^\top \theta_1 - x_t^\top \theta_2)) du \right| && \text{(By (F.1))} \\ &\leq \left\{ \|x_t\|_{\hat{H}_{\tau(t)}^{-1}} \cdot \|\theta_1 - \theta_2\|_{\hat{H}_{\tau(t)}} \right\} \int_0^1 \mathcal{L}_\mu du && \text{(Cauchy-Schwartz, } \mathcal{L}_\mu\text{-Lipschitzness)} \\ &\leq \mathcal{L}_\mu \|x_t\|_{\hat{H}_{\tau(t)}^{-1}} \left\{ \|\theta_1 - \hat{\theta}_{\tau(t)}\|_{\hat{H}_{\tau(t)}} + \|\theta_2 - \hat{\theta}_{\tau(t)}\|_{\hat{H}_{\tau(t)}} \right\} && \text{(Triangle inequality)} \\ &\leq 2\beta_{\tau(t)}(\delta') \mathcal{L}_\mu \|x_t\|_{\hat{H}_{\tau(t)}^{-1}} \\ &\leq 2\beta_{\tau(t)}(\delta') \mathcal{L}_\mu \sqrt{\frac{1}{\kappa}} \|x_t\|_{V_{\tau(t)}^{-1}} && \text{(By (E.6))} \\ &\leq 2\beta_T(\delta') \mathcal{L}_\mu \sqrt{\frac{1}{\kappa}} \|x_t\|_{V_t^{-1}} && (t \leq \tau(t) \leq T) \end{aligned}$$

\square

Lemma F.5. *Let $A, B, C, D, X \in \mathbb{R}^+$. The following implication holds:*

$$X \leq A\sqrt{B+CX} + D \implies X \leq 2(A\sqrt{B} + A^2C + D)$$

Proof. let $f : x \rightarrow x^2 - px - q$ for $p, q > 0$. Then the roots for $f(x) = 0$ are:

$$x_1, x_2 = \frac{p \pm \sqrt{p^2 + 4q}}{2}.$$

Now, as f is a convex function, $x^2 \leq px + q$ implies:

$$x \leq \max\{x_1, x_2\} \leq \frac{p + \sqrt{p^2 + 4q}}{2} \leq \frac{p + (p + 2\sqrt{q})}{2} = p + \sqrt{q}. \quad (\text{triangle inequality})$$

And accordingly, we have:

$$\begin{aligned} x \leq p\sqrt{x} + q &\implies \sqrt{x} \leq p + \sqrt{q} \\ &\implies x \leq (p + \sqrt{q})^2 \leq 2p^2 + 2q \end{aligned} \quad (\text{F.3})$$

where the inequality holds from $(x + y)^2 \leq 2(x^2 + y^2)$. Then according to the equation (F.3),

$$\begin{aligned} X \leq A\sqrt{B + CX} + D &\implies X \leq A\sqrt{C}\sqrt{X} + A\sqrt{B} + D \quad (\text{triangle inequality}) \\ &\implies X \leq 2(A\sqrt{C})^2 + 2(A\sqrt{B} + D) = 2(A\sqrt{B} + A^2C + D) \end{aligned}$$

□

F.2 Auxiliary bounding inequalities

We introduce key probabilistic inequalities and bounds frequently used in the analysis of randomized algorithms. These results provide tools to bound the probabilities of deviations and concentration of random variables, which are essential for deriving high-probability guarantees in the main analysis.

Proposition F.3 (Azuma's inequality). *If a super-martingale $(X_t)_{t \geq 0}$ corresponding to a filtration \mathcal{H}_{t-1} satisfies $|X_t - X_{t-1}| < c_t$ for some constant c_t for all $t = 1, \dots, T$ then for any $\alpha > 0$,*

$$\mathbb{P}(X_T - X_0 \geq \alpha) \leq 2 \exp \left(-\frac{\alpha^2}{2 \sum_{t=1}^T c_t^2} \right).$$

Proposition F.4 (Chernoff bound). *For a random variable X and its moment-generating function $M(t) = \mathbb{E}[e^{tX}]$,*

$$\mathbb{P}(X \geq \alpha) \leq \inf_{t > 0} M(t)e^{-t\alpha}.$$

Accordingly, for a random variable following a standard normal distribution (i.e. $X \sim \mathcal{N}(0, 1)$),

$$\mathbb{P}(X \geq \alpha) \leq \inf_{t > 0} \exp\left(\frac{t^2}{2} - t\alpha\right) = e^{-\frac{\alpha^2}{2}}.$$

Lemma F.6. *Let z be a random variable sampled from the standard Normal distribution. Then for all $\delta \in (0, 1)$,*

$$\mathbb{P}(|z| \leq \sqrt{2 \log(2/\delta)}) \geq 1 - \delta.$$

Proof. By proposition Proposition F.4, $\mathbb{P}(|z| > \alpha) = 2\mathbb{P}(z > \alpha) \leq 2e^{-\frac{\alpha^2}{2}}$. Set the right-hand side as δ , and get $\alpha = \sqrt{2 \log(2/\delta)}$. □

G Limitations

Our analysis focuses on structured settings such as generalized linear bandits (GLBs), where feature perturbation achieves both strong empirical performance and provable regret guarantees. While the same principle shows promise in more flexible or non-linear models, theoretical guarantees in these broader settings remain open. The current formulation, though practically effective, is heuristic outside the GLM framework and lacks formal justification under complex function classes. Extending the theory to overparameterized or general Lipschitz models represents an important direction for future work, where feature-level stochasticity may offer a stable alternative to parameter perturbation.

H Experimental settings and additional results

H.1 Experimental details

The GLB experiments are entirely synthetic and computationally lightweight. All GLB runs were performed on a standard CPU server equipped with an Intel Xeon Silver 4210R processor (40 threads), with each run completing within a few minutes. The neural bandit experiments were conducted using a single NVIDIA RTX 3090 GPU. Due to the moderate model size and limited data horizon, neural runs for each algorithm completed in under one hour. Overall, the compute requirements were modest, and no large-scale pretraining or extensive hyperparameter tuning was necessary for any of the experiments.

H.1.1 Generalized linear bandit settings

We consider a time-varying set of $K = 100$ arms per round over a horizon of $T = 20,000$ (linear) or $T = 10,000$ (logistic). Context vectors and the true parameter vector θ^* are sampled from a standard multivariate Gaussian distribution and normalized to satisfy the boundedness assumption. In the logistic case, we use the sigmoid link function $\mu(x) = 1/(1 + e^{-x})$ and constrain $\|\theta^*\| \leq 4$ so that the logits lie in $[-4, 4]$. For the linear bandit, the reward noise is sampled from $\mathcal{N}(0, 1)$. We vary the feature dimension d over $\{10, 20, 40\}$. The confidence level is set as $\delta = 1/T$, and the regularization parameter is fixed at $\lambda = 10^{-4}$.

Baselines and tuning. For linear bandits, we compare against ε -greedy [33], LinUCB [1], LinTS [6], LinPHE [31], and RandLinUCB [46]. For logistic bandits, we include ε -greedy, OFUL-GLB-e [37], LogTS [32], LogPHE [32], and RandUCBLog. All hyperparameters, such as inflation factors and learning rates, are tuned according to the original papers. We fix $\mathcal{L}_\mu = 0.25$ and use a minimum link derivative of 0.25 for numerical stability in the logistic setting. For TS-based algorithms and our method, we set the inflation parameter as $c_t = 1$. For ε -greedy, we use an annealing schedule $\varepsilon_t = \varepsilon\sqrt{T/t}$ with $\varepsilon = 0.05$.

H.1.2 Neural bandit settings

Objective. Following Zhang et al. [51], we evaluate neural contextual bandits using classification tasks on UCI benchmark datasets [40]: `shuttle`, `isolet`, and `mushroom`. These tasks are transformed into multi-armed bandit problems via a disjoint model construction (cf. Li et al. [38]). For a k -class classification problem with input dimension d , we construct a kd -dimensional feature representation by placing the input vector x into the corresponding class slot: $x_1 = (x; \mathbf{0}; \dots; \mathbf{0})$, $x_2 = (\mathbf{0}; x; \dots; \mathbf{0})$, and so on. Each x_i is treated as the feature vector for the i -th arm. A neural model f predicts the reward for each x_i , and the agent selects the arm with the highest predicted reward. A reward of 1 is given if the selected arm corresponds to the correct label, and 0 otherwise. Regret is measured as the cumulative number of classification errors. All experiments are repeated five times with shuffled data. The time horizon is set to $T = 10,000$ for all datasets.

Neural networks. To study the effect of model capacity, we evaluate two neural network architectures: a shallow and a deep model. The shallow network consists of a single hidden layer with a fully connected layer of size $d \times 100$, followed by a ReLU activation, a final fully connected layer of size 100×1 , and a softmax output layer, following the design used in Zhang et al. [51]. The deep network, in contrast, uses two hidden layers with a fully connected layer of size $d \times 50$, a ReLU activation, another fully connected layer of size 50×50 followed by ReLU, and a final output layer of size 50×1 with softmax. Both models are trained online using the Adam optimizer with a learning rate of 0.001 and a batch size of 32. Training is performed at each round using the most recent 32 observed examples. Figure 3(bottom) is a result with a deeper model.

Baselines and tuning. We compare against several baselines including ε -greedy, NeuralUCB [52], NeuralTS [51], and FTPL [32]. To accelerate training for NeuralUCB and NeuralTS, we approximate the confidence matrix inverse using only the reciprocals of the diagonal entries. In our proposed algorithm, DeepFP, instead of perturbing the full kd -dimensional parameter space, we selectively perturb only the d -dimensional subspace corresponding to the chosen arm. This masking avoids interference from other arms, ensuring more accurate value estimation for the selected action.

H.2 Additional experiments

We provide additional experiments by varying the norm of the true parameter S and the cardinality of the context set $|\mathcal{C}|$. The specific configurations are provided in each figure caption.

To further validate our algorithm in a neural contextual bandit setting, we conduct additional experiments using the MNIST [12] and Fashion-MNIST [48] datasets. Each instance is represented as a 28×28 grayscale image with label set size $K = 10$, naturally forming a K -armed classification bandit problem. At each round, the agent receives a context composed of K image tensors, where the i -th arm corresponds to the image being placed in the i -th channel slot, and all others set to zero. This yields a $K \times 28 \times 28$ tensor input for each arm. The model f used to estimate the expected reward is a shared convolutional neural network that takes the entire K -channel tensor as input and outputs a scalar score for each arm, promoting parameter sharing across arms.

The architecture of the model is as follows: two convolutional layers are applied with ReLU activation and 2×2 max pooling after each. The first layer uses 32 filters and the second 64 filters, both with kernel size 3×3 and stride 1. The resulting $64 \times 7 \times 7$ feature map is flattened and passed through two fully connected layers with hidden size 128, followed by ReLU and a final output scalar. This design enables efficient learning of visual features while maintaining compatibility with the bandit framework through arm-wise shared representation.

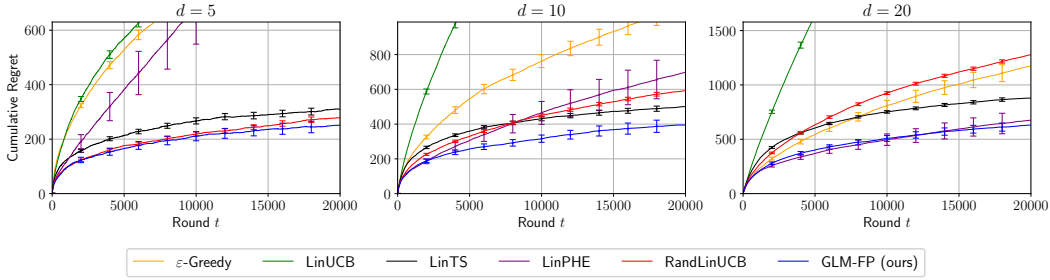


Figure H.1: Linear Bandit. $|\mathcal{C}| = 1$, $d = \{5, 10, 20\}$, $K = 100$, $S = 1$.

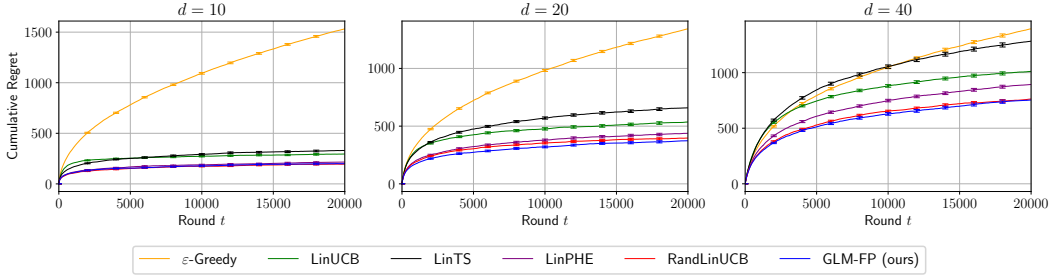


Figure H.2: Linear Bandit. $|\mathcal{C}| = T$, $d = \{10, 20, 40\}$, $K = 100$, $S = 2$.

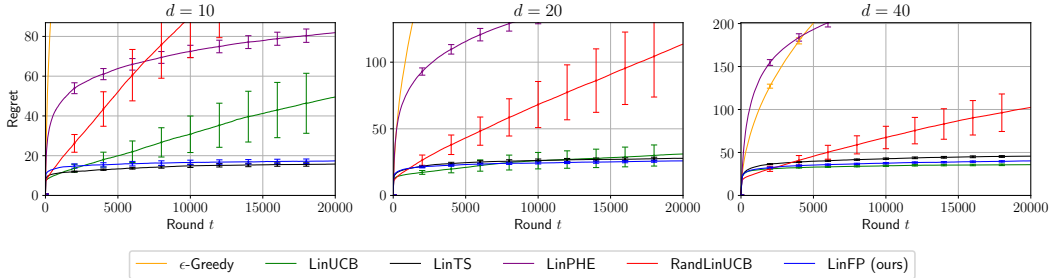


Figure H.3: Linear Bandit with noise $\mathcal{N}(0, 0.1^2)$. $|\mathcal{C}| = T$, $d = \{10, 20, 40\}$, $K = 100$, $S = 2$.

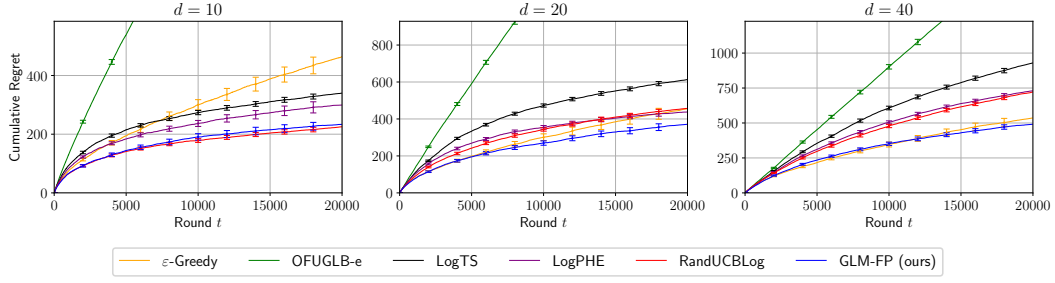


Figure H.4: Logistic Bandit. $|\mathcal{C}| = 1$, $d = \{5, 10, 20\}$, $K = 100$, $S = 1$.

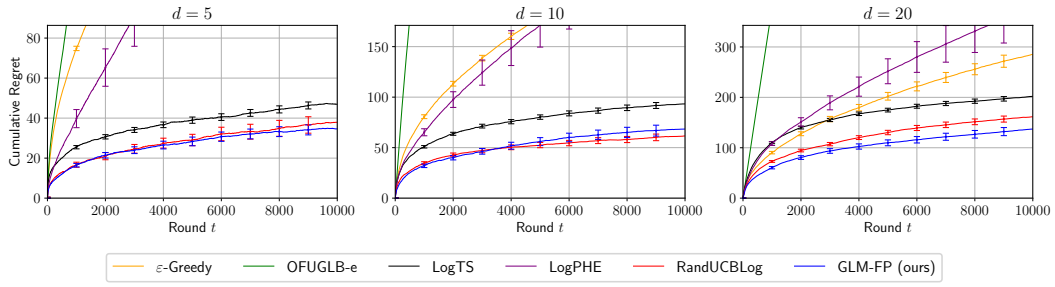


Figure H.5: Logistic Bandit. $|\mathcal{C}| = 1$, $d = \{5, 10, 20\}$, $K = 100$, $S = 4$.

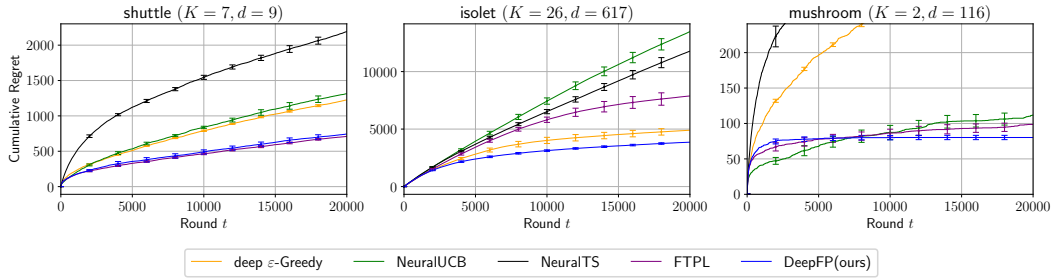


Figure H.6: Neural Bandit. Multi-layer perceptron model with one hidden layer and output layer.

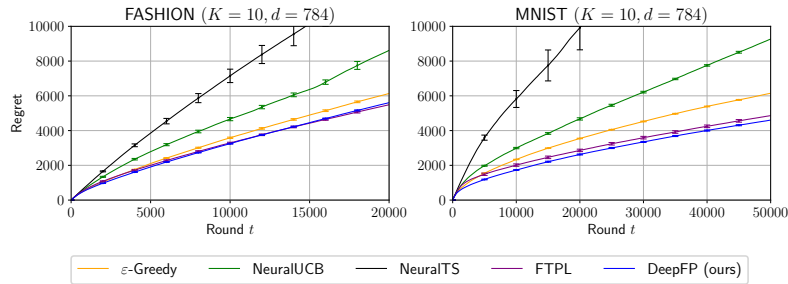


Figure H.7: Neural Bandit. Experiments on MNIST dataset with a CNN model.