Variance-Dependent Regret Lower Bounds for Contextual Bandits

Anonymous Author(s) Affiliation Address

email

Abstract

Variance-dependent regret bounds for linear contextual bandits, which improve upon the classical $\widetilde{O}(d\sqrt{K})$ regret bound to $\widetilde{O}(d\sqrt{\sum_{k=1}^K \sigma_k^2})$, where d is the context dimension, K is the number of rounds, and σ_k^2 is the noise variance in round k, has been widely studied in recent years. However, most existing works focus on the regret upper bounds instead of lower bounds. To our knowledge, the only lower bound is from Jia et al. (2024), which proved that for any eluder dimension d_{elu} and total variance budget Λ , there exists an instance with $\sum_{k=1}^K \sigma_k^2 \leq \Lambda$ for which any algorithm incurs a variance-dependent lower bound of $\Omega(\sqrt[n]{d_{\mathbf{elu}}\Lambda})$. However, this lower bound has a \sqrt{d} gap with existing upper bounds. Moreover, it only considers a fixed total variance budget Λ and does not apply to a general variance sequence $\{\sigma_1^2, \dots, \sigma_K^2\}$. In this paper, to overcome the limitations of Jia et al. (2024), we consider the general variance sequence under two settings. For a prefixed sequence, where the entire variance sequence is revealed to the learner at the beginning of the learning process, we establish a variancedependent lower bound of $\Omega(d\sqrt{\sum_{k=1}^K \sigma_k^2}/\log K)$ for linear contextual bandits. For an adaptive sequence, where an adversary can generate the variance σ_k^2 in each round k based on historical observations, we show that when the adversary must generate σ_k^2 before observing the decision set \mathcal{D}_k , a similar lower bound of $\Omega(d\sqrt{\sum_{k=1}^K \sigma_k^2/\log^6(dK)})$ holds. In both settings, our results match the upper bounds of the SAVE algorithm (Zhao et al., 2023) up to logarithmic factors. Furthermore, if the adversary can generate the variance σ_k after observing the decision set \mathcal{D}_k , we construct a counter-example showing that it is impossible to construct a variance-dependent lower bound if the adversary properly selects variances in collaboration with the learner. Our lower bound proofs use a novel peeling technique that groups rounds by variance magnitude. For each group, we construct separate instances and assign the learner distinct decision sets. We believe this proof technique may be of independent interest.

1 Introduction

2

5

6

10

11

12

13

14

15

16

17

18

19

20

21

23

24

25

26

27

28

29

30

31

35

36

We consider the linear contextual bandit problem, where each arm is represented by a feature vector and the expected reward is a linear function of this feature vector with an unknown parameter vector. Numerous studies have developed algorithms achieving optimal regret bounds for linear bandits (Chu et al., 2011; Abbasi-Yadkori et al., 2011a). However, while these works establish minimax-optimal regret bounds in the worst-case, they do not exploit additional problem-dependent structures. Our work focuses on incorporating reward variance information into the analysis, building upon a line of research studying variance-dependent regret bounds for linear bandits (Zhou et al., 2021; Zhang et al., 2021; Zhou and Gu, 2022; Zhao et al., 2022; Kim et al., 2022; Zhao et al., 2023) and general function approximation (Jia et al., 2024), which includes linear bandits as a special Submitted to 39th Conference on Neural Information Processing Systems (NeurIPS 2025). Do not distribute.

case. Notably, Zhao et al. (2023) established a near-optimal regret guarantee without requiring prior knowledge of the variances:

Theorem 1.1 (Theorem 2.3, Zhao et al. 2023). For any linear contextual bandit problem, the regret of the SAVE algorithm in the first K rounds is upper bounded by:

Regret
$$(K) \le \widetilde{O}\left(d\sqrt{\sum_{k=1}^{K} \sigma_k^2} + d\right),$$

where d is the dimension and σ_k^2 is the noise variance of the selected action in round k.

However, most of these works have focused on developing algorithms with regret upper bound guarantees, while variance-dependent lower bounds remain understudied. The only exception is Jia et al. (2024), which focuses on general function classes with finite eluder dimension $d_{\rm elu}$ and provides the following variance-dependent lower bound:

Theorem 1.2 (Theorem 5.1, Jia et al. 2024). For any dimension $d \geq 2$, action space size A, number of rounds $K \geq 2$, and total variance budget $\Lambda \in [0,K]$, there exists a contextual bandit problem with eluder dimension $d_{\mathbf{elu}} = d$, action space size A, and an adversarial sequence of variances satisfying $\sum_{k=1}^{K} \sigma_k^2 \leq \Lambda \text{ such that for any algorithm, the regret is lower bounded by:}$

$$\operatorname{Regret}(K) \ge \Omega(\min(\sqrt{d\Lambda} + d, \sqrt{AK})).$$

When restricted to the linear bandit case, where $d \geq \sqrt{A}$, the above lower bound reduces to $\sqrt{d\Lambda}$, which has a gap of \sqrt{d} factor compared with the upper bound in Zhao et al. (2023). Moreover, Jia et al. (2024) only considers instances with a fixed budget Λ and relies on carefully designed variance sequences $\{\sigma_1^2, \sigma_2^2, \ldots, \sigma_K^2\}$, failing to provide lower bounds for general variance sequences. Therefore, an open question arises:

Can we prove variance-dependent regret lower bounds for general variance sequences?

1.1 Our Contributions

56

57

76

77

78

79

In this paper, we answer this question affirmatively by constructing hard-to-learn instances in sev-58 eral different settings. For any prefixed sequence $\{\sigma_1^2,\ldots,\sigma_K^2\}$, we achieve a $\widetilde{\Omega}(d\sqrt{\sum_{k=1}^K\sigma_k^2})$ 59 variance-dependent expected lower bound, which matches the upper bound in Zhao et al. (2023) 60 up to logarithmic factors and demonstrates its optimality. For general adaptive variance sequences 61 where a weak adversary (potentially collaborating with the learner) can generate variance σ_k^2 in each 62 round k based on historical observations, our instance provides a high-probability lower bound of 63 $\widetilde{\Omega}(d\sqrt{\sum_{k=1}^K \sigma_k^2})$, which also matches the upper bound in Zhao et al. (2023) up to logarithmic fac-64 tors. To the best of our knowledge, this is the first high-probability lower bound for linear contextual 65 66

Our construction and analysis rely on the following new techniques:

- A peeling technique for prefixed variance sequences that divides rounds into groups based on variance magnitude. Through orthogonal decision set construction, each group only interacts with its corresponding parameters, allowing us to establish separate lower bounds for different variance scales and combine them effectively.
- A multi-instance framework that handles unknown group sizes in the adaptive setting. For each variance group, we maintain multiple instances designed for different possible intervals of round numbers and assign the learner to these instances in a cyclic manner, ensuring uniform visits across instances and guaranteeing the visiting times of one instance matches its designed interval.
 - A high-probability lower bound that handles adaptive group sizes through a union bound. We
 first convert expected regret bounds to constant-probability bounds through careful variance control and auxiliary algorithms, then boost these to high-probability bounds by creating multiple
 independent instances.

Furthermore, we also study the setting with a strong adversary that can generate the variance σ_k after observing the decision set \mathcal{D}_k . Under this scenario, we proposed a counter algorithm that can collaborate with the adversary by properly selecting variance, achieving an O(d) regret even the total variance $\sum_{k=1}^K \sigma_k^2 = \Omega(K)$. This implies that it is impossible to derive a variance-dependent lower bound for general variance sequence with strong adversary. As a direct extension of this result,

we also show that it is impossible to derive a variance-dependent lower bound for stochastic linear bandits, where the decision set is fixed even for a general prefixed variance sequence.

Notation We use lower case letters to denote scalars, and use lower and upper case bold face letters to denote vectors and matrices respectively. We denote by [n] the set $\{1,\ldots,n\}$. For a vector $\mathbf{x} \in \mathbb{R}^d$ and a positive semi-definite matrix $\mathbf{\Sigma} \in \mathbb{R}^{d \times d}$, we denote by $\|\mathbf{x}\|_2$ the vector's ℓ_2 norm and by $\|\mathbf{x}\|_{\mathbf{\Sigma}} = \sqrt{\mathbf{x}^\top \mathbf{\Sigma} \mathbf{x}}$ the Mahalanobis norm. For two positive sequences $\{a_n\}$ and $\{b_n\}$ with $n=1,2,\ldots$, we write $a_n=O(b_n)$ if there exists an absolute constant C>0 such that $a_n \leq Cb_n$ holds for all $n \geq 1$ and write $a_n = \Omega(b_n)$ if there exists an absolute constant C>0 such that $a_n \geq Cb_n$ holds for all $n \geq 1$. We use $O(\cdot)$ to further hide the polylogarithmic factors. We use $O(\cdot)$ to denote the indicator function.

2 Related Work

95

Heteroscedastic Linear Bandits. For linear bandit problems, the worst-case regret has been widely 96 studied (Auer, 2002; Dani et al., 2008; Li et al., 2010; Chu et al., 2011; Abbasi-Yadkori et al., 2011b; 97 Li et al., 2019), achieving $O(\sqrt{K})$ bounds in the first K rounds. Recently, a series of works has 98 considered heteroscedastic variants where noise distributions vary across rounds. Kirschner and 99 Krause (2018) first formally proposed a linear bandit model with heteroscedastic noise, assuming 100 σ_k -sub-Gaussian noise in round $k \in [K]$. Subsequently, (Zhou et al., 2021; Zhang et al., 2021; Kim 101 et al., 2021; Zhou and Gu, 2022; Dai et al., 2022; Zhao et al., 2023; Jia et al., 2024) relaxed this to 102 variance-based constraints where round k has variance σ_k^2 . Among these works, Zhou et al. (2021) 103 and Zhou and Gu (2022) obtained near-optimal regret guarantees of $\widetilde{O}(d\sqrt{\sum_{k=1}^K \sigma_k^2})$, but required 104 knowledge of σ_k after observing the reward in round k. In contrast, Zhang et al. (2021); Kim et al. 105 (2021) handled unknown variances with computationally inefficient algorithms, achieving a weaker 106 $\widetilde{O}(\mathrm{poly}(d)\sqrt{\sum_{k=1}^K \sigma_k^2})$ bound. Recently, Zhao et al. (2023) improved upon these results with an efficient algorithm (SAVE) achieving the near-optimal $\widetilde{O}(d\sqrt{\sum_{k=1}^K \sigma_k^2})$ bound without requiring 108 variance knowledge. Beyond standard linear bandits, two directions have been explored. Dai et al. 109 (2022) studied heteroscedastic sparse linear bandits, providing a framework to convert standard 110 algorithms to the sparse setting. In a different direction, Jia et al. (2024) extended the analysis 111 to contextual bandits with general function classes having finite eluder dimension, which includes 112 linear bandits as a special case, and achieved a variance-dependent regret upper bounds. 113 **Lower Bounds for Linear Contextual Bandits.** For linear contextual bandit problems, several works (Dani et al., 2008; Chu et al., 2011; Li et al., 2019) have established theoretical lower bounds 115 to illustrate the fundamental difficulty in learning process. For linear bandits with finite action sets, 116 Chu et al. (2011) established an $\Omega(\sqrt{dK})$ lower bound, matching the upper bound up to logarithmic 117 factors in the action set size and number of rounds K. For general stochastic linear bandits, Dani 118 et al. (2008) constructed an instance with $2^{\Omega(d)}$ actions and obtained an $\Omega(d\sqrt{K})$ lower bound. Later, Li et al. (2019) focused on linear contextual bandits, where the decision set can vary across 120 rounds, and provided an $\Omega(d\sqrt{K\log K})$ lower bound. However, all these works only focus on 121 worst-case regret bounds and do not consider the heteroscedastic variance information. The only 122 exception is Jia et al. (2024), which provided an $\Omega(\sqrt{d\Lambda})$ variance-dependent lower bound for a 123 fixed total variance budget Λ . Nevertheless, this work cannot handle general variance sequences and 124 leaves open the question of variance-dependent lower bounds in the general setting. 125

3 Preliminaries

126

127

128

129

130

131

132

133

135

In this work, we consider the heteroscedastic linear contextual bandit (Zhou et al., 2021; Zhang et al., 2021), where the noise variance varies across rounds. Let K be the total number of rounds. In each round $k \in [K]$, the interaction between the learner and the environment proceeds as follows:

- 1. The environment generates an arbitrary decision set $\mathcal{D}_k \subseteq \mathbb{R}^d$, where each element represents a feasible action that can be selected by the learner;
- 2. The learner observes \mathcal{D}_k and selects $\mathbf{x}_k \in \mathcal{D}_k$;
- 3. The environment generates the stochastic noise ϵ_k and reveals the stochastic reward $r_k = \langle \boldsymbol{\mu}, \mathbf{x}_k \rangle + \epsilon_k$ to the learner, where $\boldsymbol{\mu} \in \mathbb{R}^d$ is the unknown weight vector for the underlying linear reward function.

Without loss of generality, we assume the random noise ϵ_k in each round k satisfies:

$$\mathbb{P}(|\epsilon_k| \le R) = 1, \quad \mathbb{E}[\epsilon_k | \mathbf{x}_{1:k}, \epsilon_{1:k-1}] = 0, \quad \mathbb{E}[\epsilon_k^2 | \mathbf{x}_{1:k}, \epsilon_{1:k-1}] = \sigma_k^2 \le 1, \forall k \in [K]$$
 (3.1)

For any algorithm Alg and linear bandit instance \mathcal{M} , the cumulative regret is defined as follows:

$$\operatorname{Regret}_{\operatorname{Alg}}(K, \mathcal{M}) = \sum_{k \in [K]} \langle \mathbf{x}_k^*, \boldsymbol{\mu} \rangle - \langle \mathbf{x}_k, \boldsymbol{\mu} \rangle, \quad \text{where } \mathbf{x}_k^* = \underset{\mathbf{x} \in \mathcal{D}_k}{\operatorname{argmax}} \langle \mathbf{x}, \boldsymbol{\mu} \rangle. \tag{3.2}$$

For simplicity, we may omit the subscripts Alg and/or \mathcal{M} when there is no ambiguity. Additionally, with a slight abuse of notation, we may use σ_k to represent the variance σ_k^2 (which is originally the standard deviation) when there is no ambiguity. In this work, we focus on providing variance-dependent lower bounds for the regret based on the variances sequence $\{\sigma_1,...,\sigma_K\}$. We consider two settings for the variance sequence $\{\sigma_1,...,\sigma_K\}$:

- Prefixed Sequence: The variance sequence is revealed to the learner at the beginning of the learning process.
- Adaptive Sequence: An adversary (potentially collaborating with the learner) can generate
 the variance σ_k in each round k based on historical observations, with the learner receiving
 each variance at the beginning of the corresponding round. This setting can be further
 divided into two categories based on the power of the adversary:
 - Weak Adversary: The adversary must generate the variance σ_k before observing the decision set \mathcal{D}_k .
 - Strong Adversary: The adversary can generate the variance σ_k after observing the decision set \mathcal{D}_k .

Remark 3.1. Unlike the typical adversarial setting focused on maximizing regret for a specific algorithm, our work uses the idea of an "adversary" to represent the environment's inherent ability to select the variance sequence. This "adversary" might even strategically choose variance levels (σ_k) based on the **past decision sets** \mathcal{D}_k **observed so far**, potentially leading to variance levels that could temporarily improve the learner's performance or make the learning process appear easier. This seeming "cooperation," however, is ultimately aimed at exploring the fundamental lower bounds on regret that must hold for any learner in any environment. The key is that the variance is chosen without direct knowledge of the true underlying patterns μ . When this "adversary" (our "strong adversary") can adjust the variance based on the learner's actions (\mathcal{D}_k) , this strategic "cooperation," informed by past observations but blind to μ , becomes more effective in probing the true limits of learnability and challenging our lower bound results.

4 Variance-Dependent Lower Bound with Prefixed Variance Sequence

In this section, we consider the setting where the variance sequence $\{\sigma_1, \dots, \sigma_K\}$ is prefixed and fully revealed to the learner at the beginning of the learning process.

4.1 Main Results

We establish the following theorem for the variance-dependent lower bound.

Theorem 4.1. For any dimension d>1, prefixed sequence of variance $\{\sigma_1,...,\sigma_K\}$ satisfying $\sum_{k=1}^K \sigma_k^2 \ge 1 + 384d^2$ and algorithm Alg, there exists a hard linear contextual bandit instance such that each action $a \in \mathcal{D}_k$ in round k has variance bounded by σ_k . For this instance, the expected regret of algorithm Alg over K rounds is lower bounded by:

$$\mathbb{E}[\operatorname{Regret}(K)] \ge \Omega\left(d\sqrt{\sum_{i=1}^K \sigma_k^2}/(\log K)\right).$$

Remark 4.2. For a prefixed sequence $\{\sigma_1,...,\sigma_K\}$, Theorem 4.1 shows that any algorithm incurs a regret lower bounded of $\widetilde{\Omega}(d\sqrt{\sum_{k=1}^K\sigma_k^2})$, which matches the upper bound in Zhao et al. (2023) up to logarithmic factors. Compared to the lower bound in Jia et al. (2024), Theorem 4.1 focuses on the linear contextual bandit setting and achieves a \sqrt{d} improvement over the standard linear bandit setting. It is also worth noting that the lower bound in Jia et al. (2024) only considers instances with a fixed total variance $\sum_{k=1}^K \sigma_k^2$, constructed by using constant variance in the early rounds and zero variance in later rounds. In comparison, Theorem 4.1 applies to any fixed variance sequence and is more flexible.

In Theorem 4.1, we require that the total variance is no less than $\Omega(d^2)$, which reduces to $K \geq \Omega(d^2)$ 181 when all variances $\sigma_k = 1$. A similar requirement exists in standard linear bandits, since a trivial 182 lower bound of $\Omega(K)$ always holds for any algorithm, and the lower bound of $\Omega(d\sqrt{K})$ can only 183 be achieved when $K > \Omega(d^2)$. Furthermore, for general sequences of variances with total variance 184 smaller than $O(d^2)$, a large number of rounds K alone is not sufficient to establish the desired 185 lower bound. The presence of early rounds with zero variance would increase the total number of 186 rounds without affecting the fundamental complexity of the problem. This observation suggests that 187 requiring total variance no less than $\Omega(d^2)$ (or other equivalent conditions) may be necessary for 188 establishing the lower bound.

4.2 Proof of Theorem 4.1

190

191

192

193

194

195

196

197

210

211

212

227

228

In this subsection, we prove the variance-dependent lower bound in Theorem 4.1. We first start with a fixed variance threshold σ , and construct a class of hard-to-learn instances where actions are chosen from a hypercube action set $\mathcal{A} = \{-1,1\}^d$, and for any action $\mathbf{a} \in \mathcal{A}$, the reward follows a scaled Bernoulli distribution $\sigma \cdot B(1/3 + \langle \boldsymbol{\mu}, \boldsymbol{a} \rangle)$, where $\Delta = 1/\sqrt{96K}$ and $\boldsymbol{\mu} \in \{-\Delta, \Delta\}^d$. In this setting, the variance for each action is upper bounded by σ^2 , and these instances can be represented as a linear bandit problem with feature $(\sigma, \sigma \cdot \mathbf{a})$ and weight vector $\boldsymbol{\mu}' = (1/3, \boldsymbol{\mu})$. Based on these hard-to-learn instances, we have the following variance-dependent lower bound for the regret:

Lemma 4.3. For a fixed variance threshold σ and any bandit algorithm Alg, if the weight vector $\mu \in \{-\Delta, \Delta\}^d$ is uniformly random selected from $\{-\Delta, \Delta\}^d$, the variance in each round is bounded by σ^2 , and the expected regret over $K \geq 1.5 \cdot d^2$ rounds is lower bounded by:

$$\mathbb{E}_{\mu}[\operatorname{Regret}(K)] \ge d\sqrt{K\sigma^2}/8\sqrt{6}.$$

Remark 4.4. Lemma 4.3 establishes a variance-dependent lower bound for the regret with a fixed variance threshold σ . When all variances are equal ($\sigma_1 = ... = \sigma_K = \sigma$), this bound matches the upper bound in Zhao et al. (2023) up to logarithmic factors. In addition, under this fixed-variance setting, this lemma provides a tighter logarithmic dependency on the number of rounds K compared to Theorem 4.1, though it does not extend to dynamic variances.

Now, for any prefixed variance sequence $\{\sigma_1,...,\sigma_K\}$, we divide the rounds into $L = \lceil \log_2 K \rceil + 1$ different groups based on the range of their variance as follows:

$$\mathcal{K}_0 = \{k : \sigma_k \le 1/K\},\$$

 $\mathcal{K}_i = \{k : 2^{i-1}/K < \sigma_k \le 2^i/K\}, \text{ for } i = 1, \dots, L-1.$

For each group K_i with $i \in [L-1]$, we construct a bandit instance M_i with weight vector μ_i following Lemma 4.3, where:

- the variance threshold is set to be $\sigma(i) = 2^{i-1}/K$;
- the number of rounds is $K_i = |\mathcal{K}_i|$;
- the dimension is $d_i = d/L$.

For group \mathcal{K}_0 , we construct a different type of instance \mathcal{M}_0 : a d/L-armed bandit, where one randomly chosen arm gives constant reward 1 while all other arms give reward 0. Note that this instance in \mathcal{M}_0 can be equivalently represented as a $d_0 = d/L$ -dimensional linear bandit where actions are one-hot vectors \mathbf{e}_i .

Based on these sub-instances, we create a combined linear bandit instance with dimension 217 $d_0+d_1+...+d_{L-1}=d$ with weight vector $\boldsymbol{\mu}=(\boldsymbol{\mu}_0,...,\boldsymbol{\mu}_{L-1})$: At the beginning of 218 each round k, if round k belongs to group \mathcal{K}_i , then the learner receives the decision set $\mathcal{D}_k = \{(\mathbf{0}_{d_0},...,\mathbf{0}_{d_{i-1}},\mathbf{x},\mathbf{0}_{d_{i+1}},...,\mathbf{0}_{d_{L-1}}): \mathbf{x} \in \mathcal{A}_i\}$, where $\mathbf{0}_{d_j}$ corresponds to a zero vector with dimension d_j and \mathcal{A}_i is the action set in the bandit instance \mathcal{M}_i . Under this construction, for any round 219 220 221 $k \in \mathcal{K}_i$, the reward in the combined instance coincides with that of sub-instance \mathcal{M}_i . Specifically, 222 after the learner selects action x, they receive a reward drawn from a scaled Bernoulli distribution 223 with variance upper bounded by $\sigma^2(i) = (2^{i-1}/K)^2$ for $i \neq 0$, and variance 0 for i = 0. Note that in all groups, the variance is bounded by σ_k^2 . With this construction in hand, we now proceed to 224 225 prove the lower bound in Theorem 4.1.

Remark 4.5 (Linear Contextual Bandits vs. Stochastic Linear Bandits). In the proof of Theorem 4.1, we heavily rely on assigning different decision sets to rounds in the contextual bandit environment. This approach, however, does not extend to stochastic linear bandit problems, where all

rounds share the same decision set. To see this limitation, consider any prefixed variance sequence with $\sigma_1 = \cdots = \sigma_d = 0$. In this case, the learner can select canonical basis of the decision set in the first d rounds. Since these rounds have zero variance, the learner learns the exact rewards for all actions in the decision set and incurs no regret in subsequent rounds, regardless of the values of $\sigma_{d+1}, \ldots, \sigma_K$. Consequently, it is impossible to establish a lower bound of $\widetilde{\Omega}(d\sqrt{\sum_{k=1}^K \sigma_k^2})$ in this setting.

Proof of Theorem 4.1. Due to the orthogonal construction of decision sets across different groups \mathcal{K}_i , actions in group \mathcal{K}_i provide no information about the weight vector $\boldsymbol{\mu}_j$ for $j \neq i$. Consequently, the total regret can be decomposed into the sum of regrets from each sub-instance. For each sub-instance \mathcal{M}_i with $i \neq 0$, the regret is lower bounded by:

$$\mathbb{E}_{\boldsymbol{\mu}_{i}} \left[\sum_{k \in \mathcal{K}_{i}} \max_{\mathbf{x} \in \mathcal{D}_{k}} \langle \boldsymbol{\mu}_{i}, \mathbf{x} \rangle - \langle \boldsymbol{\mu}_{i}, \mathbf{x}_{k} \rangle \right] \geq \mathbb{I}(K_{i} \geq 1.5d_{i}^{2}) \cdot \frac{d_{i}\sqrt{K_{i}\sigma^{2}(i)}}{8\sqrt{6}}$$

$$\geq \frac{d_{i}\sqrt{K_{i}\sigma^{2}(i)}}{8\sqrt{6}} - \frac{d_{i}\sqrt{1.5d_{i}^{2} \cdot \sigma^{2}(i)}}{8\sqrt{6}}$$

$$\geq \frac{d_{i}\sqrt{\sum_{k \in \mathcal{K}_{i}} \sigma_{k}^{2}}}{16\sqrt{6}} - \frac{d_{i}^{2} \cdot \sigma(i)}{16}, \tag{4.1}$$

where the first inequality follows from Lemma 4.3, the second inequality holds due to $\mathbb{I}(x \geq y)\sqrt{x} \geq \sqrt{x} - \sqrt{y}$, and the last inequality follows from the definition of group \mathcal{K}_i .

Taking a summation of (4.1) over all groups, the total regret can be lower bounded as follows:

$$\mathbb{E}_{\boldsymbol{\mu}}[\operatorname{Regret}(K)] = \sum_{i=0}^{L-1} \mathbb{E}_{\boldsymbol{\mu}_{i}} \left[\sum_{k \in \mathcal{K}_{i}} \max_{\mathbf{x} \in \mathcal{D}_{k}} \langle \boldsymbol{\mu}_{i}, \mathbf{x} \rangle - \langle \boldsymbol{\mu}_{i}, \mathbf{x}_{k} \rangle \right]$$

$$\geq \sum_{i=1}^{L-1} \frac{d_{i} \sqrt{\sum_{k \in \mathcal{K}_{i}} \sigma_{k}^{2}}}{16\sqrt{6}} - \frac{d_{i}^{2} \cdot \sigma(i)}{16}$$

$$\geq \sum_{i=1}^{L-1} \frac{d \sqrt{\sum_{k \in \mathcal{K}_{i}} \sigma_{k}^{2}}}{16\sqrt{6}L} - \frac{d^{2}}{4L^{2}}$$

$$\geq \frac{d \sqrt{\sum_{i=1}^{L-1} \sum_{k \in \mathcal{K}_{i}} \sigma_{k}^{2}}}{16\sqrt{6}L} - \frac{d^{2}}{4L^{2}}, \tag{4.2}$$

where the first inequality follows from (4.1), the second inequality follows from the definition of variance threshold $\sigma(i)$ and dimension $d_i=d/L$, and the last inequality holds due to $\sum_i \sqrt{x_i} \geq \sqrt{\sum_i x_i}$. In addition, for the group \mathcal{K}_0 , we have

$$\sum_{k \in \mathcal{K}_0} \sigma_k^2 \le \sum_{k \in \mathcal{K}_0} 1/K \le 1,\tag{4.3}$$

where the first inequality follows from the definition of group K_0 and the second inequality follows from $|K_0| \le K$. Therefore, we have

$$\mathbb{E}_{\mu}[\text{Regret}(K)] \ge \frac{d\sqrt{\sum_{i=1}^{L-1} \sum_{k \in \mathcal{K}_i} \sigma_k^2}}{16\sqrt{6}L} - \frac{d^2}{4L^2}$$

$$\ge \frac{d\sqrt{\sum_{k=1}^{K} \sigma_k^2 - 1}}{16\sqrt{6}L} - \frac{d^2}{4L^2}$$

$$\ge \frac{d\sqrt{\sum_{k=1}^{K} \sigma_k^2 - 1}}{32\sqrt{6}L},$$

where the first inequality follows from (4.2), the second inequality follows from (4.3), and the last inequality follows from the fact that $\sum_{k=1}^{K} \sigma_k^2 \ge 1 + 384d^2$. Thus, we complete the proof of Theorem 4.1.

Variance-Dependent Lower Bounds with Adaptive Variance Sequence

In the previous section, we focused on the setting where the variance sequence is prefixed and revealed to the learner at the beginning of the learning process. In this section, we extend our analysis to the setting where the variance sequence can be adaptive based on historical observations, with the learner receiving the adaptive variance at the beginning of each round.

5.1 Main Results

251

252 253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271 272

276

281

282

283

284

285

286 287

288

289

290

291

292

293

294

297

298

299

300

5.1.1 Weak Adversary

We first describe the learning process and the mechanism of variance adaptation. In detail, the adaptive variance process proceeds as follows:

- 1. At the beginning of each round k, a (weak) adversary selects the variance level σ_k based on the historical observations, including actions $\{a_1, \ldots, a_{k-1}\}$, rewards $\{r_1, \ldots, r_{k-1}\}$, and decision sets $\{\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_{k-1}\}$. The adversary has access to all historical information but not to the underlying reward model parameters;
- 2. Given the selected variance level σ_k , we construct and assign a decision set \mathcal{D}_k to the learner, where the variance of the reward for each action $a \in \mathcal{D}_k$ is bounded by σ_k^2 ;
- 3. The learner observes the decision set \mathcal{D}_k and variance level σ_k , then determines an action a_k from \mathcal{D}_k based on its historical observations and current information. After selecting the action, the learner receives a reward r_k with variance bounded by σ_k^2 .

Remark 5.1. It is worth noting that our concept of adversary differs from the weak/strong adversary in Jia et al. (2024). Specifically, Jia et al. (2024) considers an adversary that attempts to hinder the learner's learning by allocating a fixed total variance budget $\sum_{k=1}^K \sigma_k^2 \le \Lambda$ across rounds to maximize regret. In contrast, our work considers an adversary that attempts to break the lower bounds themselves by collaborating with the learner. To prevent such exploitation, we must restrict the adversary from knowing the weight vector of the underlying reward model. Without this restriction, the adversary could encode each entry μ_i of the weight vector μ through the corresponding variance $\sigma_i = \mu_i$, allowing the learner to learn the weight vector after d rounds.

Under this setting, we establish the following theorem for the variance-dependent lower bound. 277

Theorem 5.2 (Weak Adversary). For any dimension d > 1, adaptive sequence of variances 278 $\{\sigma_1,\ldots,\sigma_K\}$ and algorithm Alg, there exists a hard instance such that each action $a\in\mathcal{D}_k$ in 279 round k has variance bounded by σ_k^2 . For this instance, if $\sum_{k=1}^K \sigma_k^2 \ge \Omega(d^2)$, then with probability at least 1-1/K, the regret of algorithm \mathbf{Alg} over K rounds is lower bounded by: 280

$$\operatorname{Regret}(K) \geq \Omega \Big(d \sqrt{\sum_{k=1}^K \sigma_k^2} / \log^6(dK) \Big).$$

Remark 5.3. Theorem 5.2 provides a high-probability lower bound of $\widetilde{\Omega}(d\sqrt{\sum_{k=1}^K \sigma_k^2})$, which matches the upper bound in Zhao et al. (2023) up to logarithmic factors, albeit with looser logarithmic dependencies than Theorem 4.1 due to the adaptive nature of the variance sequence. Unlike the expected lower bound in Theorem 4.1, for adaptive variance sequences, the cumulative variance $\sum_{k=1}^K \sigma_k^2$ depends on the random process and observations. This dependence makes it challenging to establish an expected variance-dependent regret bound - a fundamental difficulty that does not arise for standard $d\sqrt{K}$ -type lower bounds in linear contextual bandits. To the best of our knowledge, our result provides the first high-probability lower bound for linear contextual bandit.

5.1.2 Strong Adversary

In Theorem 5.2, we require that for each round $k \in [K]$, all actions $\mathbf{x} \in \mathcal{D}_k$ share the same adaptive variance σ_k . This is more restrictive than the setting in Zhao et al. (2023), where the variance can differ across actions $\mathbf{x} \in \mathcal{D}_k$. However, extending our lower bound to action-dependent variances is not possible in the adaptive setting. This limitation arises because we construct the decision set \mathcal{D}_k after the adversary chooses the variance σ_k , which prevents assigning specific variances to individual actions $\mathbf{x} \in \mathcal{D}_k$. Moreover, we now consider a strong adversary that can choose σ_k after observing the decision set \mathcal{D}_k . The interaction between the learner and this strong adversary proceeds as follows:

1. At the beginning of each round k, we construct and assign a decision set \mathcal{D}_k based on historical observations, including actions $\{a_1, \ldots, a_{k-1}\}$ and rewards $\{r_1, \ldots, r_{k-1}\}$;

- 2. Given the decision set \mathcal{D}_k in round k, the strong adversary selects the variance level σ_k for round k. The adversary has access to all historical information but not to the underlying reward model parameters;
- 3. The learner observes the decision set \mathcal{D}_k and variance level σ_k , then determines an action a_k from \mathcal{D}_k based on its historical observations and current information. After selecting the action, the learner receives a reward r_k with variance bounded by σ_k^2 .

The following theorem shows that under this setting, the adversary could cooperate with the learner to break the lower bound.

Theorem 5.4 (Strong Adversary). For any linear contextual bandit problem and number of rounds $K \geq 2d$, if we first provide the decision set \mathcal{D}_k and then allow an adversary to choose the variance σ_k based on the decision set \mathcal{D}_k , there exists one such type of adversary such that, there exists an algorithm whose regret in the first K rounds is upper bounded by $\operatorname{Regret}(K) \leq d$, where the total variance $\sum_{k=1}^K \sigma_k^2 \geq K/2$.

Remark 5.5. Theorem 5.4 highlights why Theorem 5.2 requires a weak adversary that set the variance sequence before seeing the learner's choices. If the adversary could see the decision set first, it could potentially choose variances that would invalidate our lower bound. This finding underscores that our construction is precise and pinpoints the exact condition under which the derived lower bound holds.

Remark 5.6. It is worth noting that Jia et al. (2024) also considered the case where the adversary assigns variances to actions after observing the decision set and action choice, and provided a variance-dependent lower bound. However, their analysis focuses on an adversary that allocates variance across rounds to maximize the regret. In contrast, our work considers an adversary that attempts to break these bounds, making it more challenging to establish lower bounds for general variance sequences. It is also worth noting that if the adversary's goal is to increase regret, choosing a prefixed sequence is a viable strategy. This case is already covered by our Theorem 4.1 for prefixed sequences, which provides a tighter lower bound than Theorem 5.2.

Theorem 5.4 suggests that it is impossible to derive a variance-dependent lower bound if the adversary can determine the variance σ_k after observing the decision set \mathcal{D}_k , which further precludes establishing a lower bound when the adversary has the ability to assign action-dependent variances for each action $\mathbf{x} \in \mathcal{D}_k$ after observing the decision set \mathcal{D}_k . This result naturally extends to stochastic linear bandit problems, where the decision set \mathcal{D} remains fixed across all rounds. In this case, since the adversary knows the decision set $\mathcal{D}_k = \mathcal{D}$ in advance, Theorem 5.4 directly implies:

Corollary 5.7. For any stochastic linear bandit problem with fixed decision set \mathcal{D} and number of rounds $K \geq 2d$, there exists a prefixed sequence $\{\sigma_1, \dots, \sigma_K\}$ such that there exists an algorithm whose regret in the first K rounds is upper bounded by: $\operatorname{Regret}_{\operatorname{Alg}}(K) \leq d$, where the total variance $\sum_{k=1}^K \sigma_k^2 \geq K/2$.

5.2 Proof Sketch of Theorem 5.2

In this section, we provide the proof sketch of Theorem 5.2. Overall, the proof follows a similar structure as Theorem 4.1, where we divide the rounds into several groups based on their variance magnitude and create hard instances for each group. The key idea is to calculate individual regret bounds for each group and combine them for the final lower bound. However, there exist several challenges when dealing with adaptive variance sequences that require careful handling.

Varying Size of Groups \mathcal{K}_i As discussed in Section 4.2, for each group \mathcal{K}_i , we create individual instance \mathcal{M}_i with fixed variance threshold $\sigma(i) = 2^{i-1}/K$ and establish a lower bound of $\widetilde{\Omega}(d_i\sqrt{\sigma^2(i)|\mathcal{K}_i|})$ on the expected regret. However, the construction of such instances relies on prior knowledge of the number of rounds $|\mathcal{K}_i|$, which can be calculated at the beginning for a prefixed variance sequence $\{\sigma_1,\ldots,\sigma_K\}$. In contrast, for general adaptive variance sequences, the number of rounds $|\mathcal{K}_i|$ is not known a priori and can even be a random variable, which creates a barrier in constructing these instances.

To address the unknown number of rounds $|\mathcal{K}_i|$, instead of constructing a single instance \mathcal{M}_i for each group, we create L instances $\mathcal{M}_{i,j}$, where $L = \lceil \log_2 K \rceil + 1$. Each instance $\mathcal{M}_{i,j}$ is designed for a specific range of round numbers, specifically $\mathcal{M}_{i,j}$ for $2^{j-1} \leq |\mathcal{K}_i| < 2^j$.

For each round k in group \mathcal{K}_i , the learner receives a decision set \mathcal{D}_i from one of the instances in $\{\mathcal{M}_{i,1},\ldots,\mathcal{M}_{i,L}\}$ in a cyclic manner. Through this sequential assignment, the number of visits to each instance $\mathcal{M}_{i,j}$ is $|\mathcal{K}_i|/L$. Consequently, we expect that the instance $\mathcal{M}_{i,j}$ corresponding to the

true range $2^{j-1} \leq |\mathcal{K}_i| < 2^j$ provides a lower bound of $\widetilde{\Omega}(d_i \sqrt{\sigma^2(i)|\mathcal{K}_i|}) = \widetilde{\Omega}(d_i \sqrt{\sigma^2(i) \cdot 2^j})$, which leads to the final lower bound of $\widetilde{\Omega}(d\sqrt{\sum_{k=1}^K \sigma_k^2})$. 357

establishing the lower bound for the triggered instance $\mathcal{M}_{i,j}$ corresponding to the true range $2^{j-1} \le$ $|\mathcal{K}_i| < 2^j$. Traditional analysis of lower bounds in linear contextual bandits has focused on the expected regret. However, when dealing with adaptive variance sequences, this approach becomes insufficient as the adversary can dynamically adjust the variance sequence to break these bounds. For instance, an adversary might continuously set $\sigma_k = 1$ until the lower bound of $\widetilde{\Omega}(d\sqrt{\sum_{i=1}^k \sigma_i^2})$ is violated at some round k, then switch to $\sigma_k = 0$ for all future rounds, causing the total variance sum $\sum_{k=1}^{K} \sigma_k^2$ to remain unchanged. In our construction, this means all rounds could fall into group $\mathcal{L}_{k=1}^{k=1}$ $\stackrel{k}{\sim}$ the adversary to adaptively change the number of rounds between different intervals $2^{j-1} \leq |\mathcal{K}_L| < 2^j$. Since the failure of the lower bound in any single instance $\mathcal{M}_{L,j}$ leads to failure of the whole construction, an expected lower bound on regret cannot guarantee robust performance against adaptive sequences. This necessitates a stronger high-probability lower bound that holds

Converting Expected Lower Bound to High-Probability Lower Bound. Another challenge is

uniformly for all instances. Unfortunately, an expectation of $\widetilde{\Omega}(d_i\sqrt{\sigma^2(i)2^j})$ in instance $\mathcal{M}_{i,j}$ only implies a low-probability 371 regret (Regret $\geq \widetilde{\Omega}(d_i \sqrt{\sigma^2(i)2^j})$) $\geq d_i \cdot 2^{-j/2}$, since the cumulative regret in \mathcal{K}_i can be up to $\sigma(i)$. 372 $|\mathcal{K}_i|$ in our instance. To solve this problem, we introduce an auxiliary algorithm that automatically 373 detects the cumulative regret and switches to the standard OFUL algorithm (Abbasi-Yadkori et al., 2011a) if the cumulative regret is larger than $\Omega(d_i\sqrt{\sigma^2(i)2^j})$. For this auxiliary algorithm, we can 375 guarantee that the upper bound is at most $\widetilde{\Omega}(d_i\sqrt{\sigma^2(i)2^j})$ while maintaining the same probability of high regret as the original algorithm. Therefore, an expectation of $\Omega(d_i \sqrt{\sigma^2(i)2^j})$ in instance $\mathcal{M}_{i,j}$ 377 implies a constant-probability regret $\mathbb{P}(\text{Regret} \geq \widetilde{\Omega}(d_i \sqrt{\sigma^2(i)2^j})) = \Omega(1)$. 378

After constructing an instance with constant-probability lower bound, we boost this probability by 379 creating $\Omega(\log^2(dK))$ independent instances. When the learner encounters instance $\mathcal{M}_{i,j}$, it is 380 assigned to one of these instances in a cyclic manner. Through this construction, with probability at 381 least 1 - 1/poly(K), the final regret is lower bounded by $\text{Regret} \geq \widetilde{\Omega}(d_i \sqrt{\sigma^2(i)2^j})$. 382

Remark 5.8. Unlike previous lower bounds for linear bandit problems which focus on expected regret, to the best of our knowledge, our result provides the first high-probability lower bound for linear contextual bandits. It is worth noting that our construction requires separate decision sets across different rounds in the random assignment process. For stochastic linear bandits with a fixed decision set, we can only derive a constant-probability lower bound. Moreover, for a fixed decision set in stochastic linear bandit problem with covering number $\log \mathcal{N} \leq O(d)$, an algorithm can randomly select one action from the covering set and perform this action in all rounds. In this case, there exists a probability of $1/\mathcal{N}=1/\exp(d)$ to achieve zero regret, which precludes the possibility of establishing high-probability lower bounds for large round numbers K. More details about the high-probability lower bound can be found in Section 5.2.

Conclusion and Future Work

358

359

360

362

363

364

365

366 367

368

369

370

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

In this paper, we study variance-dependent lower bounds for linear contextual bandits in different settings. For both prefixed and adaptive variance sequences with weak adversary, we establish tight lower bounds matching the upper bounds in Zhao et al. (2023) up to logarithmic factors. We further demonstrate a fundamental limitation: when a strong adversary can select variances after observing decision sets, it becomes impossible to establish meaningful variance-dependent lower bounds. However, our work has focused exclusively on linear bandit settings, while Jia et al. (2024) has established variance-dependent lower bounds for general function approximation with a fixed total variance budget Λ . Therefore, we leave for future work the generalization of our analysis of general variance sequence to contextual bandits with general function approximation.

¹In general settings, detecting cumulative regret is impossible as the learner lacks prior knowledge of the optimal reward and variance. However, in our lower bound construction, all instances are randomly selected from instance classes sharing the same optimal reward and variance, which are known to the learner. This knowledge enables the construction of the auxiliary algorithm.

References

- ABBASI-YADKORI, Y., PÁL, D. and SZEPESVÁRI, C. (2011a). Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*.
- ABBASI-YADKORI, Y., PÁL, D. and SZEPESVÁRI, C. (2011b). Improved algorithms for linear stochastic bandits. In *NIPS*, vol. 11.
- 408 AUER, P. (2002). Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research* **3** 397–422.
- CESA-BIANCHI, N. and LUGOSI, G. (2006). *Prediction, learning, and games*. Cambridge university press.
- 412 CHU, W., LI, L., REYZIN, L. and SCHAPIRE, R. (2011). Contextual bandits with linear payoff 413 functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence* 414 and Statistics. JMLR Workshop and Conference Proceedings.
- DAI, Y., WANG, R. and Du, S. S. (2022). Variance-aware sparse linear bandits. *arXiv preprint arXiv*:2205.13450.
- DANI, V., HAYES, T. P. and KAKADE, S. M. (2008). Stochastic linear optimization under bandit feedback.
- JIA, Z., QIAN, J., RAKHLIN, A. and WEI, C.-Y. (2024). How does variance shape the regret in contextual bandits? *arXiv preprint arXiv:2410.12713*.
- KIM, Y., YANG, I. and JUN, K.-S. (2021). Improved regret analysis for variance-adaptive linear bandits and horizon-free linear mixture mdps. *arXiv preprint arXiv:2111.03289*.
- KIM, Y., YANG, I. and JUN, K.-S. (2022). Improved regret analysis for variance-adaptive linear bandits and horizon-free linear mixture mdps. *Advances in Neural Information Processing*Systems 35 1060–1072.
- KIRSCHNER, J. and KRAUSE, A. (2018). Information directed sampling and bandits with heteroscedastic noise. In *Conference On Learning Theory*. PMLR.
- LI, L., CHU, W., LANGFORD, J. and SCHAPIRE, R. E. (2010). A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*.
- LI, Y., WANG, Y. and ZHOU, Y. (2019). Nearly minimax-optimal regret for linearly parameterized bandits. In *Conference on Learning Theory*. PMLR.
- ZHANG, Z., YANG, J., JI, X. and DU, S. S. (2021). Improved variance-aware confidence sets for
 linear bandits and linear mixture mdp. Advances in Neural Information Processing Systems 34
 4342–4355.
- ZHAO, H., HE, J., ZHOU, D., ZHANG, T. and GU, Q. (2023). Variance-dependent regret bounds
 for linear bandits and reinforcement learning: Adaptivity and computational efficiency. In *The Thirty Sixth Annual Conference on Learning Theory*. PMLR.
- ZHAO, H., ZHOU, D., HE, J. and GU, Q. (2022). Bandit learning with general function classes:
 Heteroscedastic noise and variance-dependent regret bounds. *arXiv preprint arXiv:2202.13603*.
- ZHOU, D. and Gu, Q. (2022). Computationally efficient horizon-free reinforcement learning for linear mixture mdps. *Advances in neural information processing systems* **35** 36337–36349.
- ZHOU, D., GU, Q. and SZEPESVARI, C. (2021). Nearly minimax optimal reinforcement learning for linear mixture markov decision processes. In *Conference on Learning Theory*. PMLR.

445 A Proof of Theorem 5.2

In this section, we prove the variance-dependent lower bound for adaptive variance sequences established in Theorem 5.2. We begin with the instance construction from Lemma 4.3 and establish the following constant-probability lower bound for the regret:

Lemma A.1. For a fixed variance threshold σ , number of rounds $K \geq 1.5d^2$, and any bandit algorithm Alg, for the instance constructed in Lemma 4.3, with probability at least $\Omega(1/\log(dK))$, the regret is lower bounded by

$$\operatorname{Regret}(K) \ge \frac{d\sqrt{K\sigma^2}}{16\sqrt{6}}.$$

Based on the constant-probability lower bound, we boost this probability by creating L = $\Omega(\log^2(dK))$ independent instances with dimension d'=d/L and number of rounds K'=K/L, 453 where each instance follows the structure in Lemma 4.3 with i.i.d. sampled weight vectors. Un-454 der this construction, the total dimension of all instances is d, which can be represented as a d-455 dimensional linear contextual bandit through orthogonal embedding, similar to our previous con-456 struction: for instance i, we augment its actions by padding zeros in dimensions reserved for other 457 instances, ensuring actions from different instances only interact with their corresponding param-458 eters. Here, we consider the case where the learner visits the instances in a cyclic manner and 459 establish the following high-probability regret lower bound for the constructed instance: 460

Lemma A.2. For a fixed variance threshold σ , number of rounds $K \geq 1.5d^2$, and any bandit algorithm Alg, with probability at least $\Omega(1/\log(dK))$, the regret is lower bounded by

$$\operatorname{Regret}(K) \ge \Omega(d\sqrt{K\sigma^2}/\log^3(dK)).$$

With the help of this high-probability lower regret bound from Lemma A.2, we begin the proof of Theorem 5.2. Following a similar framework to the fixed-variance case, we first divide the rounds into groups based on their variance magnitude. Specifically, for any variance sequence $\{\sigma_1,\ldots,\sigma_K\}$, we partition the rounds into $L=\lceil\log_2K\rceil+1$ groups as follows:

$$\mathcal{K}_0 = \{k : \sigma_k \le 1/K\},\$$
 $\mathcal{K}_i = \{k : 2^{i-1}/K < \sigma_k \le 2^i/K\}, \text{ for } i = 1, \dots, L-1.$

To address the unknown number of rounds $K_i = |\mathcal{K}_i|$, instead of constructing a single instance \mathcal{M}_i for each group, we create L instances $\mathcal{M}_{i,j}$, where $L = \lceil \log_2 K \rceil + 1$. Each instance $\mathcal{M}_{i,j}$ is constructed according to Lemma A.2 with dimension $d' = d/L^2$, variance $\sigma(i) = 2^{i-1}/K$ and number of rounds $K' = 2^{j-1}$. For each round k in group \mathcal{K}_i , the learner receives a decision set \mathcal{D}_i from one of the instances in $\{\mathcal{M}_{i,1},\ldots,\mathcal{M}_{i,L}\}$ in a cyclic manner.

Proof of Theorem 5.2. According to Lemma A.2, for each instance $\mathcal{M}_{i,j}$, with probability at least $1 - 1/K^3$, the regret in the first 2^{j-1} visits is lower bounded by

Regret
$$(2^{j-1}, \mathcal{M}_{i,j}) \ge \mathbb{I}(2^{j-1} \ge 1.5d'^2) \cdot \Omega(d'\sqrt{2^{j-1}\sigma^2(i)}/\log^3(d'K')),$$
 (A.1)

where the indicator reflects the requirement that $K'=2^{j-1}\geq 1.5d'^2$. For simplicity, we define \mathcal{E} as the event that (A.1) holds for all instances $\mathcal{M}_{i,j}$. By union bound, we have $\mathbb{P}(\mathcal{E})\geq 1-1/K$. Conditioned on event \mathcal{E} , for an adaptive sequence and each corresponding group \mathcal{K}_i , due to the cyclic visiting pattern, each instance $\mathcal{M}_{i,j}$ is visited $|\mathcal{K}_i|/L$ times. There exists an instance $\mathcal{M}_{i,j}$ with matching interval for the round number, i.e., $2^{j-1}\leq |\mathcal{K}_i|/L\leq 2^j$. Therefore, we have

$$\begin{split} & \sum_{k \in \mathcal{K}_i} \max_{\mathbf{x} \in \mathcal{D}_k} \langle \boldsymbol{\mu}_i, \mathbf{x} \rangle - \langle \boldsymbol{\mu}_i, \mathbf{x}_k \rangle \\ & \geq \operatorname{Regret}(2^{j-1}, \mathcal{M}_{i,j}) \\ & \geq \mathbb{I}(2^{j-1} \geq 1.5 d'^2) \cdot \Omega(d\sqrt{2^{j-1}\sigma^2(i)}/\log^3(d'K')) \\ & \geq \mathbb{I}(K_i \geq 3 d'^2 L) \cdot \Omega(d\sqrt{K_i\sigma^2(i)}/\log^4(dK)) \\ & \geq \Omega\Big(d'\sqrt{K_i\sigma^2(i)}/\log^3(dK) - d'\sqrt{3d'^2L\sigma^2(i)}/\log^4(dK)\Big) \end{split}$$

$$\geq \Omega \left(d' \sqrt{\sum_{k \in \mathcal{K}_i} \sigma_k^2 / \log^4(dK)} - \sqrt{3L} d'^2 \cdot \sigma(i) / \log^4(dK) \right), \tag{A.2}$$

where the first inequality follows from $2^{j-1} \leq |\mathcal{K}_i|/L \leq 2^j$, the second inequality holds by the definition of event \mathcal{E} , the third inequality follows from $2^{j-1} \leq |\mathcal{K}_i|/L \leq 2^j$, the fourth inequality holds due to $\mathbb{I}(x \geq y)\sqrt{x} \geq \sqrt{x} - \sqrt{y}$, and the last inequality follows from the definition of group \mathcal{K}_i .

Taking a summation of (A.2) over all groups, the total regret can be lower bounded as follows:

$$= \sum_{i=0}^{L-1} \sum_{k \in \mathcal{K}_i} \max_{\mathbf{x} \in \mathcal{D}_k} \langle \boldsymbol{\mu}_i, \mathbf{x} \rangle - \langle \boldsymbol{\mu}_i, \mathbf{x}_k \rangle$$

$$\geq \sum_{i=1}^{L-1} \Omega \left(d' \sqrt{\sum_{k \in \mathcal{K}_i} \sigma_k^2 / \log^4(dK)} - \sqrt{3L} d'^2 \cdot \sigma(i) / \log^4(dK) \right)$$

$$\geq \Omega \left(\sum_{i=1}^{L-1} d/L^2 \cdot \sqrt{\sum_{k \in \mathcal{K}_i} \sigma_k^2 / \log^4(dK)} - 2\sqrt{3L} d^2 / (L^4 \log^4(dK)) \right)$$

$$\geq \Omega \left(d/L^2 \cdot \sqrt{\sum_{i=1}^{L-1} \sum_{k \in \mathcal{K}_i} \sigma_k^2 / \log^4(dK)} - 2\sqrt{3L} d^2 / (L^4 \log^4(dK)) \right), \tag{A.3}$$

where the first inequality follows from (A.2), the second inequality follows from the definition of variance threshold $\sigma(i)$ and dimension $d'=d/L^2$, and the last inequality holds due to $\sum_i \sqrt{x_i} \ge \sqrt{\sum_i x_i}$. In addition, for the group \mathcal{K}_0 , we have

$$\sum_{k \in \mathcal{K}_0} \sigma_k^2 \le \sum_{k \in \mathcal{K}_0} 1/K \le 1,\tag{A.4}$$

where the first inequality follows from the definition of group \mathcal{K}_0 and the second inequality follows from $|\mathcal{K}_0| \leq K$. Therefore, we have

Regret(K)

$$\geq \Omega \left(d/L^2 \cdot \sqrt{\sum_{i=1}^{L-1} \sum_{k \in \mathcal{K}_i} \sigma_k^2 / \log^4(dK)} - 2\sqrt{3L} d^2 / (L^4 \log^4(dK)) \right)$$

$$\geq \Omega \left(d/L^2 \cdot \sqrt{\sum_{i=1}^{L-1} \sum_{k \in \mathcal{K}_i} \sigma_k^2 - 1 / \log^4(dK)} - 2\sqrt{3L} d^2 / (L^4 \log^4(dK)) \right)$$

$$\geq \Omega \left(d \cdot \sqrt{\sum_{i=1}^{L-1} \sum_{k \in \mathcal{K}_i} \sigma_k^2 / \log^6(dK)} \right),$$

where the first inequality follows from (A.3), the second inequality follows from (A.4), and the last inequality follows from the fact that $\sum_{k=1}^K \sigma_k^2 \geq \Omega(d^2)$. Thus, we complete the proof of Theorem 5.2.

B Proof of Theorem 5.4

492

494

495

496

497

498

493 In this subsection, we provide the proof of Theorem 5.4. We begin by describing a simple algorithm:

- 1. The learner maintains an explored action set A, which is initialized as empty.
- 2. For each decision set \mathcal{D}_k in round k, if there exists an action \mathbf{x}_k not in the spanning space of the explored action set \mathcal{A} , the learner:
 - Selects an action \mathbf{x}_k and receives reward r_k ;
- Updates the explored set: $A = A \cup \{(\mathbf{x}_k, r_k)\}.$

- 3. Otherwise, when all actions lie in the spanning space of A, the learner:
 - Estimates the reward for each action through linear combinations of $(\mathbf{x}, r) \in \mathcal{A}$;
 - Selects the action with maximum estimated reward.

It is worth noting that this algorithm assumes the received rewards r_k have no noise to provide accurate estimates in step 3. While this assumption does not hold in general, when an adversary can choose the variance σ_k based on the decision set \mathcal{D}_k , they can cooperate with the learner by setting:

- $\sigma_k = 0$ when step 2 is triggered (exploration);
- $\sigma_k = 1$ when step 3 is triggered (exploitation).

For a d-dimensional linear bandit problem, the explored action set satisfies $|\mathcal{A}| \leq d$. This implies the learner performs at most d exploration steps with zero variance, while all remaining steps have variance one. Under this construction, the regret in the first K rounds is upper bounded by:

$$\operatorname{Regret}_{\operatorname{Alg}}(K) \leq d,$$

where the total variance $\sum_{k=1}^K \sigma_k^2 = K - d \ge K/2$ (since $K \ge 2d$). Thus, through this cooperation between the adversary and learner, the $\widetilde{\Omega}(d\sqrt{\sum_{k=1}^K \sigma_k^2})$ lower bound is broken, completing the proof of Theorem 5.4.

C Proof of Key Lemmas

C.1 Proof of Lemma 4.3

499

500 501

502

503 504

505

506

513

514

526 527

528

529

530

In this subsection, we provide the proof of Lemma 4.3. When the variance threshold $\sigma=1$, our construction reduces to the standard lower bound instances for linear contextual bandits (Zhou et al., 2021). Specifically, when the number of rounds K satisfying $K \geq 1.5 \cdot d^2$, Zhou et al. (2021) provided the following variance-independent lower bound for these hard instances:

Lemma C.1 (Lemma C.8, Zhou et al. 2021). For any bandit algorithm Alg, if the weight vector $\mu \in \{-\Delta, \Delta\}^d$ is drawn uniformly at random from $\{-\Delta, \Delta\}^d$, then the expected regret over K rounds is lower bounded by:

$$\mathbb{E}_{\mu}[\operatorname{Regret}(K)] \geq \frac{d\sqrt{K}}{8\sqrt{6}}.$$

With the help of Lemma C.1, we start the proof of Lemma 4.3.

Proof of Lemma 4.3. For any algorithm Alg for linear contextual bandit with fixed variance threshold σ , we construct an auxiliary algorithm Alg1 to solve the standard linear contextual bandit problem:

- At the beginning of each round $k \in K$, Alg1 observes the decision set \mathcal{D}_k and sends it to Alg;
- Alg selects action $a_k \in \mathcal{D}_k$ based on the historical observations and delivers it to Alg1;
- Alg1 performs the action a_k , receives the reward r_k and sends the normalized reward $\sigma \cdot r_k$ to Alg.

Now, we consider the performance of auxiliary algorithm Alg1 for the standard linear contextual bandit problem. It is worth noticing that the reward/noise in bandit instances for algorithm Alg1 and algorithm Alg only differ by a scalar factor σ , therefore for each instance, we have

$$\mathbb{E}[\operatorname{Regret}_{\operatorname{Alg}}(K)] = \sigma \cdot \mathbb{E}[\operatorname{Regret}_{\operatorname{Alg1}}(K)]. \tag{C.1}$$

If we randomly select a weight parameter vector $\mu \in \{-\Delta, \Delta\}^d$, then according to Lemma C.1, the regret for Alg is lower bounded by

$$\mathbb{E}_{\boldsymbol{\mu}}[\operatorname{Regret}_{\operatorname{Alg}}(K)] = \sigma \cdot \mathbb{E}_{\boldsymbol{\mu}}[\operatorname{Regret}_{\operatorname{Alg1}}(K)] \ge \sigma \cdot \frac{d\sqrt{K}}{8\sqrt{6}} = \frac{d\sqrt{K}\sigma^2}{8\sqrt{6}},$$

where the equation holds due to (C.1) and the inequality holds due to Lemma C.1. Thus, we complete the proof of Lemma 4.3.

C.2 Proof of Lemma A.1

538

552

553

554

555

556

557

558

559

560

561

562

In this subsection, we provide the proof of Lemma A.1. We begin by recalling the OFUL algorithm in Abbasi-Yadkori et al. (2011a) and its corresponding upper bound for the regret:

Lemma C.2 (Theorem 3 in Abbasi-Yadkori et al. 2011a). For any linear contextual bandit problem, with probability at least $1-\delta$, the regret for OFUL algorithm in the first K rounds is upper bounded by $\operatorname{Regret}(K) \leq \widetilde{O}(d\sqrt{K\log(dK/\delta)})$.

It is worth noting that the reward/noise in the instance construction from Lemma 4.3 only differs by a scalar factor σ from the standard bandit. Therefore, as discussed in Section C.1, the regret in these two cases also only differs by a scalar factor σ . This leads to the following corollary:

Corollary C.3. For the instance construction from Lemma 4.3, there exists a constant C such that with probability at least $1-\delta$, the regret for OFUL algorithm in the first K rounds is upper bounded by $\operatorname{Regret}(K) \leq C d \sqrt{K \sigma^2 \log(dK/\delta)}$.

550 With the help of Corollary C.3, we can begin the proof of Lemma A.1.

551 Proof of Lemma A.1. For any algorithm Alg, we construct an auxiliary algorithm Alg1 as follows:

- At the beginning of each round k ∈ [K], Alg1 observes the decision set D_k and sends it to Alg;
- Alg selects action $a_k \in \mathcal{D}_k$ based on the historical observations and delivers it to Alg1;
- Alg1 performs the action a_k and receives the reward r_k ;
- Alg1 calculates the pseudo regret as:

Regret'(k) =
$$\sum_{i=1}^{k} \frac{1}{3} + \frac{d}{\sqrt{96K}} - r_k$$
.

If the pseudo regret is larger than $d\sqrt{K\sigma^2}/(8\sqrt{6}) + \sigma\sqrt{2K\log(2K/\delta)}$, Alg1 removes all previous information and performs the OFUL algorithm in all future rounds.

Based on the construction of the instances, whatever the weight vector μ is, the optimal action is to select an action in the same direction as the weight vector, obtaining an expected reward of $1/3 + d/\sqrt{96K}$. Under this scenario, with probability at least $1 - \delta$, for any round $k \in [K]$, the difference between pseudo regret Regret'(k) and true regret Regret(k) can be upper bounded by

$$\left| \operatorname{Regret}(k) - \operatorname{Regret}'(k) \right| = \left| \sum_{i=1}^{k} \epsilon_i \right| \le \sigma \sqrt{2K \log(2K/\delta)},$$
 (C.2)

where the inequality holds due to Lemma D.1 with the fact that the noise satisfies $\mathbb{E}[\epsilon_k|a_{1:k},r_{1:k-1}]=0$ and $|\epsilon_k|\leq \sigma$. Thus, according to the criterion of auxiliary algorithm Alg1, with probability at least $1-\delta$, the regret of Alg1 before transitioning to OFUL is up to $d\sqrt{K\sigma^2}/(8\sqrt{6})+2\sigma\sqrt{2K\log(2K/\delta)}$. On the other hand, for the stage after transitioning to OFUL, Corollary C.3 suggests that with probability at least $1-\delta$, the regret is no more than $Cd\sqrt{K\sigma^2\log(dK/\delta)}$. Therefore, with a selection of $\delta=1/K$, we have

$$\mathbb{P}\big[\mathrm{Regret}_{\mathrm{Alg}_1}(K) \geq C d \sqrt{K \sigma^2 \log(dK^2)} + d \sqrt{K \sigma^2} / (8\sqrt{6}) + 2\sigma \sqrt{2K \log(2K^2)}\big] \leq 2/K. \tag{C.3}$$

For simplicity, let $R = Cd\sqrt{K\sigma^2\log(dK^2)} + d\sqrt{K\sigma^2}/(8\sqrt{6}) + 2\sigma\sqrt{2K\log(2K^2)}$ and we have $\mathbb{E}_{\mu}[\operatorname{Regret}_{\operatorname{Alg}_{+}}(K)]$

$$\begin{split} & \leq \mathbb{P} \big[\mathrm{Regret}_{\mathrm{Alg}_1}(K) \geq R \big] \cdot K \sigma + \mathbb{P} \big[\mathrm{Regret}_{\mathrm{Alg}_1}(K) \geq d \sqrt{K \sigma^2} / (16 \sqrt{6}) \big] \cdot R \\ & + \mathbb{P} \big[\mathrm{Regret}_{\mathrm{Alg}_1}(K) \geq 0 \big] \cdot d \sqrt{K \sigma^2} / (16 \sqrt{6}) \end{split}$$

$$\leq 2\sigma + \mathbb{P}\big[\mathrm{Regret}_{\mathrm{Alg}_1}(K) \geq d\sqrt{K\sigma^2}/(16\sqrt{6})\big] \cdot \widetilde{O}(d\sqrt{K\sigma^2\log(dK)}) + d\sqrt{K\sigma^2}/(16\sqrt{6}),$$

where the first inequality holds due to $\mathbb{E}[X] \leq \mathbb{P}(X \geq x_1) \cdot R + \mathbb{P}(X \geq x_2) \cdot x_1 + \mathbb{P}(X \geq 0) \cdot x_2$ for $0 \leq X \leq R$ and $x_1 > x_2 > 0$, and the second inequality holds due to (C.3). Combining this result with the lower bound of expected regret in Lemma 4.1, we have

$$d\sqrt{K\sigma^2}/(8\sqrt{6}) \geq 2\sigma + \mathbb{P}\big[\mathrm{Regret}_{\mathrm{Alg}_1}(K) \geq d\sqrt{K\sigma^2}/(16\sqrt{6})\big] \cdot \widetilde{O}(d\sqrt{K\sigma^2\log(dK)})$$

$$+d\sqrt{K\sigma^2}/(16\sqrt{6}),$$

573 which implies that

$$\mathbb{P}\big[\mathrm{Regret}_{\mathrm{Alg}_1}(K) \geq d\sqrt{K\sigma^2}/(16\sqrt{6})\big] \geq \Omega(1/\log(dK)). \tag{C.4}$$

In addition, according to the criterion of auxiliary algorithm Alg1 with (C.2), with probability at least $1 - \delta = 1 - 1/K$, Alg1 will not switch to the OFUL algorithm until the cumulative regret is

larger than $d\sqrt{K\sigma^2}/(8\sqrt{6})$, which implies that

$$\begin{split} \mathbb{P}\big[\mathrm{Regret}_{\mathrm{Alg}}(K) \geq d\sqrt{K\sigma^2}/(16\sqrt{6})\big] \geq \mathbb{P}\big[\mathrm{Regret}_{\mathrm{Alg}_1}(K) \geq d\sqrt{K\sigma^2}/(16\sqrt{6})\big] - 1/K \\ &= \Omega(1/\log(dK)). \end{split}$$

Thus, we complete the proof of Lemma A.1.

578 C.3 Proof of Lemma A.2

In this subsection, we provide the proof of Lemma A.2.

580 Proof of Lemma A.2. Since the learner visits the instances in a cyclic manner, over all K rounds,

each instance \mathcal{M}_i $(i=1,2,\ldots,L)$ is visited K'=K/L times. As actions from different instances

only interact with their corresponding parameters, according to Lemma A.1, for each instance \mathcal{M}_i ,

with probability at least $\Omega(1/\log(dK))$, the regret is lower bounded by

Regret
$$(K', \mathcal{M}_i) \ge \frac{d'\sqrt{K'\sigma^2}}{16\sqrt{6}} = \frac{d\sqrt{K\sigma^2}}{16\sqrt{6} \cdot L^{1.5}}.$$

Note that the weight vectors for each instance are independently sampled, hence the probability that at least one instance has regret no less than $d\sqrt{K\sigma^2}/16\sqrt{6} \cdot L^{1.5}$ is at least

$$1 - (1 - \Omega(1/\log(dK)))^{L} \ge 1 - 1/K^{3}$$
 Qingyue: ????

Under this condition, the total regret can be lower bounded as:

$$\operatorname{Regret}(K) = \sum_{i=1}^{L} \operatorname{Regret}(K', \mathcal{M}_i) \ge \frac{d\sqrt{K\sigma^2}}{16\sqrt{6} \cdot L^{0.5}}.$$
 (C.5)

Thus, we obtain a high-probability lower bound and complete the proof of Lemma A.2. \Box

588 D Auxiliary Lemmas

Lemma D.1 (Azuma–Hoeffding inequality, Cesa-Bianchi and Lugosi 2006). Let $\{\eta_k\}_{k=1}^K$ be a martingale difference sequence with respect to a filtration $\{\mathcal{G}_k\}$ satisfying $|\eta_k| \leq R$ for some constant R, η_k is \mathcal{G}_{k+1} -measurable, $\mathbb{E}\big[\eta_k|\mathcal{G}_k\big] = 0$. Then for any $0 < \delta < 1$, with high probability at least $1 - \delta$, we have

$$\sum_{k=1}^{K} \eta_k \le R\sqrt{2K\log(1/\delta)}.$$

NeurIPS Paper Checklist

1. Claims

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

Answer: [Yes]

Justification: In both abstract and introduction, we highlight the contribution in our paper. The proposed algorithm and the corresponding theoretical results are discussed in the followed sections

Guidelines:

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

2. Limitations

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

Answer: [Yes]

Justification: We explicitly list all the necessary assumptions for our theoretical analysis.

Guidelines:

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

3. Theory assumptions and proofs

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

Answer: [Yes]

Justification: The complete set of assumptions for our analysis is presented in Section 3, with the detailed proofs of all our claims provided in a later section.

Guidelines

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

4. Experimental result reproducibility

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

Answer: [NA]

Justification: The paper does not include experiments.

Guidelines:

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

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [NA]

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

715

716

717

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741 742

746

747

748

749

750

Justification: The paper does not include experiments.

Guidelines:

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

6. Experimental setting/details

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

Answer: [NA]

Justification: The paper does not include experiments.

Guidelines:

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

7. Experiment statistical significance

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

Answer: [NA]

Justification: The paper does not include experiments.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)

- The assumptions made should be given (e.g., Normally distributed errors).
 - It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
 - It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
 - If error bars are reported in tables or plots, The authors should explain in the text how
 they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

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

Answer: [NA]

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778 779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

Justification: The paper does not include experiments.

Guidelines:

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

9. Code of ethics

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

Answer: [Yes]

Justification: The research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics.

Guidelines:

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

10. Broader impacts

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

Answer: [NA]

Justification: The paper is a theoretical work with no societal impact.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

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

Answer: [NA]

Justification: The paper is a theoretical work and poses no such risks

Guidelines:

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

12. Licenses for existing assets

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

Answer: [Yes]

Justification: We have described the related works, especially those work which our work is based on with proper citations in corresponding sections.

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

 If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872 873

874

875

876

877

879

880

881 882

883

884

885

886

887

888

890

891

892

893

894

895

896

897

898 899

900

901

902

903

904

905

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

Answer: [No]

Justification: This is a theoretical paper without experiments.

Guidelines:

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

14. Crowdsourcing and research with human subjects

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

Answer: [NA]

Justification: This paper does not include crowdsourcing or human subjects.

Guidelines:

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

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

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

Answer: [NA]

Justification: This paper does not include crowdsourcing or human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions
 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
 guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [No]

910

911

912

913

914

915

916

917

Justification: We only used an LLM to rephrase the writing, which did not affect the core methodology, scientific rigor, or originality of the research.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.