
Linear Contextual Bandits with Adversarial Corruptions

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study the linear contextual bandits problem in the presence of adversarial
2 corruption, where the interaction between the player and a possibly infinite decision
3 set is contaminated by an adversary that can corrupt the reward up to a corruption
4 level C measured by the sum of the largest alteration on rewards in each round.
5 We present a variance-aware algorithm that is adaptive to the level of adversarial
6 contamination C . The key algorithmic design includes (1) a multi-level partition
7 scheme of the observed data, (2) a cascade of confidence sets that are adaptive to
8 the level of the corruption, and (3) a variance-aware confidence set construction
9 that can take advantage of low-variance reward. We further prove that the regret
10 of the proposed algorithm is $\tilde{O}(C^2 d \sqrt{\sum_{t=1}^T \sigma_t^2} + C^2 \sqrt{dT} + CR\sqrt{dT})$, where
11 d is the dimension of context vectors, T is the number of rounds, R is the range
12 of noise and $\sigma_t^2, t = 1 \dots, T$ are the variances of instantaneous reward. We also
13 prove a gap-dependent regret bound for the proposed algorithm, which is instance-
14 dependent and thus leads to better performance on good practical instances. To the
15 best of our knowledge, this is the first variance-aware corruption robust algorithm
16 for contextual bandits.

17 1 Introduction

18 Multi-armed bandits algorithms are widely applied in online advertising (Li et al., 2010), clinical
19 trials (Villar et al., 2015), recommendation system (Deshpande and Montanari, 2012) and many other
20 real-world tasks. In the model of multi-armed bandits, the algorithm needs to decide which action
21 (or arm) to take (or pull) at each round and receive a reward for the chosen action. In the stochastic
22 setting, the reward is subject to a fixed but unknown distribution for each action. In reality, however,
23 these rewards can easily be “corrupted” by some malicious users. A typical example is click fraud
24 (Lykouris et al., 2018), where botnets simulate the legitimate users clicking on an ad to fool the
25 recommendation systems. This motivates the studies of the bandits algorithms that are robust to
26 adversarial corruptions.

27 For example, Lykouris et al. (2018) introduced a bandit model in which an adversary could corrupt
28 the stochastic reward generated by an arm pull. They proposed an algorithm and show that the
29 regret of this “middle ground” scenario degrades smoothly with the amount of corruption injected
30 by the adversary. Gupta et al. (2019) proposed an alternative algorithm which gives a significant
31 improvement in regret.

32 While the algorithms that are robust to the corruptions have been studied in the setting of multi-armed
33 bandits in a number of prior works, they are still understudied in the setting of linear contextual
34 bandits. The linear contextual bandits problem can be regarded as an extension of the multi-armed
35 bandit problem to linear optimization, in order to tackle an unfixed and possibly infinite set of feasible

36 actions. There is a large body of literature on efficient algorithms for linear contextual bandits with
 37 no corruptions (Abe et al., 2003; Auer, 2002; Chu et al., 2011; Dani et al., 2008; Rusmevichientong
 38 and Tsitsiklis, 2010; Abbasi-Yadkori et al., 2011; Li et al., 2019b), to mention a few. The significance
 39 of this setting lies in the fact that linear regression approaches are widely used in recommendation
 40 systems and advertising (Li et al., 2010; Jhalani et al., 2016; Deshpande and Montanari, 2012). Linear
 41 contextual bandits with adversarial corruptions is an arguably more challenging setting since most of
 42 the previous corruption-robust algorithms are based on the idea of action elimination (Lykouris et al.,
 43 2018; Gupta et al., 2019; Bogunovic et al., 2021), which is not applicable to the contextual bandits
 44 settings where the decision set is time varying and possibly infinite at each round. In Garcelon et al.
 45 (2020), it is shown that a malicious agent can force a linear contextual bandit algorithm to take any
 46 desired action $T - o(T)$ times over T rounds, while applying adversarial corruptions to rewards with
 47 a cumulative cost that only grow logarithmically. This poses a big challenge for designing corruption
 48 robust algorithms for linear contextual bandits.

49 In this paper, we make a first attempt to study a linear contextual bandit model where an adversary
 50 can corrupt the rewards up to a corruption level C , which is defined as the the sum of biggest
 51 alteration the adversary made on rewards in each round. We propose a linear contextual bandits
 52 algorithm that is robust to reward corruption, dubbed multi-level optimism-in-the-face-of-uncertainty
 53 weighted learning (Multi-level OFUL). More specifically, our algorithm consists of the following
 54 novel techniques: (1) We design a multi-level partition scheme and adopt the idea of *sub-sampling* to
 55 do the robust estimation of the model parameters; (2) We maintain a cascade of *candidate confidence*
 56 *sets* corresponding to different corruption level (which is unknown) and randomly select a confidence
 57 set at each round to take the action; and (3) We design confidence sets that depend on the *variances*
 58 *of rewards*, which lead to a potentially tighter regret bound.

59 Our contributions are summarized as follows:

- 60 • We propose a variance-aware algorithm which is adaptive to the amount of adversarial corruptions
 61 C . To the best of our knowledge, it is the first algorithm for the setting of linear contextual bandits
 62 with adversarial corruptions which does not rely on the finite number of actions and other additional
 63 assumptions.
- 64 • We prove that the regret of our algorithm is in $\tilde{O}\left(C^2 d \sqrt{\sum_{t=1}^T \sigma_t^2} + C^2 \sqrt{dT} + CR\sqrt{dT}\right)$, where
 65 d is the dimension of context vectors, T is the number of rounds, R is the range of noise and $\sigma_t^2, t =$
 66 $1 \dots, T$ are the variances of instantaneous reward. Our regret upper bound has a multiplicative
 67 dependence on C^2 which indicates that our algorithm achieves a sub-linear regret when the
 68 corruption level satisfies $C = o(T^{1/4})$.
- 69 • We also derive a gap-dependent regret bound $\tilde{O}\left(\frac{1}{\Delta} \cdot C^2 R^2 d + \frac{1}{\Delta} \cdot d^2 C^2 \max_{t \in [T]} \sigma_t^2\right)$ for our
 70 proposed algorithm, which is instance-dependent and thus leads to a better performance on good
 71 practical instances.

72 **Notation.** We use lower case letters to denote scalars, and use lower and upper case bold face letters
 73 to denote vectors and matrices respectively. We denote by $[n]$ the set $\{1, \dots, n\}$. For a vector $\mathbf{x} \in \mathbb{R}^d$
 74 and matrix $\Sigma \in \mathbb{R}^{d \times d}$, a positive semi-definite matrix, we denote by $\|\mathbf{x}\|_2$ the vector’s Euclidean
 75 norm and define $\|\mathbf{x}\|_{\Sigma} = \sqrt{\mathbf{x}^{\top} \Sigma \mathbf{x}}$. For two positive sequences $\{a_n\}$ and $\{b_n\}$ with $n = 1, 2, \dots$,
 76 we write $a_n = O(b_n)$ if there exists an absolute constant $C > 0$ such that $a_n \leq C b_n$ holds for all
 77 $n \geq 1$ and write $a_n = \Omega(b_n)$ if there exists an absolute constant $C > 0$ such that $a_n \geq C b_n$ holds
 78 for all $n \geq 1$. We use $\tilde{O}(\cdot)$ to further hide the polylogarithmic factors. We use $\mathbb{1}(\cdot)$ to denote the
 79 indicator function.

80 2 Related Work

81 **Bandits with Adversarial Attacks:** There is a large body of literature on the problems of multi-
 82 armed bandits with adversarial corruptions. Most research in this area aims to design algorithms that
 83 achieve desirable regret bound in both stochastic multi-armed bandits and adversarial bandits, known
 84 as “the best of both worlds” guarantees (Bubeck and Slivkins, 2012; Seldin and Slivkins, 2014; Auer
 85 and Chiang, 2016; Seldin and Lugosi, 2017; Zimmert and Seldin, 2019). These works mainly focus
 86 on achieving bounds in the worst case and the case where there is no adversary. As a result, these

87 algorithms are either not robust to instances that moderate amount of corruptions occur, or suffer
 88 from restrictive assumptions on adversarial corruptions. Distinctive from the above line of research,
 89 [Lykouris et al. \(2018\)](#) focus on a variant of classic multi-armed bandit model in which each pull of an
 90 arm generates a stochastic reward that may be contaminated by an adversary before it is revealed
 91 to the player. In their work, the corruption level C is defined as $C = \sum_t \max_a |r^t(a) - r_S^t(a)|$
 92 where $r_S^t(a)$ is the stochastic reward of arm a and $r^t(a)$ is the corrupted reward of arm a at round t .
 93 They develop algorithms adaptive to the unknown corruption level, which achieves an $O(K^{1.5}C\sqrt{T})$
 94 regret bound. [Gupta et al. \(2019\)](#) proposed an improved algorithm that can achieve a regret bound
 95 with only additive dependence on C .

96 On the other hand, many research efforts have also been devoted into designing adversarial attacks
 97 that cause standard algorithms to fail ([Jun et al., 2018](#); [Liu and Shroff, 2019](#); [Gupta et al., 2019](#);
 98 [Garcelon et al., 2020](#)).

99 **Linear Bandits with Corruptions:** [Li et al. \(2019a\)](#) studied stochastic linear bandits with adversarial
 100 corruptions and achieved $\tilde{O}(\frac{1}{\Delta} \cdot d^{5/2}C + \frac{1}{\Delta^2} \cdot d^6)$ regret bound where d is the dimension of the context
 101 vectors, Δ is the gap between the rewards of the best and the second best action in the decision
 102 set \mathcal{D} . The distinction between [Li et al. \(2019a\)](#) and our work is that [Li et al. \(2019a\)](#) considers a
 103 fixed decision set \mathcal{D} throughout all T rounds, while we consider contextual bandits with changing
 104 decision set observed before each round. [Bogunovic et al. \(2021\)](#) also studied linear bandits with
 105 adversarial corruptions and considered the setting under the assumption that context vectors undergo
 106 small random perturbations, which is previously introduced by [Kannan et al. \(2018\)](#). Aside from
 107 the additional assumption, another major distinction in [Bogunovic et al. \(2021\)](#) is that the number
 108 of actions k is finite and the regret bound depends on k in the contextual setting with unknown
 109 corruption level C . Recently, [Lee et al. \(2021\)](#) considered corrupted linear bandits with a finite and
 110 fixed decision set and achieve an instance-independent regret of $\tilde{O}(d\sqrt{T} + C)$. Though both their
 111 work and the work by [Li et al. \(2019a\)](#) focus on corrupted linear stochastic bandits, [Lee et al. \(2021\)](#)
 112 have a slightly different definition of regret and adopt a strong assumption on corruptions that in
 113 each round t , the corruptions on rewards are linear in the actions. [Neu and Olkhovskaya \(2020\)](#)
 114 studied linear contextual bandits with a finite decision set (i.e., K actions) and an adversary. Unlike
 115 our model, they assume that the adversary can add an arbitrary noise to the loss under a limited
 116 amount ϵ and prove an $\tilde{O}((Kd)^{\frac{1}{3}}T^{\frac{2}{3}}) + \epsilon \cdot \sqrt{dT}$ regret bound for their proposed algorithm. [Kapoor
 117 et al. \(2019\)](#) considered the corrupted linear contextual bandits setting under a strong assumption on
 118 corruptions that for any prefix, at most an η fraction of the rounds are corrupted.

119 3 Preliminaries

120 In this paper, we study linear contextual bandits with adversarial corruptions. We will introduce our
 121 model and some basic concepts in this section.

122 **Corrupted linear contextual bandits.** We consider the the linear contextual bandits model studied
 123 in [Abbasi-Yadkori et al. \(2011\)](#) under the same corruption studied by [Lykouris et al. \(2018\)](#). In detail,
 124 distinctive from the linear contextual bandits [Abbasi-Yadkori et al. \(2011\)](#), the interaction between
 125 the agent and the environment is now contaminated by an adversary. The protocol between the agent
 126 and the adversary at each round $t \in [T]$ can be described as follows:

- 127 1. At the beginning of round t , the environment generates an arbitrary decision set $\mathcal{D}_t \subseteq \mathbb{R}^d$ where
 128 each element represents a feasible action that can be selected by the agent.
- 129 2. The environment generates stochastic reward function $r'_t(\mathbf{a}) = \langle \mathbf{a}, \boldsymbol{\mu}^* \rangle + \epsilon_t(\mathbf{a})$ together with an
 130 upper bound on the standard variance of $\epsilon_t(\mathbf{a})$, i.e., $\sigma_t(\mathbf{a})$ for all $\mathbf{a} \in \mathcal{D}_t$.
- 131 3. The adversary observes $\mathcal{D}_t, r'_t(\mathbf{a}), \sigma_t(\mathbf{a})$ for all $\mathbf{a} \in \mathcal{D}_t$ and decides a corrupted reward function
 132 r_t defined over \mathcal{D}_t .
- 133 4. The agent observes \mathcal{D}_t and selects $\mathbf{a}_t \in \mathcal{D}_t$.
- 134 5. The adversary observes \mathbf{a}_t and then returns $r_t(\mathbf{a}_t)$ and $\sigma_t(\mathbf{a}_t)$.
- 135 6. The agent observes $r_t(\mathbf{a}_t), \sigma_t(\mathbf{a}_t)$.

136 Let \mathcal{F}_t be the σ -algebra generated by $\mathcal{D}_{1:t}, \mathbf{a}_{1:t-1}, \epsilon_{1:t-1}, r_{1:t-1}$ and $\sigma_{1:t-1}$.

137 At step 2, $\boldsymbol{\mu}^*$ is a hidden vector unknown to the agent which can be observed by the adversary at the
 138 beginning. We assume that for all $t \geq 1$ and all $\mathbf{a} \in \mathcal{D}_t$, $\|\mathbf{a}\|_2 \leq A$, $|\langle \mathbf{a}, \boldsymbol{\mu}^* \rangle| \leq 1$ and $\|\boldsymbol{\mu}^*\|_2 \leq B$
 139 almost surely. $\epsilon_t(\mathbf{a})$ can be any form of random noise as long as it satisfies

$$\forall t \geq 1, \forall \mathbf{a} \in \mathcal{D}_t, |\epsilon_t(\mathbf{a})| \leq R, \quad \mathbb{E}[\epsilon_t(\mathbf{a})|\mathcal{F}_t] = 0, \quad \mathbb{E}[\epsilon_t^2(\mathbf{a})|\mathcal{F}_t] \leq \sigma_t^2(\mathbf{a}). \quad (3.1)$$

140 This assumption on ϵ_t is a variant of that in Zhou et al. (2020): We now require the noise to be
 141 generated for all $\mathbf{a} \in \mathcal{D}_t$ in advance before the adversary decides the corrupted reward function. Our
 142 assumption on noises is more general than those in (Li et al., 2019a; Bogunovic et al., 2021; Kapoor
 143 et al., 2019) where they are assumed to be 1-sub-Gaussian or Gaussian.

144 At step 3, we assume that the adversary has observed all the previous information and thus may
 145 predict which policy the agent will take at the current round. However, since the agent can take a
 146 randomized policy, the adversary may not know exactly which action the agent will take.

147 **Corruption level.** We define corruption level

$$C = \frac{1}{R+1} \sum_{t=1}^T \sup_{\mathbf{a} \in \mathcal{D}_t} |r'_t(\mathbf{a}) - r_t(\mathbf{a})|. \quad (3.2)$$

148 to indicate the level of adversarial contamination. We say a model is C -corrupted if the corruption
 149 level is no larger than C .

150 Our definition of corruption level is equivalent to the counterpart in Lykouris et al. (2018) and Gupta
 151 et al. (2019) where they define $C = \sum_{t=1}^T \max_{\mathbf{a}} |r'_t(\mathbf{a}) - r_t(\mathbf{a})|$ in our notation of rewards. We
 152 introduce a factor of $\frac{1}{R+1}$ since the noise is of range R in our model, while they assume all the
 153 rewards are in range $[0, 1]$.

154 **Regret.** Since the actions selected by the agent may not be deterministic, we define the regret for this
 155 model as follows:

$$\mathbf{Regret}(T) = \sum_{t=1}^T \langle \mathbf{a}_t^*, \boldsymbol{\mu}^* \rangle - \mathbb{E} \left[\sum_{t=1}^T \langle \mathbf{a}_t, \boldsymbol{\mu}^* \rangle \right]. \quad (3.3)$$

156 Our definition follows from the definition in Gupta et al. (2019) where the standard metric in stochastic
 157 multi-armed bandit models of pseudo-regret is adopted. But note that we need to take the expectation
 158 on $\sum_{t=1}^T r'_t(\mathbf{a}_t)$ (the second term in (3.3)), since a randomized policy is applied in each round.

159 **Gap.** Let Δ_t be the gap between the rewards of the best and the second best action in the decision set
 160 \mathcal{D}_t as defined in Dani et al. (2008) which can be formally written as

$$\Delta_t = \min_{\mathbf{a} \in \mathcal{D}_t, \mathbf{a} \notin \mathcal{A}_t^*} (\langle \mathbf{a}_t^*, \boldsymbol{\mu}^* \rangle - \langle \mathbf{a}, \boldsymbol{\mu}^* \rangle). \quad (3.4)$$

161 where $\mathcal{A}_t^* = \operatorname{argmax}_{\mathbf{a} \in \mathcal{D}_t} \langle \mathbf{a}, \boldsymbol{\mu}^* \rangle$ and \mathbf{a}_t^* is an arbitrary element in \mathcal{A}_t^* . Let Δ denotes the smallest
 162 gap $\min_{t \in [T]} \Delta_t$.

163 4 The Proposed Algorithm

164 In this section, we propose a variance-aware algorithm, Multi-level OFUL, in Algorithm 1, to tackle
 165 the corrupted linear contextual bandits problem. At the core of our algorithm is an action partition
 166 scheme to group historical selected actions and use them to select the future actions in different
 167 groups with different probabilities. Such a scheme is introduced to deal with the unknown corruption
 168 level. For simplicity, we denote $r_t(\mathbf{a}_t), \sigma_t(\mathbf{a}_t)$ in Section 3 by r_t, σ_t in our algorithm.

169 **Main difficulty in our setting.** We begin with the main difficulty that prevents us from applying
 170 existing algorithms to our setting. Consider a simpler setting where the agent knows the corruption
 171 level C in prior, and we have $\sigma_t = R$ for all t . Then we can apply OFUL (Abbasi-Yadkori et al.,
 172 2011) to solve our problem. In detail, in each round we estimate $\boldsymbol{\mu}^*$ by $\boldsymbol{\mu}_t$, which is the minimizer of
 173 the following ridge regression problem:

$$\boldsymbol{\mu}_t = \operatorname{argmin}_{\boldsymbol{\mu} \in \mathbb{R}^d} \lambda \|\boldsymbol{\mu}\|_2^2 + \sum_{i=1}^{t-1} [(\langle \boldsymbol{\mu}, \mathbf{a}_i \rangle - r_i)]^2. \quad (4.1)$$

Algorithm 1 Multi-level OFUL

```

1: Set the largest level of confidence sets:  $\ell_{\max} \leftarrow \lceil \log_2 2T \rceil$ .
2: For  $\ell \in [\ell_{\max}]$ , set  $\Sigma_{1,\ell} \leftarrow \lambda \mathbf{I}$ ,  $\boldsymbol{\mu}_{1,\ell} \leftarrow \mathbf{0}$ ,  $\mathbf{c}_{1,\ell} \leftarrow \mathbf{0}$ .
3: Set  $\Sigma_1 \leftarrow \lambda \mathbf{I}$ ,  $\boldsymbol{\mu}_1 \leftarrow \mathbf{0}$ ,  $\mathbf{c}_1 \leftarrow \mathbf{0}$ .
4: for  $t = 1, \dots, T$  do
5:   Observe  $\mathcal{D}_t$ .
6:   for  $\ell = 1, \dots, \ell_{\max}$  do
7:     Set  $\beta_{t,\ell}$  and  $\gamma_{t,\ell}$  as defined in (4.5) and (4.6).
8:      $\mathcal{C}'_{t,\ell} \leftarrow \{\boldsymbol{\mu} \mid \|\boldsymbol{\mu} - \boldsymbol{\mu}_t\|_{\Sigma_t} \leq \beta_{t,\ell}\} \cap \{\boldsymbol{\mu} \mid \|\boldsymbol{\mu} - \boldsymbol{\mu}_{t,\ell}\|_{\Sigma_{t,\ell}} \leq \gamma_{t,\ell}\}$ .
9:      $\mathcal{C}_{t,\ell} \leftarrow \begin{cases} \mathcal{C}'_{t,\ell}, & \mathcal{C}'_{t,\ell} \neq \emptyset \\ \mathcal{C}_{t,\ell+1}, & \text{otherwise} \end{cases}$ .
10:  end for
11:  Set  $f(t) = \begin{cases} \ell & \text{with probability } 2^{-\ell} \quad 1 < \ell \leq \ell_{\max} \\ 1 & \text{otherwise} \end{cases}$ .
12:  Select  $\mathbf{a}_t \leftarrow \operatorname{argmax}_{\mathbf{a} \in \mathcal{D}_t} \max_{\boldsymbol{\mu} \in \mathcal{C}_{t,f(t)}} \langle \boldsymbol{\mu}, \mathbf{a} \rangle$  and observe  $r_t, \sigma_t$ .
13:  Set  $\bar{\sigma}_t = \max\{(R+1)/\sqrt{d}, \sigma_t\}$ .
14:   $\Sigma_{t+1} \leftarrow \Sigma_t + \mathbf{a}_t \mathbf{a}_t^\top / \bar{\sigma}_t^2$ ,  $\mathbf{c}_{t+1} \leftarrow \mathbf{c}_t + r_t \mathbf{a}_t / \bar{\sigma}_t^2$ ,  $\boldsymbol{\mu}_{t+1} \leftarrow \Sigma_{t+1}^{-1} \mathbf{c}_{t+1}$ .
15:  for  $\ell \neq f(t)$  do
16:     $\Sigma_{t+1,\ell} \leftarrow \Sigma_{t,\ell}$ ,  $\mathbf{c}_{t+1,\ell} \leftarrow \mathbf{c}_{t,\ell}$ ,  $\boldsymbol{\mu}_{t+1,\ell} \leftarrow \boldsymbol{\mu}_{t,\ell}$ .
17:  end for
18:   $\Sigma_{t+1,f(t)} \leftarrow \Sigma_{t,f(t)} + \mathbf{a}_t \mathbf{a}_t^\top / \bar{\sigma}_t^2$ ,  $\mathbf{c}_{t+1,f(t)} \leftarrow \mathbf{c}_{t,f(t)} + r_t \mathbf{a}_t / \bar{\sigma}_t^2$ .
19:   $\boldsymbol{\mu}_{t+1,f(t)} \leftarrow \Sigma_{t+1,f(t)}^{-1} \mathbf{c}_{t+1,f(t)}$ .
20: end for

```

174 By slightly modifying the self-normalized martingale concentration inequality proposed in [Abbasi-](#)
175 [Yadkori et al. \(2011\)](#), we can conclude that $\boldsymbol{\mu}^*$ belongs to the ellipsoid $\|\boldsymbol{\mu} - \boldsymbol{\mu}_t\|_{\Sigma_t^{-1}} \leq \beta_t$ with high
176 probability, where $\beta_t = \tilde{O}(R\sqrt{d} + C\sqrt{d})$. Such a confidence bound leads to a final regret which
177 has a polynomial dependence on C . However, such a simple approach have two limitations. First,
178 the agent does not know C apriori in our setting, thus it is impossible to set β_t to be dependent on
179 C . Second, vanilla ridge regression estimator does not consider different variances σ_t in each round,
180 thus it only gives a very conservative estimation.

181 **Action partition scheme.** To address the unknown C issue, besides the original estimator $\boldsymbol{\mu}_t$ which
182 uses all previous data, Algorithm 1 maintains several additional learners to learn $\boldsymbol{\mu}^*$ at different
183 accuracy level simultaneously, and it *randomly* selects one of the learners with different probabilities
184 at each round. Such a ‘‘parallel learning’’ idea is inspired by [Lykouris et al. \(2018\)](#). In detail,
185 we partition the observed data into ℓ_{\max} levels indexed by $[\ell_{\max}]$ and maintain ℓ_{\max} sub-sampled
186 estimators $\boldsymbol{\mu}_{t,1}, \dots, \boldsymbol{\mu}_{t,\ell_{\max}}$. According to line 11, the observed data in round t goes into level ℓ with
187 probability $2^{-\ell}$ if $1 < \ell \leq \ell_{\max}$ and it goes to level 1 with probability $1 - \sum_{\ell=2}^{\ell_{\max}} 2^{-\ell} = 1/2 + 2^{-\ell_{\max}}$.
188 The intuition is that if $2^\ell \geq C$, then the corruption level experienced by level ℓ

$$\text{Corruption}_{t,\ell} = \sum_{i=1}^t \frac{\mathbb{1}(f(i) = \ell)}{R+1} \cdot \sup_{\mathbf{a} \in \mathcal{D}_i} |r_i(\mathbf{a}) - r'_i(\mathbf{a})| \quad (4.2)$$

189 can be bounded by some quantity that is *independent of* C . That says, the individual learners whose
190 level is greater than $\log C$ can learn $\boldsymbol{\mu}^*$ successfully, even with the corruption. For the learners whose
191 level is less than $\log C$, we can also control the error by controlling the probability for the agent to
192 select them.

193 **Weighted regression estimator.** After introducing the partition scheme, we still need to deal
194 with the varying variance (heteroscedastic) case. Similar to ([Kirschner and Krause, 2018](#); [Zhou](#)
195 [et al., 2020](#)), we proposed the following *weighted ridge regression estimator*, which incorporates the
196 variance information of the rewards into estimation:

$$\boldsymbol{\mu}_t = \operatorname{argmin}_{\boldsymbol{\mu} \in \mathbb{R}^d} \lambda \|\boldsymbol{\mu}\|_2^2 + \sum_{i=1}^{t-1} [\langle \boldsymbol{\mu}, \mathbf{a}_i \rangle - r_i]^2 / \bar{\sigma}_i^2. \quad (4.3)$$

197 Here $\bar{\sigma}_t$ is defined as the upper bound of the true variance σ_t in line 13. The closed-form solution to
 198 (4.3) is calculated at each round in line 14. The use of $\bar{\sigma}_t$, as we will show later, makes our estimator
 199 more efficient in the heteroscedastic case. Meanwhile, we also apply our weighted regression
 200 estimator to each individual learner, and their estimator $\boldsymbol{\mu}_{t,\ell}$ can be written as follows:

$$\boldsymbol{\mu}_{t,\ell} = \operatorname{argmin}_{\boldsymbol{\mu} \in \mathbb{R}^d} \lambda \|\boldsymbol{\mu}\|_2^2 + \sum_{i=1}^{t-1} \mathbb{1}(f(i) = \ell) \cdot [\langle \boldsymbol{\mu}, \mathbf{a}_i \rangle - r_i]^2 / \bar{\sigma}_i^2. \quad (4.4)$$

201 The closed-form solution to (4.4) is calculated at each round in lines 15–20.

202 **Final Multi-Level confidence sets.** With the estimators $\boldsymbol{\mu}_t, \boldsymbol{\mu}_{t,1}, \dots, \boldsymbol{\mu}_{t,\ell_{\max}}$ at the beginning of
 203 round t , we define a cascade of candidate confidence sets as in lines 6–10, where

$$\beta_{t,\ell} = 8\sqrt{d \log \frac{(R+1)^2 \lambda + tA^2}{(R+1)^2 \lambda} \log(4t^2/\delta) + 4\sqrt{d} \log(4t^2/\delta) + 2^\ell \sqrt{d} + \sqrt{\lambda} B}, \quad (4.5)$$

$$\gamma_{t,\ell} = 8\sqrt{d \log \frac{(R+1)^2 \lambda + tA^2}{(R+1)^2 \lambda} \log(8t^2 T/\delta) + 4\sqrt{d} \log(8t^2 T/\delta) + \bar{C}_\ell \sqrt{d} + \sqrt{\lambda} B}, \quad (4.6)$$

204 with $\bar{C}_\ell = \log(2\ell^2/\delta) + 3$. For simplicity, we define

$$\ell^* = \max\{2, \lceil \log_2 C \rceil\} \quad (4.7)$$

205 as an important threshold in our later proof for regret bound analysis. Later we will prove that $\mathcal{C}_{t,\ell}$
 206 contains $\boldsymbol{\mu}^*$ for all $\ell \geq \ell^*, t \geq 1$ with high probability.

207 Note that each candidate confidence set can be written as the intersection of two ellipsoids. The
 208 intuition behind our construction of candidate confidence sets is that we hope that $\mathcal{C}_{t,\ell}$ is robust
 209 enough to handle the 2^ℓ -corrupted case, i.e., $\boldsymbol{\mu}^* \in \mathcal{C}_{t,\ell}$ with high probability. To achieve this, the first
 210 ellipsoid makes use of the global information and the “radius” $\beta_{t,\ell}$ need to contain a factor of 2^ℓ to
 211 tolerate a corruption level of 2^ℓ , and the second ellipsoid makes use of the observed data in level ℓ
 212 since this level only contain a few times of corruptions in 2^ℓ -corrupted case.

213 **Action selection.** With the candidate confidence sets, we use line 11 to randomly decide one
 214 confidence set and select an action based on the optimism-in-the-face-of-uncertainty (OFU) principle
 215 in line 12. Then we update the estimators for the next round $t + 1$.

216 **Remark 4.1.** Our algorithm shares a similar strategy for partitioning the observed data with the
 217 algorithm in Lykouris et al. (2018) but note that there is a major difference in that: Lykouris et al.
 218 (2018) regard the partition scheme as a “layer structure”, i.e., their algorithm further uses different
 219 estimators in layers of parallel learners and do action elimination layer by layer in each round. In
 220 contrast, the sub-sampled estimators in our algorithm are used independently, i.e., the selected action
 221 only relies on one of the partitions. As a result, Algorithm 1 does not need to do action elimination,
 222 thus is capable of handling the cases where the number of actions is huge or even infinite.

223 5 Main Results

224 In this section we present our main theorem, which establishes the regret bound for Multi-level
 225 OFUL.

Theorem 5.1. Set $\lambda = 1/B^2$. Suppose that $C = \Omega(1), R = \Omega(1)$, for all $t \geq 1$ and all $\mathbf{a} \in \mathcal{D}_t$,
 $\langle \mathbf{a}, \boldsymbol{\mu}^* \rangle \in [-1, 1]$. Then with probability at least $1 - 3\delta$, the regret of Algorithm 1 is bounded as
 follows:

$$\mathbf{Regret}(T) = \tilde{O} \left(C^2 d \sqrt{\sum_{t=1}^T \sigma_t^2} + C^2 \sqrt{dT} + CR\sqrt{dT} \right).$$

226 **Remark 5.2.** When $\sigma_t, R = \Omega(1)$, the regret bound in Theorem 5.1 matches the regret bound of
 227 OFUL proposed in Zhou et al. (2020) when the corruption level C is a constant.

228 **Remark 5.3.** Compared with the $\tilde{O}(d\sqrt{T} + C)$ result in Lee et al. (2021), our result has a multi-
 229 plicative quadratic dependence on C , which seems to be worse. However, we want to emphasize that
 230 we focus on a more challenging contextual bandits setting where the decision sets \mathcal{D}_t at each round
 231 are not identical, which is different from that in Lee et al. (2021). Therefore, our result and that in
 232 Lee et al. (2021) are not directly comparable.

233 **Remark 5.4.** Note that this instance-independent regret upper bound also holds in a stronger model
 234 than the one described in Section 3, where the adversary can even decide the decision set \mathcal{D}_t at each
 235 round t since our regret bound can hold without any assumption on the decision sets.

Corollary 5.5. Under the same conditions as in Theorem 5.1, if σ_t given by the environment are all R , the regret of Algorithm 1 is bounded by:

$$\mathbf{Regret}(T) = \tilde{O}\left(C^2 d R \sqrt{T}\right).$$

236 We also provide a gap-dependent regret bound.

Theorem 5.6. Suppose that $C = \Omega(1)$, $R = \Omega(1)$, for all $t \geq 1$ and all $\mathbf{a} \in \mathcal{D}_t$, $\langle \mathbf{a}, \boldsymbol{\mu}^* \rangle \in [-1, 1]$. Then with probability at least $1 - 3\delta$, the regret of Algorithm 1 is bounded as follows:

$$\mathbf{Regret}(T) = \tilde{O}\left(\frac{1}{\Delta} \cdot C^2 R^2 d + \frac{1}{\Delta} \cdot d^2 C^2 \max_{t \in [T]} \sigma_t^2\right).$$

237 **Remark 5.7.** Theorem 5.6 automatically suggests an $\tilde{O}(R^2 d^2 C^2 / \Delta)$ regret bound, by the fact
 238 $\sigma_t = O(R)$. Compared with previous result $\tilde{O}(d^{5/2} C / \Delta + d^6 / \Delta^2)$ (Lee et al., 2021), our result
 239 has a better dependence on the dimension d but a worse dependence on the corruption level C . As
 240 Remark 5.3 suggests, we focus on a more challenging contextual bandits setting, and the worse
 241 dependence on C might be due to this.

242 6 Proof Outline

243 First we have the following lemma which is a corruption-tolerant variant of Bernstein inequality for
 244 self-normalized vector-valued martingales introduced in Zhou et al. (2020).

Lemma 6.1 (Bernstein inequality for vector-valued martingales with corruptions). Let $\{\mathcal{G}_t\}_{t=1}^\infty$ be a filtration, $\{\mathbf{x}_t, \eta_t\}_{t \geq 1}$ a stochastic process so that $\mathbf{x}_t \in \mathbb{R}^d$ is \mathcal{G}_t -measurable and $\eta_t \in \mathbb{R}$ is \mathcal{G}_{t+1} -measurable. Fix $R, L, \sigma, \lambda > 0$, $\boldsymbol{\mu}^* \in \mathbb{R}^d$. For $t \geq 1$ let $y_t^{\text{stoch}} = \langle \boldsymbol{\mu}^*, \mathbf{x}_t \rangle + \eta_t$ and suppose that η_t, \mathbf{x}_t also satisfy

$$|\eta_t| \leq R, \mathbb{E}[\eta_t | \mathcal{G}_t] = 0, \mathbb{E}[\eta_t^2 | \mathcal{G}_t] \leq \sigma^2, \|\mathbf{x}_t\|_2 \leq L.$$

245 Suppose $\{y_t\}$ is a sequence such that $\sum_{i=1}^t |y_i - y_i^{\text{stoch}}| = C(t)$ for all $t \geq 1$. Then, for any
 246 $0 < \delta < 1$, with probability at least $1 - \delta$ we have $\forall t > 0$,

$$\|\boldsymbol{\mu}_t - \boldsymbol{\mu}^*\|_{\mathbf{Z}_t} \leq \beta_t + C(t) + \sqrt{\lambda} \|\boldsymbol{\mu}^*\|_2,$$

where for $t \geq 1$, $\boldsymbol{\mu}_t = \mathbf{Z}_t^{-1} \mathbf{b}_t$, $\mathbf{Z}_t = \lambda \mathbf{I} + \sum_{i=1}^t \mathbf{x}_i \mathbf{x}_i^\top$, $\mathbf{b}_t = \sum_{i=1}^t y_i \mathbf{x}_i$, and

$$\beta_t = 8\sigma \sqrt{d \log \frac{d\lambda + tL^2}{d\lambda} \log(4t^2/\delta) + 4R \log(4t^2/\delta)}.$$

247 Next, we have that with high probability, all the level $\ell \geq \ell^*$ only influenced by limited amount of
 248 corruptions as mentioned in Section 4.

Lemma 6.2. Let $\text{Corruption}_{t,\ell}$ be defined in (4.2). Then we have with probability at least $1 - \delta$, for all $\ell \geq \ell^*$, $t \geq 1$:

$$\text{Corruption}_{t,\ell} \leq \bar{C}_\ell = \log(2\ell^2/\delta) + 3.$$

249 We denote by \mathcal{E}_{sub} the event that the above inequality holds.

250 We define the following event to further show that our candidate confidence sets with $\ell \geq \ell^*$ are
 251 “robust” enough, i.e. $\mathcal{C}_{t,\ell}$ contains $\boldsymbol{\mu}^*$ with high probability.

252 **Definition 6.3.** Let ℓ^* be defined in (4.7). We introduce the event \mathcal{E}_1 as follows.

$$\mathcal{E}_1 := \left\{ \forall \ell \geq \ell^* \text{ and } t \geq 1, \|\boldsymbol{\mu}^* - \boldsymbol{\mu}_t\|_{\boldsymbol{\Sigma}_t} \leq \beta_{t,\ell} \text{ and } \|\boldsymbol{\mu}^* - \boldsymbol{\mu}_{t,\ell}\|_{\boldsymbol{\Sigma}_{t,\ell}} \leq \gamma_{t,\ell} \right\}. \quad (6.1)$$

253 where $\beta_{t,\ell}, \gamma_{t,\ell}$ are defined in (4.5) and (4.6).

254 Next lemma suggests that the event \mathcal{E}_1 happens with high probability.

255 **Lemma 6.4.** Let \mathcal{E}_1 be defined in (6.1). For any $0 < \delta < 1/3$, we have $\mathbb{P}(\mathcal{E}_1) \geq 1 - 3\delta$.

256 For simplicity, we define $\mathbf{a}_{t,\ell} = \operatorname{argmax}_{\mathbf{a} \in \mathcal{D}_t} \max_{\boldsymbol{\mu} \in \mathcal{C}_{t,\ell}} \langle \boldsymbol{\mu}, \mathbf{a} \rangle$ for each level ℓ . \mathbf{a}_t can be seen as an
 257 action vector randomly chosen from $\mathbf{a}_{t,\ell}$, $\ell \in [\ell_{\max}]$. Next two lemmas suggest that under event \mathcal{E}_1 ,
 258 at each round, the gap between the optimal reward and the selected reward can be upper bounded by
 259 some bonus terms related to $\mathbf{a}_{t,\ell}$.

260 **Lemma 6.5.** Suppose \mathcal{E}_1 occurs. If $f(t) \leq \ell^*$, we have $\langle \mathbf{a}_t^* - \mathbf{a}_t, \boldsymbol{\mu}^* \rangle \leq 2\beta_{t,\ell^*} \|\mathbf{a}_t\|_{\Sigma_t^{-1}} +$
 261 $2\beta_{t,\ell^*} \|\mathbf{a}_{t,\ell^*}\|_{\Sigma_t^{-1}}$.

262 **Lemma 6.6.** On event \mathcal{E}_1 , if $f(t) = \ell > \ell^*$, we have $\langle \mathbf{a}_t^* - \mathbf{a}_t, \boldsymbol{\mu}^* \rangle \leq 2\gamma_{t,\ell} \|\mathbf{a}_t\|_{\Sigma_t^{-1}}$.

263 Now we provide the proof sketch of Theorem 5.1.

264 *Proof sketch of Theorem 5.1.* Suppose \mathcal{E}_1 occurs. The main idea to bound the regret is to decompose
 265 the total rounds $[T]$ into two non-overlapping parts, based on which individual learner is selected at
 266 that round. In detail, we have

$$\begin{aligned} \text{Regret}(T) &= \mathbb{E} \left[\sum_{t=1}^T (\langle \mathbf{a}_t^*, \boldsymbol{\mu}^* \rangle - \langle \mathbf{a}_t, \boldsymbol{\mu}^* \rangle) \right] \\ &= \mathbb{E} \left[\underbrace{\sum_{t=1}^T \mathbb{1}(f(t) \leq \ell^*) (\langle \mathbf{a}_t^*, \boldsymbol{\mu}^* \rangle - \langle \mathbf{a}_t, \boldsymbol{\mu}^* \rangle)}_{I_1} \right] \\ &\quad + \sum_{\ell=\ell^*+1}^{\ell_{\max}} \mathbb{E} \left[\underbrace{\sum_{t=1}^T \mathbb{1}(f(t) = \ell) (\langle \mathbf{a}_t^*, \boldsymbol{\mu}^* \rangle - \langle \mathbf{a}_t, \boldsymbol{\mu}^* \rangle)}_{I_2(\ell)} \right]. \end{aligned} \quad (6.2)$$

267 Here I_1 represents the regret where the "low-level" learner is selected, where the corruption level
 268 is beyond the learner level. For this case, by Lemma 6.5, we can directly show that

$$I_1 \leq \mathbb{E} \left[\sum_{t=1}^T \mathbb{1}(f(t) \leq \ell^*) \min \left\{ 2, 2\beta_{t,\ell^*} \|\mathbf{a}_{t,\ell^*}\|_{\Sigma_t^{-1}} + 2\beta_{t,\ell^*} \|\mathbf{a}_t\|_{\Sigma_t^{-1}} \right\} \right]. \quad (6.3)$$

269 We further bound (6.3). Let \mathcal{F}_t be the σ -algebra generated by $\mathbf{a}_s, r_s, \sigma_s, f(s)$ for $s \leq t-1$.
 270 Then by the property of our partition scheme (note that $\mathbb{P}(f(t) = \ell^*) = 2^{-\ell^*}$), we can show that
 271 $\mathbb{E} \left[\mathbb{1}(f(t) \leq \ell^*) \|\mathbf{a}_{t,\ell^*}\|_{\Sigma_t^{-1}} | \mathcal{F}_t \right] \leq 2^{\ell^*} \mathbb{E} \left[\|\mathbf{a}_t\|_{\Sigma_t^{-1}} | \mathcal{F}_t \right]$. Therefore, we can further bound I_1 by

$$I_1 \leq 4 \cdot 2^{\ell^*} \mathbb{E} \left[\underbrace{\sum_{t=1}^T \min \left\{ 2, \beta_{T,\ell^*} \|\mathbf{a}_t\|_{\Sigma_t^{-1}} \right\}}_{I_3} \right]. \quad (6.4)$$

272 To further bound I_3 , we split $[T]$ into 2 parts, $\mathcal{I}_1 = \{t \in [T] \mid \|\mathbf{a}_t/\bar{\sigma}_t\|_{\Sigma_t^{-1}} > 1\}$, $\mathcal{I}_2 = \{t \in$
 273 $[T] \mid \|\mathbf{a}_t/\bar{\sigma}_t\|_{\Sigma_t^{-1}} \leq 1\}$ to bound I_3 . The intuition here is that the cardinality of \mathcal{I}_1 is bounded, and
 274 the sum of terms with $t \in \mathcal{I}_2$ can be bounded using Cauchy-Schwarz inequality.

$$\sum_{t \in \mathcal{I}_1} \min \left\{ 2, \beta_{T,\ell^*} \|\mathbf{a}_t\|_{\Sigma_t^{-1}} \right\} \leq 2|\mathcal{I}_1| \leq 2 \sum_{t=1}^T \min \left\{ 1, \|\mathbf{a}_t/\bar{\sigma}_t\|_{\Sigma_t^{-1}}^2 \right\} \leq 4d \log \frac{(R+1)^2 \lambda + TA^2}{(R+1)^2 \lambda}, \quad (6.5)$$

275 where the first inequality holds since $\min \left\{ 2, \beta_{T,\ell^*} \|\mathbf{a}_t\|_{\Sigma_t^{-1}} \right\} \leq 2$, the second inequality follows
 276 from the definition of \mathcal{I}_1 , the third inequality holds by Lemma C.2.

$$\begin{aligned}
\sum_{t \in \mathcal{I}_2} \min \left\{ 2, \beta_{T, \ell^*} \|\mathbf{a}_t\|_{\Sigma_t^{-1}} \right\} &\leq \beta_{T, \ell^*} \sqrt{\sum_{t \in \mathcal{I}_2} \bar{\sigma}_t^2} \cdot \sqrt{\sum_{t \in \mathcal{I}_2} \min \left\{ 1, \|\mathbf{a}_t / \bar{\sigma}_t\|_{\Sigma_t^{-1}}^2 \right\}} \\
&\leq \beta_{T, \ell^*} \sqrt{(R+1)^2 T / d + \sum_{t=1}^T \sigma_t^2} \cdot \sqrt{2d \log \frac{(R+1)^2 \lambda + T A^2}{(R+1)^2 \lambda}},
\end{aligned} \tag{6.6}$$

277 where the first inequality follows from Cauchy-Schwarz inequality, the second inequality follows
278 from the definition of $\bar{\sigma}_t$ and Lemma C.2.

279 Substituting (6.5) and (6.6) into (6.3), we have

$$I_1 = \tilde{O} \left(C^2 d \sqrt{\sum_{t=1}^T \sigma_t^2} + C^2 \sqrt{dT} + CR \sqrt{dT} \right). \tag{6.7}$$

280 Now it remains to bound $I_2(\ell)$. By Lemma 6.6, we have

$$I_2(\ell) \leq 2\mathbb{E} \underbrace{\left[\sum_{t=1}^T \mathbb{1}(f(t) = \ell) \min \left\{ 1, \gamma_{t, \ell} \|\mathbf{a}_{t, \ell}\|_{\Sigma_{t, \ell}^{-1}} \right\} \right]}_{I_4} = \tilde{O} \left(R\sqrt{Td} + d \sqrt{\sum_{t=1}^T \sigma_t^2} \right), \tag{6.8}$$

281 where the second equality can be proved by analysis similar to that of (6.5) and (6.6). Finally,
282 substituting (6.7) and (6.8) into (6.2) ends our proof.

283 □

284 7 Conclusion and Future Work

285 In this paper, we have considered the linear contextual bandits problem in the presence of adversarial
286 corruptions. We propose a Multi-level OFUL algorithm, which is provably robust to the adversarial
287 attacks. We prove a gap-independent regret bound of $\tilde{O} \left(C^2 d \sqrt{\sum_{t=1}^T \sigma_t^2} + C^2 \sqrt{dT} + CR \sqrt{dT} \right)$
288 together with a gap-dependent bound of $\tilde{O} \left(\frac{1}{\Delta} \cdot C^2 R^2 d + \frac{1}{\Delta} \cdot d^2 C^2 \max_{t \in [T]} \sigma_t^2 \right)$.

289 We leave it as an open question that whether the multiplicative dependence on C^2 in the regret upper
290 bounds can be removed without making additional assumptions in our setting.

291 References

- 292 ABBASI-YADKORI, Y., PÁL, D. and SZEPESVÁRI, C. (2011). Improved algorithms for linear
293 stochastic bandits. In *NIPS*, vol. 11.
- 294 ABE, N., BIERMANN, A. W. and LONG, P. M. (2003). Reinforcement learning with immediate
295 rewards and linear hypotheses. *Algorithmica* **37** 263–293.
- 296 AUER, P. (2002). Using confidence bounds for exploitation-exploration trade-offs. *Journal of*
297 *Machine Learning Research* **3** 397–422.
- 298 AUER, P. and CHIANG, C.-K. (2016). An algorithm with nearly optimal pseudo-regret for both
299 stochastic and adversarial bandits. In *Conference on Learning Theory*. PMLR.
- 300 BOGUNOVIC, I., LOSALKA, A., KRAUSE, A. and SCARLETT, J. (2021). Stochastic linear bandits
301 robust to adversarial attacks. In *International Conference on Artificial Intelligence and Statistics*.
302 PMLR.
- 303 BUBECK, S. and SLIVKINS, A. (2012). The best of both worlds: Stochastic and adversarial bandits.
304 In *Conference on Learning Theory*. JMLR Workshop and Conference Proceedings.

- 305 CHU, W., LI, L., REYZIN, L. and SCHAPIRE, R. (2011). Contextual bandits with linear payoff
306 functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and*
307 *Statistics*. JMLR Workshop and Conference Proceedings.
- 308 DANI, V., HAYES, T. P. and KAKADE, S. (2008). Stochastic linear optimization under bandit
309 feedback. In *COLT*.
- 310 DESHPANDE, Y. and MONTANARI, A. (2012). Linear bandits in high dimension and recommendation
311 systems. In *2012 50th Annual Allerton Conference on Communication, Control, and Computing*
312 *(Allerton)*. IEEE.
- 313 GARCELON, E., ROZIERE, B., MEUNIER, L., TARBOURIECH, J., TEYTAUD, O., LAZARIC, A.
314 and PIROTTA, M. (2020). Adversarial attacks on linear contextual bandits. *arXiv preprint*
315 *arXiv:2002.03839*.
- 316 GUPTA, A., KOREN, T. and TALWAR, K. (2019). Better algorithms for stochastic bandits with
317 adversarial corruptions. In *Conference on Learning Theory*. PMLR.
- 318 JHALANI, T., KANT, V. and DWIVEDI, P. (2016). A linear regression approach to multi-criteria
319 recommender system. In *International Conference on Data Mining and Big Data*. Springer.
- 320 JUN, K.-S., LI, L., MA, Y. and ZHU, X. J. (2018). Adversarial attacks on stochastic bandits. In
321 *NeurIPS*.
- 322 KANNAN, S., MORGENSTERN, J. H., ROTH, A., WAGGONER, B. and WU, Z. (2018). A smoothed
323 analysis of the greedy algorithm for the linear contextual bandit problem. In *NeurIPS*.
- 324 KAPOOR, S., PATEL, K. K. and KAR, P. (2019). Corruption-tolerant bandit learning. *Machine*
325 *Learning* **108** 687–715.
- 326 KIRSCHNER, J. and KRAUSE, A. (2018). Information directed sampling and bandits with het-
327 eroscedastic noise. In *Conference On Learning Theory*. PMLR.
- 328 LEE, C.-W., LUO, H., WEI, C.-Y., ZHANG, M. and ZHANG, X. (2021). Achieving near instance-
329 optimality and minimax-optimality in stochastic and adversarial linear bandits simultaneously.
330 *arXiv preprint arXiv:2102.05858*.
- 331 LI, L., CHU, W., LANGFORD, J. and SCHAPIRE, R. E. (2010). A contextual-bandit approach to
332 personalized news article recommendation. In *Proceedings of the 19th international conference on*
333 *World wide web*.
- 334 LI, Y., LOU, E. Y. and SHAN, L. (2019a). Stochastic linear optimization with adversarial corruption.
335 *arXiv preprint arXiv:1909.02109*.
- 336 LI, Y., WANG, Y. and ZHOU, Y. (2019b). Nearly minimax-optimal regret for linearly parameterized
337 bandits. In *Conference on Learning Theory*. PMLR.
- 338 LIU, F. and SHROFF, N. (2019). Data poisoning attacks on stochastic bandits. In *International*
339 *Conference on Machine Learning*. PMLR.
- 340 LYKOURIS, T., MIRROKNI, V. and PAES LEME, R. (2018). Stochastic bandits robust to adversarial
341 corruptions. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*.
- 342 NEU, G. and OLKHOVSKAYA, J. (2020). Efficient and robust algorithms for adversarial linear
343 contextual bandits. In *Conference on Learning Theory*. PMLR.
- 344 RUSMEVICHIENTONG, P. and TSITSIKLIS, J. N. (2010). Linearly parameterized bandits. *Mathemat-*
345 *ics of Operations Research* **35** 395–411.
- 346 SELDIN, Y. and LUGOSI, G. (2017). An improved parametrization and analysis of the exp3++
347 algorithm for stochastic and adversarial bandits. In *Conference on Learning Theory*. PMLR.
- 348 SELDIN, Y. and SLIVKINS, A. (2014). One practical algorithm for both stochastic and adversarial
349 bandits. In *International Conference on Machine Learning*. PMLR.

- 350 VILLAR, S. S., BOWDEN, J. and WASON, J. (2015). Multi-armed bandit models for the optimal
 351 design of clinical trials: benefits and challenges. *Statistical science: a review journal of the*
 352 *Institute of Mathematical Statistics* **30** 199.
- 353 ZHOU, D., GU, Q. and SZEPESVARI, C. (2020). Nearly minimax optimal reinforcement learning for
 354 linear mixture markov decision processes. *arXiv preprint arXiv:2012.08507* .
- 355 ZIMMERT, J. and SELDIN, Y. (2019). An optimal algorithm for stochastic and adversarial bandits.
 356 In *The 22nd International Conference on Artificial Intelligence and Statistics*. PMLR.

357 Checklist

- 358 1. For all authors...
- 359 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
 360 contributions and scope? [Yes]
- 361 (b) Did you describe the limitations of your work? [Yes]
- 362 (c) Did you discuss any potential negative societal impacts of your work? [N/A] Our work
 363 studies the regret bounds for contextual linear bandits with corruption. That is a pure
 364 theoretical problem, thus it does not have any negative social impact.
- 365 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 366 them? [Yes]
- 367 2. If you are including theoretical results...
- 368 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 369 (b) Did you include complete proofs of all theoretical results? [Yes]
- 370 3. If you ran experiments...
- 371 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 372 mental results (either in the supplemental material or as a URL)? [N/A]
- 373 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 374 were chosen)? [N/A]
- 375 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 376 ments multiple times)? [N/A]
- 377 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 378 of GPUs, internal cluster, or cloud provider)? [N/A]
- 379 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 380 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 381 (b) Did you mention the license of the assets? [N/A]
- 382 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 383
- 384 (d) Did you discuss whether and how consent was obtained from people whose data you're
 385 using/curating? [N/A]
- 386 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 387 information or offensive content? [N/A]
- 388 5. If you used crowdsourcing or conducted research with human subjects...
- 389 (a) Did you include the full text of instructions given to participants and screenshots, if
 390 applicable? [N/A]
- 391 (b) Did you describe any potential participant risks, with links to Institutional Review
 392 Board (IRB) approvals, if applicable? [N/A]
- 393 (c) Did you include the estimated hourly wage paid to participants and the total amount
 394 spent on participant compensation? [N/A]