
Nearly Optimal Algorithms for Contextual Dueling Bandits from Adversarial Feedback

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Learning from human feedback plays an important role in aligning generative
2 models, such as large language models (LLM). However, the effectiveness of
3 this approach can be influenced by adversaries, who may intentionally provide
4 misleading preferences to manipulate the output in an undesirable or harmful
5 direction. To tackle this challenge, we study a specific model within this problem
6 domain—contextual dueling bandits with adversarial feedback, where the true
7 preference label can be flipped by an adversary. We propose an algorithm namely
8 robust contextual dueling bandits (RCDB), which is based on uncertainty-weighted
9 maximum likelihood estimation. Our algorithm achieves an $\tilde{O}(d\sqrt{T} + dC)$ regret
10 bound, where T is the number of rounds, d is the dimension of the context, and
11 $0 \leq C \leq T$ is the total number of adversarial feedback. We also prove a lower
12 bound to show that our regret bound is nearly optimal, both in scenarios with and
13 without ($C = 0$) adversarial feedback. Additionally, we conduct experiments to
14 evaluate our proposed algorithm against various types of adversarial feedback.
15 Experimental results demonstrate its superiority over the state-of-the-art dueling
16 bandit algorithms in the presence of adversarial feedback.

17 1 Introduction

18 Acquiring an appropriate reward proves challenging in numerous real-world applications, often
19 necessitating intricate instrumentation (Zhu et al., 2020) and time-consuming calibration (Yu et al.,
20 2020) to achieve satisfactory levels of sample efficiency. For instance, in training large language
21 models (LLM) using reinforcement learning from human feedback (RLHF), the diverse values and
22 perspectives of humans can lead to uncalibrated and noisy rewards (Ouyang et al., 2022). In contrast,
23 preference-based data, which involves comparing or ranking various actions, is a more straightforward
24 method for capturing human judgments and decisions. In this context, the dueling bandit model
25 (Yue et al., 2012) provides a problem framework that focuses on optimal decision-making through
26 pairwise comparisons, rather than relying on the absolute reward for each action.
27 However, human feedback may not always be reliable. In real-world applications, human feedback
28 is particularly vulnerable to manipulation through preference label flip. Adversarial feedback can
29 significantly increase the risk of misleading a large language model (LLM) into erroneously prioritiz-
30 ing harmful content, under the false belief that it reflects human preference. Despite the significant
31 influence of adversarial feedback, there is limited existing research on the impact of adversarial
32 feedback specifically within the context of dueling bandits. A notable exception is Agarwal et al.
33 (2021), which studies dueling bandits when an adversary can flip some of the preference labels
34 received by the learner. They proposed an algorithm that is agnostic to the amount of adversarial
35 feedback introduced by the adversary. However, their setting has the following two limitations.
36 First, their study was confined to a finite-armed setting, which renders their results less applicable
37 to modern applications such as RLHF. Second, their adversarial feedback is defined on the whole
38 comparison matrix. In each round, the adversary observes the outcomes of all pairwise comparisons
39 and then decides to corrupt some of the pairs before the agent selects the actions. This assumption

40 does not align well with the real-world scenario, where the adversary often flips the preference label
 41 based on the information of the selected actions.

42 In this paper, to address the above challenge, we aim to develop contextual dueling bandit algorithms
 43 that are robust to adversarial feedback. This enables us to effectively tackle problems involving
 44 a large number of actions while also taking advantage of contextual information. We specifically
 45 consider a scenario where the adversary knows the selected action pair and the true preference of
 46 their comparison. In this setting, the adversary’s only decision is whether to flip the preference label
 47 or not. We highlight our contributions as follows:

- 48 • We propose a new algorithm called robust contextual dueling bandits (RCDB), which integrates
 49 uncertainty-dependent weights into the Maximum Likelihood Estimator (MLE). Intuitively, our
 50 choice of weight is designed to induce a higher degree of skepticism about potentially “untrust-
 51 worthy” feedback. The agent is encouraged to focus more on feedback that is more likely to be
 52 genuine, effectively diminishing the impact of any adversarial feedback.
- 53 • We analyze the regret of our algorithm under at most C number of adversarial feedback. Our result
 54 consists of two terms: a C -independent term $\tilde{O}(d\sqrt{T})$, which matches the lower bound established
 55 in Bengs et al. (2022) for uncorrupted linear contextual dueling bandits, and a C -dependent term
 56 $\tilde{O}(dC)$. Furthermore, we establish a lower bound for dueling bandits with adversarial feedback,
 57 demonstrating the optimality of our adversarial term. Consequently, our algorithm for dueling
 58 bandits attains the optimal regret in both scenarios, with and without adversarial feedback.
- 59 • We conduct extensive experiments to validate the effectiveness of our algorithm RCDB. To compre-
 60 hensively assess RCDB’s robustness against adversarial feedback, we evaluate its performance under
 61 various types of adversarial feedback and compare the results with state-of-the-art dueling bandit
 62 algorithms. Experimental results demonstrate the superiority of our algorithm in the presence of
 63 adversarial feedback, which corroborate our theoretical analysis.

Table 1: Comparison of algorithms for robust bandits and dueling bandits.

Model	Algorithm	Setting	Regret
Bandits	Multi-layer Active Arm Elimination Race (Lykouris et al., 2018)	K -armed Bandits	$\tilde{O}(K^{1.5}C\sqrt{T})$
	BARBAR (Gupta et al., 2019)	K -armed Bandits	$\tilde{O}(\sqrt{KT} + KC)$
	SBE (Li et al., 2019)	Linear Bandits	$\tilde{O}(d^2C/\Delta + d^5/\Delta^2)$
	Robust Phased Elimination (Bogunovic et al., 2021)	Linear Bandits	$\tilde{O}(\sqrt{dT} + d^{1.5}C + C^2)$
	Robust weighted OFUL (Zhao et al., 2021)	Linear Contextual Bandits	$\tilde{O}(dC\sqrt{T})$
	CW-OFUL (He et al., 2022)	Linear Contextual Bandits	$\tilde{O}(d\sqrt{T} + dC)$
Dueling Bandits	WIWR (Agarwal et al., 2021)	K -armed Dueling Bandits	$\tilde{O}(K^2C/\Delta_{\min} + \sum_{i \neq i^*} K^2/\Delta_i^2)$
	Versatile-DB (Saha and Gaillard, 2022)	K -armed Dueling Bandits	$\tilde{O}(C + \sum_{i \neq i^*} 1/\Delta_i + \sqrt{K})$
	RCDB (Our work)	Contextual Dueling Bandits	$\tilde{O}(d\sqrt{T} + dC)$

64 **Notation.** In this paper, we use plain letters such as x to denote scalars, lowercase bold letters such
 65 as \mathbf{x} to denote vectors and uppercase bold letters such as \mathbf{X} to denote matrices. For a vector \mathbf{x} , $\|\mathbf{x}\|_2$
 66 denotes its ℓ_2 -norm. The weighted ℓ_2 -norm associated with a positive-definite matrix \mathbf{A} is defined
 67 as $\|\mathbf{x}\|_{\mathbf{A}} = \sqrt{\mathbf{x}^\top \mathbf{A} \mathbf{x}}$. For two symmetric matrices \mathbf{A} and \mathbf{B} , we use $\mathbf{A} \succeq \mathbf{B}$ to denote $\mathbf{A} - \mathbf{B}$ is
 68 positive semidefinite. We use $\mathbb{1}$ to denote the indicator function and $\mathbf{0}$ to denote the zero vector. For
 69 two actions a, b , we use $a \succ b$ to denote a is more preferable to b . For a positive integer N , we use
 70 $[N]$ to denote $\{1, 2, \dots, N\}$. We use standard asymptotic notations including $O(\cdot)$, $\Omega(\cdot)$, $\Theta(\cdot)$, and
 71 $\tilde{O}(\cdot)$, $\tilde{\Omega}(\cdot)$, $\tilde{\Theta}(\cdot)$ will hide logarithmic factors.

72 2 Related Work

73 **Bandits with Adversarial Reward.** The multi-armed bandit problem, involving an agent making
 74 sequential decisions among multiple arms, has been studied with both stochastic rewards (Lai
 75 et al., 1985; Lai, 1987; Auer, 2002; Auer et al., 2002a; Kalyanakrishnan et al., 2012; Lattimore and
 76 Szepesvári, 2020; Agrawal and Goyal, 2012), and adversarial rewards (Auer et al., 2002b; Bubeck
 77 et al., 2012). Moreover, a line of works focuses on designing algorithms that can achieve near-optimal
 78 regret bounds for both stochastic bandits and adversarial bandits simultaneously (Bubeck and Slivkins,
 79 2012; Seldin and Slivkins, 2014; Auer and Chiang, 2016; Seldin and Lugosi, 2017; Zimmert and
 80 Seldin, 2019; Lee et al., 2021), which is known as “the best of both worlds” guarantee. Distinct from

81 fully stochastic and fully adversarial models, Lykouris et al. (2018) studied a setting, where only a
82 portion of the rewards is subject to corruption. They proposed an algorithm with a regret dependent
83 on the corruption level C , defined as the cumulative sum of the corruption magnitudes in each round.
84 Their result is C times worse than the regret without corruption. Gupta et al. (2019) improved the
85 result by providing a regret guarantee comprising two terms, a corruption-independent term that
86 matches the regret lower bound without corruption, and a corruption-dependent term that is linear in
87 C . In addition, Gupta et al. (2019) proved a lower bound demonstrating the optimality of the linear
88 dependency on C .

89 **Contextual Bandits with Corruption.** Li et al. (2019) studied stochastic linear bandits with
90 corruption and presented an instance-dependent regret bound linearly dependent on the corruption
91 level C . Bogunovic et al. (2021) studied the same problem and proposed an algorithm with near-
92 optimal regret in the non-corrupted case. Lee et al. (2021) studied this problem in a different setting,
93 where the adversarial corruptions are generated through the inner product of a corrupted vector
94 and the context vector. For linear contextual bandits, Bogunovic et al. (2021) proved that under an
95 additional context diversity assumption, the regret of a simple greedy algorithm is nearly optimal
96 with an additive corruption term. Zhao et al. (2021) and Ding et al. (2022) extended the OFUL
97 algorithm (Abbasi-Yadkori et al., 2011) and proved a regret with a corruption term polynomially
98 dependent on the total number of rounds T . He et al. (2022) proposed an algorithm for known
99 corruption level C to remove the polynomial dependency on T in the corruption term, which only
100 has a linear dependency on C . They also proved a lower bound showing the optimality of linear
101 dependency on C for linear contextual bandits with a known corruption level. Additionally, He et al.
102 (2022) extended the proposed algorithm to an unknown corruption level and provided a near-optimal
103 performance guarantee that matches the lower bound. For more extensions, Kuroki et al. (2023)
104 studied best-of-both-worlds algorithms for linear contextual bandits. Ye et al. (2023) proposed a
105 corruption robust algorithm for nonlinear contextual bandits.

106 **Dueling Bandits and Logistic Bandits.** The dueling bandit model was first proposed in Yue
107 et al. (2012). Compared with bandits, the agent will select two arms and receive the preference
108 feedback between the two arms from the environment. For general preference, there may not exist
109 the “best” arm that always wins in the pairwise comparison. Therefore, various alternative winners
110 are considered, including Condorcet winner (Zoghi et al., 2014; Komiyama et al., 2015), Copeland
111 winner (Zoghi et al., 2015; Wu and Liu, 2016; Komiyama et al., 2016), Borda winner (Jamieson et al.,
112 2015; Falahatgar et al., 2017; Heckel et al., 2018; Saha et al., 2021; Wu et al., 2023) and von Neumann
113 winner (Ramamohan et al., 2016; Dudík et al., 2015; Balsubramani et al., 2016), along with their
114 corresponding performance metrics. To handle potentially large action space or context information,
115 Saha (2021) studied a structured contextual dueling bandit setting. In this setting, each arm possesses
116 an unknown intrinsic reward. The comparison is determined based on a logistic function of the relative
117 rewards. In a similar setting, Bengs et al. (2022) studied contextual linear stochastic transitivity
118 model with contextualized utilities. Di et al. (2023) proposed a layered algorithm with variance
119 aware regret bound. Another line of works does not make the reward assumption. Instead, they
120 assume the preference feedback can be represented by a function class. Saha and Krishnamurthy
121 (2022) designed an algorithm that achieves the optimal regret for K -armed contextual dueling bandit
122 problem. Sekhari et al. (2023) studied contextual dueling bandits in a more general setting and
123 proposed an algorithm that provides guarantees for both regret and the number of queries. Another
124 related area of research is the logistic bandits, where the agent selects one arm in each round and
125 receives a Bernoulli reward. Faury et al. (2020) studied the dependency with respect to the degree
126 of non-linearity of the logistic function κ . They proposed an algorithm with no dependency in κ .
127 Abeille et al. (2021) further improved the dependency on κ and proved a problem dependent lower
128 bound. Faury et al. (2022) proposed a computationally efficient algorithm with regret performance
129 still matching the lower-bound proved in Abeille et al. (2021).

130 **Dueling Bandits with Adversarial Feedback.** A line of work has focused on dueling bandits with
131 adversarial feedback or corruption. Gajane et al. (2015) studied a fully adversarial utility-based
132 version of dueling bandits, which was proposed in Ailon et al. (2014). Saha et al. (2021) considered
133 the Borda regret for adversarial dueling bandits without the assumption of utility. In a setting
134 parallel to that in Lykouris et al. (2018); Gupta et al. (2019), Agarwal et al. (2021) studied K -armed
135 dueling bandits in a scenario where an adversary has the capability to corrupt part of the feedback
136 received by the learner. They designed an algorithm whose regret comprises two terms: one that
137 is optimal in uncorrupted scenarios, and another that is linearly dependent on the total times of
138 adversarial feedback C . Later on, Saha and Gaillard (2022) achieved “best-of-both world” result for
139 noncontextual dueling bandits and improved the adversarial term of Agarwal et al. (2021) in the same

140 setting. For contextual dueling bandits, Wu et al. (2023) proposed an EXP3-type algorithm for the
 141 adversarial linear setting using Borda regret. For a comparison of the most related works for robust
 142 bandits and dueling bandits, please refer to Table 1. In this paper, we study the influence of adversarial
 143 feedback within contextual dueling bandits, particularly in a setting where only a minority of the
 144 feedback is adversarial. Compared to previous studies, most studies have focused on the multi-armed
 145 dueling bandit framework without integrating context information. The notable exception is Wu et al.
 146 (2023); however, this study does not provide guarantees regarding the dependency on the number of
 147 adversarial feedback instances.

148 3 Preliminaries

149 In this work, we study linear contextual dueling bandits with adversarial feedback. In each round
 150 $t \in [T]$, the agent observes the context information x_t from a context set \mathcal{X} and the corresponding
 151 action set \mathcal{A} . Utilizing this context information, the agent selects two actions, a_t and b_t . Subsequently,
 152 the environment will generate a binary feedback (i.e., preference label) $l_t = \mathbb{1}(a_t \succ b_t) \in \{0, 1\}$
 153 indicating the preferable action. We assume the existence of a reward function $r^*(x, a)$ dependent on
 154 the context information x and action a , and a monotonically increasing link function σ satisfying
 155 $\sigma(x) + \sigma(-x) = 1$. The preference probability will be determined by the link function and the
 156 difference between the rewards of the selected arms, i.e.,

$$\mathbb{P}(a \succ b|x) = \sigma(r^*(x, a) - r^*(x, b)). \quad (3.1)$$

157 We assume that the reward function is linear with respect to some known feature map $\phi(x, a)$. To be
 158 more specific, we make the following assumption:

159 **Assumption 3.1.** Let $\phi : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}^d$ be a known feature map, with $\|\phi(x, a)\|_2 \leq 1$ for any
 160 $(x, a) \in \mathcal{X} \times \mathcal{A}$. We define the reward function r_θ parameterized by $\theta \in \mathbb{R}^d$, with $r_\theta(x, a) =$
 161 $\langle \theta, \phi(x, a) \rangle$. Moreover, there exists θ^* satisfying $r_{\theta^*} = r^*$, with $\|\theta^*\|_2 \leq B$.

162 Similar assumptions have been made in the literature of dueling bandits (Saha, 2021; Bengs et al.,
 163 2022; Xiong et al., 2023). We also make an assumption on the derivative of the link function, which
 164 is common in the study of generalized linear models for bandits (Filippi et al., 2010).

165 **Assumption 3.2.** The link function σ is differentiable. Furthermore, its first-order derivative satisfies:

$$\dot{\sigma}(\cdot) \geq \kappa$$

166 for some constant $\kappa > 0$.

167 In our setting, however, the agent does not directly observe the true binary feedback. Instead, an
 168 adversary will see both the choice of the agent and the true feedback. Based on the information, the
 169 adversary can decide whether to corrupt the binary feedback or not.¹ We represent the adversary’s
 170 decision in round t by an adversarial indicator c_t , which takes values from the set $\{0, 1\}$. If the
 171 adversary chooses not to corrupt the result, we have $c_t = 0$. Otherwise, we have $c_t = 1$, which means
 172 adversarial feedback in this round. As a result, the agent will observe a flipped preference label, i.e.,
 173 the observation $o_t = 1 - l_t$. We define C as the total level of adversarial feedback, i.e.,

$$\sum_{t=1}^T c_t \leq C.$$

174 **Remark 3.3.** Adversarial corruption has been firstly studied in bandits (Lykouris et al., 2018), where
 175 in each round t , the agent selects an action a_t and the environment generates a numerical reward
 176 $r_t(a_t)$. The adversary observes the reward and returns a corrupted reward \bar{r}_t . The corruption level
 177 C is defined by $\sum_{t=1}^T |r_t(a_t) - \bar{r}_t| \leq C$. Compared with the continuous perturbation of rewards
 178 in bandits, the adversary’s label flipping attack method in our model is quite different. The cost
 179 of obtaining adversarial feedback is uniformly 1, unlike in bandits where the cost depends on the
 180 intensity of the perturbation. Additionally, adversarial feedback in our setting involves comparing two
 181 arms, whereas in bandits it pertains to the reward of a single arm. The only previous work that studied
 182 label-flipping is (Agarwal et al., 2021), where the adversary cannot observe the action selected by the
 183 agent. In contrast, our setting focuses on scenarios where this information is available to adversaries,
 184 which is common in many real-life applications.

185 As the context is changing, the optimal action is different in each round, denoted by $a_t^* =$
 186 $\operatorname{argmax}_{a \in \mathcal{A}} r^*(x_t, a)$. The goal of our algorithm is to minimize the cumulative gap between the
 187 rewards of both selected actions and the optimal action

$$\operatorname{Regret}(T) = \sum_{t=1}^T 2r^*(x_t, a_t^*) - r^*(x_t, a_t) - r^*(x_t, b_t). \quad (3.2)$$

¹Such adversary is referred to as strong adversary (He et al., 2022), compared with the weak adversary who cannot obtain the information before the decision.

188 This regret definition is the same as that in Saha (2021) and the average regret defined in Bengs et al.
 189 (2022). It is typically stronger than weak regret defined in Bengs et al. (2022), which only considers
 190 the reward gap of the better action.

191 **4 Algorithm**

192 In this section, we present our new algorithm RCDB, designed for learning contextual linear dueling
 193 bandits. The main algorithm is illustrated in Algorithm 1. At a high level, we incorporate uncertainty-
 194 dependent weighting into the Maximum Likelihood Estimator (MLE) to counter adversarial feedback.
 195 Specifically, in each round $t \in [T]$, we construct the estimator of parameter θ by solving the following
 196 equation:

$$\lambda\kappa\theta + \sum_{i=1}^{t-1} w_i (\sigma(\phi_i^\top \theta) - o_i) \phi_i = \mathbf{0}, \quad (4.1)$$

197 where we denote $\phi_i = \phi(x_i, a_i) - \phi(x_i, b_i)$ for simplicity, w_i is the uncertainty weight we are going
 198 to choose. To obtain an intuitive understanding of our weight, we consider any action-observation
 199 sequence $(x_1, a_1, b_1, o_1, x_2, a_2, b_2, o_2, \dots, x_t, a_t, b_t, o_t)$ up to round t . For simplicity, we denote
 200 $\mathcal{F}_t = \sigma(x_1, a_1, b_1, o_1, x_2, a_2, b_2, o_2, \dots, x_t, a_t, b_t)$ as the filtration. Suppose the estimated parameter
 201 θ_t is the solution to the unweighted version equation of (4.1), i.e.,

$$\lambda\kappa\theta_t + \sum_{i=1}^t (\sigma(\phi_i^\top \theta_t) - o_i) \phi_i = \mathbf{0}. \quad (4.2)$$

202 When we receive $\phi_t = \phi(x_t, a_t) - \phi(x_t, b_t)$, the probability of receiving $l_t = 1$ can be estimated
 203 by $\sigma(\phi_t^\top \theta_t)$. We consider the conditional variance of the estimated probability $\sigma(\phi_t^\top \theta_t)$ in round t ,
 204 i.e., $\text{Var}[\sigma(\phi_t^\top \theta_t) | \mathcal{F}_t]$, involving a posterior estimate of the prediction's variance. First, we have

$$\begin{aligned} \mathbb{E}[\sigma(\phi_t^\top \theta_t) | \mathcal{F}_t] &\approx \mathbb{E}[\sigma(\phi_t^\top \theta^*) + \sigma'(\phi_t^\top \theta^*) \phi_t^\top (\theta_t - \theta^*) | \mathcal{F}_t] \\ &= \mathbb{E}[\underbrace{\sigma(\phi_t^\top \theta^*) - \sigma'(\phi_t^\top \theta^*) \phi_t^\top \theta^*}_{\mathcal{F}_t\text{-measurable}} | \mathcal{F}_t] + \mathbb{E}[\sigma'(\phi_t^\top \theta^*) \phi_t^\top \theta_t | \mathcal{F}_t]. \end{aligned}$$

205 Moreover, using the Taylor's expansion to (4.2), we have

$$\begin{aligned} \mathbf{0} &= \lambda\kappa\theta_t + \sum_{i=1}^t (\sigma(\phi_i^\top \theta_t) - o_i) \phi_i \\ &\approx \left(\lambda\kappa\mathbf{I} + \sum_{i=1}^t \sigma'(\phi_i^\top \theta^*) \phi_i \phi_i^\top \right) \theta_t + \sum_{i=1}^t (\sigma(\phi_i^\top \theta^*) - o_i) \phi_i - \sum_{i=1}^t \sigma'(\phi_i^\top \theta^*) \phi_i \phi_i^\top \theta^*. \end{aligned}$$

206 Let $\Lambda_t = \lambda\kappa\mathbf{I} + \sum_{i=1}^t \sigma'(\phi_i^\top \theta^*) \phi_i \phi_i^\top$, we have

$$\begin{aligned} \theta_t &\approx \Lambda_t^{-1} \left[\sum_{i=1}^t \sigma'(\phi_i^\top \theta^*) \phi_i \phi_i^\top \theta^* - \sum_{i=1}^t (\sigma(\phi_i^\top \theta^*) - o_i) \phi_i \right] \\ &= \underbrace{\Lambda_t^{-1} \left[\sum_{i=1}^t \sigma'(\phi_i^\top \theta^*) \phi_i \phi_i^\top \theta^* - \sum_{i=1}^{t-1} (\sigma(\phi_i^\top \theta^*) - o_i) \phi_i - \sigma(\phi_t^\top \theta^*) \right]}_{\mathcal{F}_t\text{-measurable}} + o_t \Lambda_t^{-1} \phi_t \end{aligned}$$

207 Therefore, the variance of the estimated preference probability can be approximated by

$$\begin{aligned} \text{Var}[\sigma(\phi_t^\top \theta_t) | \mathcal{F}_t] &= \mathbb{E}[(\sigma(\phi_t^\top \theta_t) - \mathbb{E}[\sigma(\phi_t^\top \theta_t) | \mathcal{F}_t])^2 | \mathcal{F}_t] \\ &\approx \mathbb{E} \left[\left(\mathbb{E}[o_t \sigma'(\phi_t^\top \theta^*) \phi_t^\top \Lambda_t^{-1} \phi_t | \mathcal{F}_t] \right)^2 \middle| \mathcal{F}_t \right] \\ &\leq \mathbb{E}[o_t [\sigma'(\phi_t^\top \theta^*)]^2 \|\phi_t\|_{\Lambda_t^{-1}}^2 | \mathcal{F}_t] \leq [\sigma'(\phi_t^\top \theta^*)]^2 \|\phi_t\|_{\Lambda_t^{-1}}^2, \end{aligned}$$

208 where the first inequality holds due to the Jensen's inequality and $o_t^2 = o_t$, and the last inequality
 209 holds due to $\mathbb{E}[o_t | \mathcal{F}_t] \leq 1$. Using $\sigma'(\phi_t^\top \theta^*) \leq 1$, $\phi_t^\top \theta^* \leq 1$, $\Lambda_t \geq \kappa \Sigma_{t+1} \geq \kappa \Sigma_t$, we can see that
 210 $\text{Var}[\sigma(\phi_t^\top \theta_t) | \mathcal{F}_t] \leq \kappa^{-1} \|\phi_t\|_{\Sigma_t^{-1}}^2$. Since higher variance leads to larger uncertainty, which harms
 211 the credibility of the data, it is natural to assign a smaller weight to the data with high uncertainty.
 212 Thus, we choose the weight to cancel out the uncertainty as follows

$$w_i = \min\{1, \alpha / \|\phi_i\|_{\Sigma_i^{-1}}\}, \quad (4.3)$$

213 where $\alpha / \|\phi_i\|_{\Sigma_i^{-1}}$ normalizes the variance of the estimated probability. To prevent excessively
 214 large weights, we apply truncation to this value. A similar weight has been used in He et al. (2022)
 215 for linear contextual bandits under corruption. Different from their setting where the weight is an
 216 estimate of the variance of the linear model, our weight is an estimate of a generalized linear model.

217 Furthermore, by selecting a proper threshold parameter, e.g., $\alpha = \sqrt{d}/C$, the weighted MLE shares
 218 the same confidence radius with that of the no-adversary scenario.
 219 After constructing the estimator θ_t from the weighted MLE, the sum of the estimated reward for
 220 each duel (a, b) can be calculated as $(\phi(x_t, a) + \phi(x_t, b))^\top \theta_t$. To encourage the exploration of duel
 221 (a, b) with high uncertainty during the learning process, we introduce an exploration bonus with the
 222 following $\beta \|\phi(x_t, a) - \phi(x_t, b)\|_{\Sigma_t^{-1}}$, which follows a similar spirit to the bonus term in the context
 223 of linear bandit problems (Abbasi-Yadkori et al., 2011). However, the reward term and the bonus term
 224 exhibit different combinations of the feature maps $\phi(x_t, a)$ and $\phi(x_t, b)$, which is the key difference
 225 between bandits and dueling bandits. The selection of action pairs (a, b) is subsequently determined
 226 by maximizing the estimated reward with the exploration bonus term, i.e.,

$$(\phi(x_t, a) + \phi(x_t, b))^\top \theta_t + \beta \|\phi(x_t, a) - \phi(x_t, b)\|_{\Sigma_t^{-1}}.$$

227 More discussion about the selection rule was discussed in Appendix A of Di et al. (2023).

Algorithm 1 Robust Contextual Dueling Bandit (RCDB)

1: **Require:** $\alpha > 0$, Regularization parameter λ , confidence radius β .

2: **for** $t = 1, \dots, T$ **do**

3: Compute $\Sigma_t = \lambda \mathbf{I} + \sum_{i=1}^{t-1} w_i (\phi(x_i, a_i) - \phi(x_i, b_i)) (\phi(x_i, a_i) - \phi(x_i, b_i))^\top$.

4: Calculate the MLE θ_t by solving the following equation:

$$\lambda \kappa \theta + \sum_{i=1}^{t-1} w_i \left[\sigma \left((\phi(x_i, a_i) - \phi(x_i, b_i))^\top \theta \right) - o_i \right] (\phi(x_i, a_i) - \phi(x_i, b_i)) = \mathbf{0}. \quad (4.4)$$

5: Observe the context vector x_t .

6: Choose $a_t, b_t = \operatorname{argmax}_{a,b} \left\{ (\phi(x_t, a) + \phi(x_t, b))^\top \theta_t + \beta \|\phi(x_t, a) - \phi(x_t, b)\|_{\Sigma_t^{-1}} \right\}$.

7: The adversary sees the feedback $l_t = \mathbf{1}(a_t \succ b_t)$ and decides the indicator c_t . Observe $o_t = l_t$
 when $c_t = 0$, otherwise observe $o_t = 1 - l_t$.

8: Set weight w_t as (4.3).

9: **end for**

228 5 Main Results

229 5.1 Known Number of Adversarial Feedback

230 At the center of our algorithm design is the uncertainty-weighted MLE. When faced with adversarial
 231 feedback, the estimation error of the weighted MLE θ_t can be characterized by the following lemma.

232 **Lemma 5.1.** If we set $\beta = \sqrt{\lambda}B + (\alpha C + \sqrt{d \log((1 + 2T/\lambda)/\delta)})/\kappa$, then with probability at
 233 least $1 - \delta$, for any $t \in [T]$, we have

$$\|\theta_t - \theta^*\|_{\Sigma_t} \leq \beta.$$

234 **Remark 5.2.** If we set $\alpha = (\sqrt{d} + \sqrt{\lambda}B)/C$, then the bonus radius β has no direct dependency on
 235 the number of adversarial feedback C . This observation plays a key role in proving the adversarial
 236 term in the regret without polynomial dependence on the total number of rounds T .

237 With Lemma 5.1, we can present the following regret guarantee of our algorithm RCDB in the dueling
 238 bandit framework.

239 **Theorem 5.3.** Under Assumption 3.1 and 3.2, let $0 < \delta < 1$, the total number of adversarial feedback
 240 be C . If we set the bonus radius to be

$$\beta = \sqrt{\lambda}B + (\alpha C + \sqrt{d \log((1 + 2T/\lambda)/\delta)})/\kappa,$$

241 then with probability at least $1 - \delta$, the regret in the first t rounds can be upper bounded by

$$\begin{aligned} \text{Regret}(T) &\leq 4(\sqrt{\lambda}B + \alpha C/\kappa) \sqrt{dT \log(1 + 2T/\lambda)} \\ &\quad + 4d(\sqrt{T}/\kappa + \sqrt{\lambda}B/\alpha + 4C/\kappa) \log((1 + 2T/\lambda)/\delta) \\ &\quad + 4d^{1.5} \sqrt{\log^3((1 + 2T/\lambda)/\delta)} / (\alpha \kappa). \end{aligned}$$

242 Moreover, if we set $\alpha = (\sqrt{d} + \sqrt{\lambda}B)/C$, $\lambda = 1/B^2$, the regret upper bound can be simplified to

$$\text{Regret}(T) = \tilde{O}(d\sqrt{T}/\kappa + dC/\kappa).$$

243 **Remark 5.4.** Our regret bound consists of two terms. The first one is a C -independent term $\tilde{O}(d\sqrt{T})$,
 244 which matches the lower bound $\tilde{\Omega}(d\sqrt{T})$ proved in Bengs et al. (2022). This indicates that our result
 245 is optimal in scenarios without adversarial feedback ($C = 0$). Additionally, our result includes an
 246 additive term that is linearly dependent on the number of adversarial feedback C . When $C = O(\sqrt{T})$,
 247 the order of regret will be the same as the stochastic setting. It indicates the robustness of our algorithm
 248 to adversarial feedback. Additionally, the following theorem we present establishes a lower bound
 249 for this adversarial term, indicating that our dependency on the number of adversarial feedback C
 250 and the context dimension d is also optimal.

251 **Theorem 5.5.** For any dimension d , there exists an instance of dueling bandits with $|\mathcal{A}| = d$, such
 252 that any algorithm with the knowledge of the number of adversarial feedback C must incur $\Omega(dC)$
 253 regret with probability at least $1/2$.

254 **Remark 5.6.** The proof of Theorem 5.5 follows Bogunovic et al. (2021). In the constructed instances,
 255 only one action has reward 1, while others have 0. Compared with linear bandits, where the feedback
 256 is an exact reward, dueling bandits deal with the comparison between a pair of actions. A critical
 257 observation from our preference model, as formulated in (3.1), is that two actions with identical
 258 rewards result in a pair that is challenging to differentiate. The lower bound can be proved by
 259 corrupting every comparison into a random guess until the total times of adversarial feedback have
 260 been used up. For detailed proof, please refer to Section B.2. Our proved lower bound $\Omega(dC)$ shows
 261 that our result is nearly optimal because of the linear dependency on C, d and only logarithmic
 262 dependency on the total number of rounds T .

263 5.2 Unknown Number of Adversarial Feedback

264 In our previous analysis, the selection of parameters depends on having prior knowledge of the total
 265 number of adversarial feedback C . In this subsection, we extend our previous result to address
 266 the challenge posed by an unknown number of adversarial feedback C . Our approach to tackle
 267 this uncertainty follows He et al. (2022), we introduce an adversarial tolerance threshold \bar{C} for the
 268 adversary count. This threshold can be regarded as an optimistic estimator of the actual number of
 269 adversarial feedback C . Under this situation, the subsequent theorem provides an upper bound for
 270 regret of Algorithm 1 in the case of an unknown number of adversarial feedback C .

271 **Theorem 5.7.** Under Assumptions 3.1 and 3.2, if we set the the confidence radius as

$$\beta = \sqrt{\lambda}B + [\alpha\bar{C} + \sqrt{d \log((1 + 2T/\lambda)/\delta)}] / \kappa,$$

272 with the pre-defined adversarial tolerance threshold \bar{C} and $\alpha = (\sqrt{d} + \sqrt{\lambda}B) / \bar{C}$, then with probability
 273 at least $1 - \delta$, the regret of Algorithm 1 can be upper bounded as following:

- 274 • If the actual number of adversarial feedback C is smaller than the adversarial tolerance threshold
 275 \bar{C} , then we have

$$\text{Regret}(T) = \tilde{O}(d\sqrt{T}/\kappa + d\bar{C}/\kappa).$$

- 276 • If the actual number of adversarial feedback C is larger than the adversarial tolerance threshold \bar{C} ,
 277 then we have $\text{Regret}(T) = O(T)$.

278 **Remark 5.8.** The COBE framework (Wei et al., 2022) converts any algorithm with the known
 279 adversarial level to an algorithm in the unknown case. However, such a framework only works for
 280 weak adversaries and does not work in our strong adversary setting. In fact, He et al. (2022) proved
 281 that any algorithm cannot simultaneously achieve near-optimal regret when uncorrupted and maintain
 282 sublinear regret with corruption level $C = \Omega(\sqrt{T})$. Therefore, there exists a trade-off between robust
 283 adversarial defense and near-optimal algorithmic performance. Our algorithm achieves the same
 284 nearly optimal $\tilde{O}(d\sqrt{T})$ regret as the no-adversary case even when $C = \Theta(\sqrt{T})$, which indicates
 285 that our results are optimal in the presence of an unknown number of adversarial feedback.

286 6 Experiments

287 6.1 Experiment Setup

288 **Preference Model.** We study the effect of adversarial feedback with the preference model deter-
 289 mined by (3.1), where $\sigma(x) = 1/(1 + e^{-x})$. We randomly generate the underlying parameter in
 290 $[-0.5, 0.5]^d$ and normalize it to be a vector with $\|\theta^*\|_2 = 2$. Then, we set it to be the underlying
 291 parameter and construct the reward utilized in the preference model as $r^*(x, a) = \langle \theta^*, \phi(x, a) \rangle$.
 292 We set the action set $\mathcal{A} = \{-1/\sqrt{d}, 1/\sqrt{d}\}^d$. For simplicity, we assume $\phi(x, a) = a$. In our
 293 experiment, we set the dimension $d = 5$, with the size of action set $|\mathcal{A}| = 2^d = 32$.

294 **Adversarial Attack Methods.** We study the performance of our algorithm using different adversarial attack methods. We categorize the first two methods as “weak” primarily because the adversary in these scenarios does not utilize information about the agent’s actions. In contrast, we classify the latter two methods as “strong” attacks. In these cases, the adversary leverages a broader scope of information, including knowledge of the actions selected by the agent and the true preference model. This enables it to devise more targeted adversarial methods.

- 300 • “Greedy Attack”: The adversary will flip the preference label for the first C rounds. After that, it will not corrupt the result anymore.
- 302 • “Random Attack”: In each round, the adversary will flip the preference label with the probability of $0 < p < 1$, until the times of adversarial feedback reach C .
- 304 • “Adversarial Attack”: The adversary can have access to the true preference model. It will only flip the preference label when it aligns with the preference model, i.e., the probability for the preference model to make that decision is larger than 0.5, until the times of adversarial feedback reach C .
- 307 • “Misleading Attack”: The adversary selects a suboptimal action. It will make sure this arm is always the winner in the comparison until the times of adversarial feedback reach C . In this way, it will mislead the agent to believe this action is the optimal one.

310 **Experiment Setup.** For each experiment instance, we simulate the interaction with the environment for $T = 2000$ rounds. In each round, the feedback for the action pair selected by the algorithm is generated according to the defined preference model. Subsequently, the adversary observes both the selected actions and their corresponding feedback and then engages in one of the previously described adversarial attack methods. We report the regret defined in (3.2) averaged across 10 random runs.

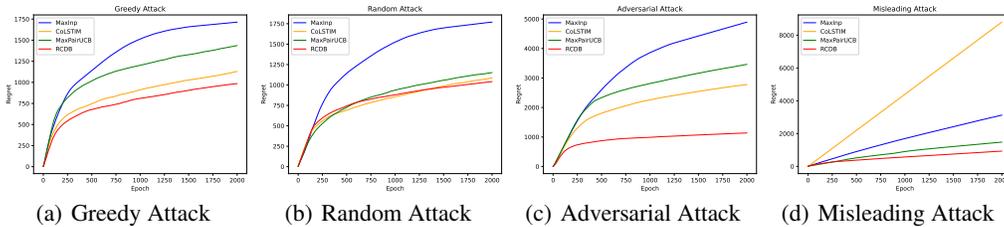


Figure 1: Comparison of RCDB (Our Algorithm 1), MaxInP (Saha, 2021), CoLSTIM (Bengs et al., 2022) and MaxPairUCB (Di et al., 2023). We report the cumulative regret with various adversarial attack methods (Greedy, Random, Adversarial, Misleading). For the baselines, the parameters are carefully tuned to achieve better results with different attack methods. The total number of adversarial feedback is $C = \lceil \sqrt{T} \rceil$.

314

315 6.2 Performance Comparison

316

We first introduce the algorithms studied in this section.

317

• **MaxInP:** Maximum Informative Pair by Saha (2021). It involves maintaining a standard MLE. With the estimated model, it then identifies a set of promising arms possible to beat the rest. The selection of arm pairs is then strategically designed to maximize the uncertainty in the difference between the two arms within this promising set, referred to as “maximum informative”.

318

319

320

321

• **CoLSTIM:** The method by Bengs et al. (2022). It involves maintaining a standard MLE for the estimated model. Based on this model, the first arm is selected as the one with the highest estimated reward, implying it is the most likely to prevail over competitors. The second arm is selected to be the first arm’s toughest competitor, with an added uncertainty bonus.

322

323

324

325

• **MaxPairUCB:** This algorithm was proposed in Di et al. (2023). It uses the regularized MLE to estimate the parameter θ^* . Then it selects the actions based on a symmetric action selection rule, i.e. the actions with the largest estimated reward plus some uncertainty bonus.

326

327

328

• **RCDB:** Algorithm 1 proposed in this paper. The key difference from the other algorithms is the use of uncertainty weight in the calculation of MLE (4.4). The we use the same symmetric action selection rule as MaxPairUCB. Our experiment results show that the uncertainty weight is critical in the face of adversarial feedback.

329

330

331

332

333

334

Our results are demonstrated in Figure 1. In Figure 1(a) and Figure 1(b), we observe scenarios where the adversary is “weak” due to the lack of access to information regarding the selected actions and the underlying preference model. Notably, in these situations, our algorithm RCDB outperforms all other

335 baseline algorithms, demonstrating its robustness. Among the other algorithms, CoLSTIM performs
 336 as the strongest competitor.

337 In Figure 1(c), the adversary employs a 'stronger' adversarial method. Due to the inherent randomness
 338 of the model, some labels may naturally be 'incorrect'. An adversary with knowledge of the selected
 339 actions and the preference model can strategically neglect these naturally incorrect labels and
 340 selectively flip the others. This method proves catastrophic for algorithms to learn the true model,
 341 as it results in the agent encountering only incorrect preference labels at the beginning. Our results
 342 indicate that this leads to significantly higher regret. However, it's noteworthy that our algorithm
 343 RCDB demonstrates considerable robustness.

344 In Figure 1(d), the adversary employs a strategy aimed at misleading algorithms into believing a
 345 suboptimal action is the best choice. The algorithm CoLSTIM appears to be the most susceptible to
 346 being cheated by this method. Despite the deployment of 'strong' adversarial methods, as shown
 347 in both Figure 1(c) and Figure 1(d), our algorithm, RCDB, consistently demonstrates exceptional
 348 robustness against these attacks. A significant advantage of RCDB lies in that our parameter is selected
 349 solely based on the number of adversarial feedback C , irrespective of the nature of the adversarial
 350 methods employed. This contrasts with other algorithms where parameter tuning must be specifically
 adapted for each distinct adversarial method.

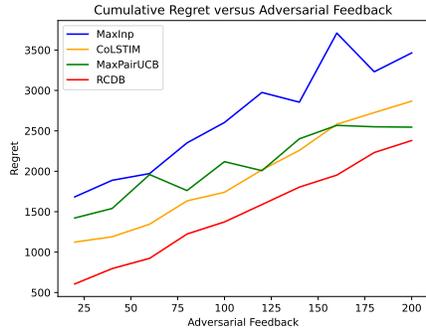


Figure 2: The relationship between cumulative regret and the number of adversarial feedback C . For this specific experiment, we employ the "greedy attack" method to generate the adversarial feedback. C is selected from the set [20, 40, 60, 80, 100, 120, 140, 160, 180, 200] (10 adversarial levels).

351

352 6.3 Robustness to Different Numbers of Adversarial Feedback

353 In this section, we test the performance of algorithms with increasing times of adversarial feedback.
 354 Our results show a linear dependency on the number of adversarial feedback C , which is consistent
 355 with the theoretical results we have proved in Theorem 5.3 and 5.5. In comparison to other algorithms,
 356 RCDB demonstrates superior robustness against adversarial feedback, as evidenced by its notably
 357 smaller regret.

358 7 Conclusion

359 In this paper, we focus on the contextual dueling bandit problem from adversarial feedback. We
 360 introduce a novel algorithm, RCDB, which utilizes an uncertainty-weighted Maximum Likelihood
 361 Estimator (MLE) approach. This algorithm not only achieves optimal theoretical results in scenarios
 362 with and without adversarial feedback but also demonstrates superior performance with synthetic
 363 data. For future direction, we aim to extend our uncertainty-weighted method to encompass more
 364 general settings involving preference-based data. A particularly promising future direction of our
 365 research lies in addressing adversarial feedback within the process of aligning large language models
 366 using Reinforcement Learning from Human Feedback (RLHF).

367 **Limitations.** We assume that the reward is linear with respect to some known feature maps. Although
 368 this setting is common in the literature, we observe that some recent works on dueling bandits can
 369 deal with nonlinear rewards (Li et al., 2024). Therefore, it's possible to extend our results to a more
 370 general setting. Another assumption concerns the lower bound of the derivative of the link function.
 371 Notably, in the logistic bandit model, which shares similarities with our setting through Bernoulli
 372 variables, some work (Abeille et al., 2021; Fauray et al., 2022) can improve the dependency of κ from
 373 $1/\kappa$ to $\sqrt{\kappa}$. A similar improvement might be achieved in our setting as well.

374 **References**

- 375 ABBASI-YADKORI, Y., PÁL, D. and SZEPESVÁRI, C. (2011). Improved algorithms for linear
376 stochastic bandits. In *Advances in Neural Information Processing Systems*.
- 377 ABEILLE, M., FAURY, L. and CALAUZÈNES, C. (2021). Instance-wise minimax-optimal algorithms
378 for logistic bandits. In *International Conference on Artificial Intelligence and Statistics*. PMLR.
- 379 AGARWAL, A., AGARWAL, S. and PATIL, P. (2021). Stochastic dueling bandits with adversarial
380 corruption. In *Algorithmic Learning Theory*. PMLR.
- 381 AGRAWAL, S. and GOYAL, N. (2012). Analysis of thompson sampling for the multi-armed bandit
382 problem. In *Conference on learning theory*. JMLR Workshop and Conference Proceedings.
- 383 AILON, N., KARNIN, Z. and JOACHIMS, T. (2014). Reducing dueling bandits to cardinal bandits.
384 In *International Conference on Machine Learning*. PMLR.
- 385 AUER, P. (2002). Using confidence bounds for exploitation-exploration trade-offs. *Journal of*
386 *Machine Learning Research* **3** 397–422.
- 387 AUER, P., CESA-BIANCHI, N. and FISCHER, P. (2002a). Finite-time analysis of the multiarmed
388 bandit problem. *Machine Learning* **47** 235–256.
- 389 AUER, P., CESA-BIANCHI, N., FREUND, Y. and SCHAPIRE, R. E. (2002b). The nonstochastic
390 multiarmed bandit problem. *SIAM journal on computing* **32** 48–77.
- 391 AUER, P. and CHIANG, C.-K. (2016). An algorithm with nearly optimal pseudo-regret for both
392 stochastic and adversarial bandits. In *Conference on Learning Theory*. PMLR.
- 393 BALSUBRAMANI, A., KARNIN, Z., SCHAPIRE, R. E. and ZOGHI, M. (2016). Instance-dependent
394 regret bounds for dueling bandits. In *Conference on Learning Theory*. PMLR.
- 395 BENGIS, V., SAHA, A. and HÜLLERMEIER, E. (2022). Stochastic contextual dueling bandits under
396 linear stochastic transitivity models. In *International Conference on Machine Learning*. PMLR.
- 397 BOGUNOVIC, I., LOSALKA, A., KRAUSE, A. and SCARLETT, J. (2021). Stochastic linear bandits
398 robust to adversarial attacks. In *International Conference on Artificial Intelligence and Statistics*.
399 PMLR.
- 400 BUBECK, S., CESA-BIANCHI, N. ET AL. (2012). Regret analysis of stochastic and nonstochastic
401 multi-armed bandit problems. *Foundations and Trends® in Machine Learning* **5** 1–122.
- 402 BUBECK, S. and SLIVKINS, A. (2012). The best of both worlds: Stochastic and adversarial bandits.
403 In *Conference on Learning Theory*. JMLR Workshop and Conference Proceedings.
- 404 CESA-BIANCHI, N. and LUGOSI, G. (2006). *Prediction, learning, and games*. Cambridge university
405 press.
- 406 DI, Q., JIN, T., WU, Y., ZHAO, H., FARNOUD, F. and GU, Q. (2023). Variance-aware regret bounds
407 for stochastic contextual dueling bandits. *arXiv preprint arXiv:2310.00968* .
- 408 DING, Q., HSIEH, C.-J. and SHARPNACK, J. (2022). Robust stochastic linear contextual bandits
409 under adversarial attacks. In *International Conference on Artificial Intelligence and Statistics*.
410 PMLR.
- 411 DUDÍK, M., HOFMANN, K., SCHAPIRE, R. E., SLIVKINS, A. and ZOGHI, M. (2015). Contextual
412 dueling bandits. In *Conference on Learning Theory*. PMLR.
- 413 FALAHATGAR, M., HAO, Y., ORLITSKY, A., PICHAPATI, V. and RAVINDRAKUMAR, V. (2017).
414 Maxing and ranking with few assumptions. *Advances in Neural Information Processing Systems*
415 **30**.
- 416 FAURY, L., ABEILLE, M., CALAUZÈNES, C. and FERCOQ, O. (2020). Improved optimistic
417 algorithms for logistic bandits. In *International Conference on Machine Learning*. PMLR.

- 418 FAURY, L., ABEILLE, M., JUN, K.-S. and CALAUZÈNES, C. (2022). Jointly efficient and optimal
419 algorithms for logistic bandits. In *International Conference on Artificial Intelligence and Statistics*.
420 PMLR.
- 421 FILIPPI, S., CAPPE, O., GARIVIER, A. and SZEPESVÁRI, C. (2010). Parametric bandits: The
422 generalized linear case. *Advances in Neural Information Processing Systems* **23**.
- 423 GAJANE, P., URVOY, T. and CLÉROT, F. (2015). A relative exponential weighing algorithm for
424 adversarial utility-based dueling bandits. In *International Conference on Machine Learning*.
425 PMLR.
- 426 GUPTA, A., KOREN, T. and TALWAR, K. (2019). Better algorithms for stochastic bandits with
427 adversarial corruptions. In *Conference on Learning Theory*. PMLR.
- 428 HE, J., ZHOU, D., ZHANG, T. and GU, Q. (2022). Nearly optimal algorithms for linear contextual
429 bandits with adversarial corruptions. *Advances in Neural Information Processing Systems* **35**
430 34614–34625.
- 431 HECKEL, R., SIMCHOWITZ, M., RAMCHANDRAN, K. and WAINWRIGHT, M. (2018). Approximate
432 ranking from pairwise comparisons. In *International Conference on Artificial Intelligence and*
433 *Statistics*. PMLR.
- 434 JAMIESON, K., KATARIYA, S., DESHPANDE, A. and NOWAK, R. (2015). Sparse dueling bandits.
435 In *Artificial Intelligence and Statistics*. PMLR.
- 436 KALYANAKRISHNAN, S., TEWARI, A., AUER, P. and STONE, P. (2012). Pac subset selection in
437 stochastic multi-armed bandits. In *ICML*, vol. 12.
- 438 KOMIYAMA, J., HONDA, J., KASHIMA, H. and NAKAGAWA, H. (2015). Regret lower bound and
439 optimal algorithm in dueling bandit problem. In *Conference on learning theory*. PMLR.
- 440 KOMIYAMA, J., HONDA, J. and NAKAGAWA, H. (2016). Copeland dueling bandit problem:
441 Regret lower bound, optimal algorithm, and computationally efficient algorithm. In *International*
442 *Conference on Machine Learning*. PMLR.
- 443 KUROKI, Y., RUMI, A., TSUCHIYA, T., VITALE, F. and CESA-BIANCHI, N. (2023). Best-of-both-
444 worlds algorithms for linear contextual bandits. *arXiv preprint arXiv:2312.15433* .
- 445 LAI, T. L. (1987). Adaptive treatment allocation and the multi-armed bandit problem. *The annals of*
446 *statistics* 1091–1114.
- 447 LAI, T. L., ROBBINS, H. ET AL. (1985). Asymptotically efficient adaptive allocation rules. *Advances*
448 *in applied mathematics* **6** 4–22.
- 449 LATTIMORE, T. and SZEPESVÁRI, C. (2020). *Bandit Algorithms*. Cambridge University Press.
- 450 LEE, C.-W., LUO, H., WEI, C.-Y., ZHANG, M. and ZHANG, X. (2021). Achieving near instance-
451 optimality and minimax-optimality in stochastic and adversarial linear bandits simultaneously. In
452 *International Conference on Machine Learning*. PMLR.
- 453 LI, L., LU, Y. and ZHOU, D. (2017). Provably optimal algorithms for generalized linear contextual
454 bandits. In *International Conference on Machine Learning*. PMLR.
- 455 LI, X., ZHAO, H. and GU, Q. (2024). Feel-good thompson sampling for contextual dueling bandits.
456 *arXiv preprint arXiv:2404.06013* .
- 457 LI, Y., LOU, E. Y. and SHAN, L. (2019). Stochastic linear optimization with adversarial corruption.
458 *arXiv preprint arXiv:1909.02109* .
- 459 LYKOURIS, T., MIRROKNI, V. and PAES LEME, R. (2018). Stochastic bandits robust to adversarial
460 corruptions. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*.
- 461 OUYANG, L., WU, J., JIANG, X., ALMEIDA, D., WAINWRIGHT, C., MISHKIN, P., ZHANG,
462 C., AGARWAL, S., SLAMA, K., RAY, A. ET AL. (2022). Training language models to follow
463 instructions with human feedback. *Advances in Neural Information Processing Systems* **35** 27730–
464 27744.

- 465 RAMAMOCHAN, S. Y., RAJKUMAR, A. and AGARWAL, S. (2016). Dueling bandits: Beyond
466 condorcet winners to general tournament solutions. *Advances in Neural Information Processing*
467 *Systems* **29**.
- 468 SAHA, A. (2021). Optimal algorithms for stochastic contextual preference bandits. *Advances in*
469 *Neural Information Processing Systems* **34** 30050–30062.
- 470 SAHA, A. and GAILLARD, P. (2022). Versatile dueling bandits: Best-of-both world analyses for
471 learning from relative preferences. In *International Conference on Machine Learning*. PMLR.
- 472 SAHA, A., KOREN, T. and MANSOUR, Y. (2021). Adversarial dueling bandits. In *International*
473 *Conference on Machine Learning*. PMLR.
- 474 SAHA, A. and KRISHNAMURTHY, A. (2022). Efficient and optimal algorithms for contextual dueling
475 bandits under realizability. In *International Conference on Algorithmic Learning Theory*. PMLR.
- 476 SEKHARI, A., SRIDHARAN, K., SUN, W. and WU, R. (2023). Contextual bandits and imitation
477 learning via preference-based active queries. *arXiv preprint arXiv:2307.12926* .
- 478 SELDIN, Y. and LUGOSI, G. (2017). An improved parametrization and analysis of the exp3++
479 algorithm for stochastic and adversarial bandits. In *Conference on Learning Theory*. PMLR.
- 480 SELDIN, Y. and SLIVKINS, A. (2014). One practical algorithm for both stochastic and adversarial
481 bandits. In *International Conference on Machine Learning*. PMLR.
- 482 WEI, C.-Y., DANN, C. and ZIMMERT, J. (2022). A model selection approach for corruption robust
483 reinforcement learning. In *International Conference on Algorithmic Learning Theory*. PMLR.
- 484 WU, H. and LIU, X. (2016). Double thompson sampling for dueling bandits. *Advances in neural*
485 *information processing systems* **29**.
- 486 WU, Y., JIN, T., LOU, H., FARNOUD, F. and GU, Q. (2023). Borda regret minimization for
487 generalized linear dueling bandits. *arXiv preprint arXiv:2303.08816* .
- 488 XIONG, W., DONG, H., YE, C., ZHONG, H., JIANG, N. and ZHANG, T. (2023). Gibbs sam-
489 pling from human feedback: A provable kl-constrained framework for rlhf. *arXiv preprint*
490 *arXiv:2312.11456* .
- 491 YE, C., XIONG, W., GU, Q. and ZHANG, T. (2023). Corruption-robust algorithms with uncertainty
492 weighting for nonlinear contextual bandits and markov decision processes. In *International*
493 *Conference on Machine Learning*. PMLR.
- 494 YU, T., QUILLEN, D., HE, Z., JULIAN, R., HAUSMAN, K., FINN, C. and LEVINE, S. (2020).
495 Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In
496 *Conference on robot learning*. PMLR.
- 497 YUE, Y., BRODER, J., KLEINBERG, R. and JOACHIMS, T. (2012). The k-armed dueling bandits
498 problem. *Journal of Computer and System Sciences* **78** 1538–1556.
- 499 ZHAO, H., ZHOU, D. and GU, Q. (2021). Linear contextual bandits with adversarial corruptions.
500 *arXiv preprint arXiv:2110.12615* .
- 501 ZHU, H., YU, J., GUPTA, A., SHAH, D., HARTIKAINEN, K., SINGH, A., KUMAR, V. and
502 LEVINE, S. (2020). The ingredients of real-world robotic reinforcement learning. *arXiv preprint*
503 *arXiv:2004.12570* .
- 504 ZIMMERT, J. and SELDIN, Y. (2019). An optimal algorithm for stochastic and adversarial bandits.
505 In *The 22nd International Conference on Artificial Intelligence and Statistics*. PMLR.
- 506 ZOGHI, M., KARNIN, Z. S., WHITESON, S. and DE RIJKE, M. (2015). Copeland dueling bandits.
507 *Advances in neural information processing systems* **28**.
- 508 ZOGHI, M., WHITESON, S., MUNOS, R. and RIJKE, M. (2014). Relative upper confidence bound
509 for the k-armed dueling bandit problem. In *International conference on machine learning*. PMLR.

510 **Broader Impact**

511 This paper studies contextual dueling bandits with adversarial feedback. Our primary objective is
 512 to propel advancements in bandit theory by introducing a more robust algorithm backed by solid
 513 theoretical guarantees. The uncertainty-weighted approach we have developed for dueling bandits
 514 holds significant potential to address the issue of adversarial feedback in preference-based data, which
 515 could be instrumental in enhancing the robustness of generative models against adversarial attacks,
 516 thereby contributing positively to the societal impact and reliability of machine learning applications.

517 **A Roadmap of the Proof**

518 **A.1 Uncertainty-weighted MLE with Adversarial Feedback**

519 In this section, we offer an overview of the proof for Lemma 5.1. The general proof idea for
 520 the uncertainty-weighted MLE with adversarial feedback lies in decomposing the estimation error
 521 into three terms, a stochastic error term, an adversarial term, and an additional regularization term.
 522 Following the analysis of standard (weighted) MLE (Li et al., 2017), we introduce an auxiliary
 523 function:

$$G_t(\boldsymbol{\theta}) = \lambda\kappa\boldsymbol{\theta} + \sum_{i=1}^{t-1} w_i \left[\sigma \left((\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i))^\top \boldsymbol{\theta} \right) - \sigma \left((\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i))^\top \boldsymbol{\theta}^* \right) \right] (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i)).$$

524 It satisfies two conditions: First, for the true parameter value $\boldsymbol{\theta}^*$, $G_t(\boldsymbol{\theta}^*)$ has a simple expression, i.e.,

$$G_t(\boldsymbol{\theta}^*) = \lambda\kappa\boldsymbol{\theta}^*.$$

525 Second, according to (4.4), we can get the value of function G_t at the MLE $\boldsymbol{\theta}_t$,

$$G_t(\boldsymbol{\theta}_t) = \sum_{i=1}^{t-1} w_i \gamma_i (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i)), \quad (\text{A.1})$$

526 where $\gamma_i = o_i - \sigma \left((\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i))^\top \boldsymbol{\theta}^* \right)$. To connect the desired estimation error with the
 527 function G_t , we use the mean value theorem. This leads to an upper bound of the estimation error:

$$\begin{aligned} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|_{\Sigma_t} &\leq \frac{1}{\kappa} \|G_t(\boldsymbol{\theta}_t) - G_t(\boldsymbol{\theta}^*)\|_{\Sigma_t^{-1}} \\ &\leq \underbrace{\frac{1}{\kappa} \lambda \|\boldsymbol{\theta}^*\|_{\Sigma_t^{-1}}}_{\text{Regularization term}} + \underbrace{\frac{1}{\kappa} \|G_t(\boldsymbol{\theta}_t)\|_{\Sigma_t^{-1}}}_{I_1}. \end{aligned}$$

528 For term I_1 , we can decompose the summation in (A.1) based on the adversarial feedback c_t , i.e.,

$$G_t(\boldsymbol{\theta}_t) = \sum_{i < t: c_i = 0} w_i \gamma_i (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i)) + \underbrace{\sum_{i < t: c_i = 1} w_i \gamma_i (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i))}_{I_2},$$

529 where I_2 can be further decomposed as

$$I_2 = \sum_{i < t: c_i = 1} w_i \epsilon_i (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i)) + \sum_{i < t: c_i = 1} w_i (\gamma_i - \epsilon_i) (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i)).$$

530 where $\epsilon_i = l_i - \sigma \left((\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i))^\top \boldsymbol{\theta}^* \right)$. With our notation of adversarial feedback, when
 531 $c_i = 0$, we have $\gamma_i = \epsilon_i$. Therefore, we have $|\gamma_i - \epsilon_i| \leq 1$ and

$$I_1 \leq \underbrace{\frac{1}{\kappa} \left\| \sum_{i=1}^{t-1} w_i \epsilon_i (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i)) \right\|_{\Sigma_t^{-1}}}_{\text{Stochastic term}} + \underbrace{\frac{1}{\kappa} \left\| \sum_{i < t: c_i = 1} w_i (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i)) \right\|_{\Sigma_t^{-1}}}_{\text{Adversarial term}}.$$

532 The stochastic term can be upper bounded with the concentration inequality (Lemma D.2). Addition-
 533 ally, by employing our specifically chosen weight, as (4.3), we can control the adversarial term, with
 534 $w_i \|\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i)\|_{\Sigma_t^{-1}} \leq \alpha$. Therefore, the adversarial term can be bounded by $\alpha C / \kappa$.

535 **A.2 Regret Upper Bound**

536 With a similar discussion of the symmetric arm selection rule to Di et al. (2023), the regret defined in
 537 (3.2) can be bounded by

$$\text{Regret}(T) \leq \sum_{t=1}^T \min \left\{ 4, 2\beta \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \right\}.$$

538 Note that in our selection of weight w_t , it has two possible values. We decompose the summation
 539 based on the two cases separately. We have

$$\begin{aligned} \text{Regret}(T) \leq & \underbrace{\sum_{w_t=1} \min \left\{ 4, 2\beta \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \right\}}_{J_1} \\ & + \underbrace{\sum_{w_t < 1} \min \left\{ 4, 2\beta \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \right\}}_{J_2}. \end{aligned}$$

540 We consider J_1, J_2 separately. For the term J_1 , we define $\Lambda_t = \lambda \mathbf{I} + \sum_{i \leq t-1, w_i=1} (\phi(x_i, a_i) -$
 541 $\phi(x_i, b_i))(\phi(x_i, a_i) - \phi(x_i, b_i))^\top$. Then, we have $\Sigma_t \succeq \Lambda_t$, and therefore

$$\|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \leq \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Lambda_t^{-1}}.$$

542 Using Lemma D.3 with $\mathbf{x}_t = \phi(x_t, a_t) - \phi(x_t, b_t)$, we have

$$J_1 \leq 4\beta \sqrt{dT \log(1 + 2T/\lambda)}. \quad (\text{A.2})$$

543 For term J_2 , we note that $w_t < 1$ implies that $w_t = \alpha / \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}}$. Therefore, we
 544 have

$$J_2 \leq \sum_{t=1}^T \frac{4\beta}{\alpha} \min \left\{ 1, \|\sqrt{w_t}(\phi(x_t, a_t) - \phi(x_t, b_t))\|_{\Sigma_t^{-1}}^2 \right\}.$$

545 Using Lemma D.3 with $\mathbf{x}'_t = \sqrt{w_t}(\phi(x_t, a_t) - \phi(x_t, b_t))$, we have

$$J_2 \leq \frac{4d\beta \log(1 + 2T/\lambda)}{\alpha}. \quad (\text{A.3})$$

546 We conclude the proof of regret by combining (A.2) and (A.3).

547 **B Proof of Theorems in Section 5**

548 **B.1 Proof of Theorem 5.3**

549 In this subsection, we provide the proof of Theorem 5.3. We condition on the high-probability event
 550 in Lemma 5.1

$$\mathcal{E} = \left\{ \|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|_{\Sigma_t} \leq \beta, \forall t \in [T] \right\}.$$

551 Let $r_t = 2r^*(x_t, a_t^*) - r^*(x_t, a_t) - r^*(x_t, b_t)$ be the regret incurred in round t . The following lemma
 552 provides the upper bound of r_t .

553 **Lemma B.1.** Let $0 < \delta < 1$. If we set $\beta = \sqrt{\lambda}B + (\alpha C + \sqrt{d \log((1 + 2T/\lambda)/\delta)})/\kappa$, on event \mathcal{E} ,
 554 the regret of Algorithm 1 incurred in round t can be upper bounded by

$$r_t \leq \min \left\{ 4, 2\beta \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \right\}.$$

555 Moreover, the regret can be upper bounded by

$$\text{Regret}(T) \leq \sum_{t=1}^T \min \left\{ 4, 2\beta \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \right\}.$$

556 With Lemma B.1, we can provide the proof of Theorem 5.3.

557 *Proof of Theorem 5.3.* Using Lemma B.1, the total regret can be upper bounded by

$$\text{Regret}(T) \leq \sum_{t=1}^T \min \left\{ 4, 2\beta \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \right\}.$$

558 Our weight w_t has two possible values. We decompose the summation based on the two cases
559 separately. We have

$$\begin{aligned} \text{Regret}(T) &\leq \underbrace{\sum_{w_t=1} \min \left\{ 4, 2\beta \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \right\}}_{J_1} \\ &\quad + \underbrace{\sum_{w_t < 1} \min \left\{ 4, 2\beta \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \right\}}_{J_2}. \end{aligned}$$

560 For the term J_1 , we consider a partial summation in rounds when $w_t = 1$. Let $\Lambda_t = \lambda \mathbf{I} +$
561 $\sum_{i \leq k-1, w_i=1} (\phi(x_i, a_i) - \phi(x_i, b_i))(\phi(x_i, a_i) - \phi(x_i, b_i))^\top$. Then we have

$$\begin{aligned} J_1 &\leq 4\beta \sum_{t:w_t=1} \min \left\{ 1, \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \right\} \\ &\leq 4\beta \sum_{t:w_t=1} \min \left\{ 1, \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Lambda_t^{-1}} \right\} \\ &\leq 4\beta \sqrt{T \sum_{t:w_t=1} \min \left\{ 1, \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Lambda_t^{-1}}^2 \right\}} \\ &\leq 4\beta \sqrt{dT \log(1 + 2T/\lambda)}, \end{aligned} \tag{B.1}$$

562 where the second inequality holds due to $\Sigma_t \succeq \Lambda_t$. The third inequality holds due to the Cauchy-
563 Schwartz inequality, The last inequality holds due to Lemma D.3.

564 For the term J_2 , the weight in this summation satisfies $w_t < 1$, and therefore $w_t = \alpha / \|\phi(x_t, a_t) -$
565 $\phi(x_t, b_t)\|_{\Sigma_t^{-1}}$. Then we have

$$\begin{aligned} J_2 &= \sum_{w_t < 1} \min \left\{ 4, 2\beta \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} w_t \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} / \alpha \right\} \\ &\leq \sum_{t=1}^T \min \left\{ 4, 2\beta / \alpha \|\sqrt{w_t}(\phi(x_t, a_t) - \phi(x_t, b_t))\|_{\Sigma_t^{-1}}^2 \right\} \\ &\leq \sum_{t=1}^T \frac{4\beta}{\alpha} \min \left\{ 1, \|\sqrt{w_t}(\phi(x_t, a_t) - \phi(x_t, b_t))\|_{\Sigma_t^{-1}}^2 \right\} \\ &\leq \frac{4d\beta \log(1 + 2T/\lambda)}{\alpha}, \end{aligned} \tag{B.2}$$

566 where the first equality holds due to the choice of w_t . The first inequality holds because each term in
567 the summation is positive. The last inequality holds due to Lemma D.3. Combining (B.1) and (B.2),
568 we complete the proof of Theorem 5.3. \square

569 B.2 Proof of Theorem 5.5

570 *Proof of Theorem 5.5.* Our proof adapts the argument in Bogunovic et al. (2021) to dueling bandits.
571 For any dimension d , we construct d instances, each with $\theta_i = \mathbf{e}_i$, where \mathbf{e}_i is the i -th standard basis
572 vector. We set the action set $\mathcal{A} = \{\mathbf{e}_i\}_{i=1}^d$. Therefore, in the i -th instance, the reward for the i -th
573 action will be 1. For the other actions, it will be 0. Therefore, the i -th action will be more preferable
574 to any other action. While for other pairs, the feedback is simply a random guess.
575 Consider an adversary that knows the exact instance. When the comparison involves the i -th action,
576 it will corrupt the feedback with a random guess. Otherwise, it will not corrupt. In the i -th instance,
577 the adversary stops the adversarial attack only after C times of comparison involving the i -th action.
578 However, after $Cd/4$ rounds, at least $d/2$ actions have not been compared for C times. For the
579 instances corresponding to these actions, the agent learns no information and suffers from $\Omega(dC)$
580 regret. This completes the proof of Theorem 5.5. \square

581 **B.3 Proof of Theorem 5.7**

582 *Proof of Theorem 5.7.* Here, based on the relationship between C and the threshold \bar{C} , we discuss
583 two distinct cases separately.

- 584 • In the scenario where $\bar{C} < C$, Algorithm 1 can ensure a trivial regret bound, with the guarantee
585 that $\text{Regret}(T) \leq 2T$.
- 586 • In the scenario where $C \leq \bar{C}$, we know that \bar{C} remains a valid upper bound on the number of
587 adversarial feedback. Under this situation, Algorithm 1 operates successfully with \bar{C} adversarial
588 feedback. Therefore, according to Theorem 5.3, the regret is upper bounded by

$$\text{Regret}(T) \leq \tilde{O}(d\sqrt{T} + d\bar{C}).$$

589

□

590 **C Proof of Lemmas 5.1 and B.1**

591 **C.1 Proof of Lemma 5.1**

592 *Proof of Lemma 5.1.* Using a similar reasoning in Li et al. (2017), we define some auxiliary quantities

$$\begin{aligned} G_t(\boldsymbol{\theta}) &= \lambda\kappa\boldsymbol{\theta} + \sum_{i=1}^{t-1} w_i \left[\sigma \left((\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i))^\top \boldsymbol{\theta} \right) \right. \\ &\quad \left. - \sigma \left((\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i))^\top \boldsymbol{\theta}^* \right) \right] (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i)), \\ \epsilon_t &= l_t - \sigma \left((\boldsymbol{\phi}(x_t, a_t) - \boldsymbol{\phi}(x_t, b_t))^\top \boldsymbol{\theta}^* \right), \\ \gamma_t &= o_t - \sigma \left((\boldsymbol{\phi}(x_t, a_t) - \boldsymbol{\phi}(x_t, b_t))^\top \boldsymbol{\theta}^* \right), \\ Z_t &= \sum_{i=1}^{t-1} w_i \gamma_i (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i)). \end{aligned}$$

593 In Algorithm 1, $\boldsymbol{\theta}_t$ is chosen to be the solution to the following equation,

$$\lambda\kappa\boldsymbol{\theta}_t + \sum_{i=1}^{t-1} w_i \left[\sigma \left((\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i))^\top \boldsymbol{\theta}_t \right) - o_i \right] (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i)) = \mathbf{0}. \quad (\text{C.1})$$

594 Then we have

$$\begin{aligned} G_t(\boldsymbol{\theta}_t) &= \lambda\kappa\boldsymbol{\theta}_t + \sum_{i=1}^{t-1} w_i \left[\sigma \left((\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i))^\top \boldsymbol{\theta}_t \right) \right. \\ &\quad \left. - \sigma \left((\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i))^\top \boldsymbol{\theta}^* \right) \right] (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i)) \\ &= \sum_{i=1}^{t-1} w_i \left[o_i - \sigma \left((\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i))^\top \boldsymbol{\theta}^* \right) \right] (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i)) \\ &= Z_t. \end{aligned}$$

595 The analysis in Li et al. (2017); Di et al. (2023) shows that this equation has a unique solution, with
596 $\boldsymbol{\theta}_t = G_t^{-1}(Z_t)$. Using the mean value theorem, for any $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}^d$, there exists $m \in [0, 1]$ and
597 $\bar{\boldsymbol{\theta}} = m\boldsymbol{\theta}_1 + (1 - m)\boldsymbol{\theta}_2$, such that the following equation holds,

$$\begin{aligned} G_t(\boldsymbol{\theta}_1) - G_t(\boldsymbol{\theta}_2) &= \lambda\kappa(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2) + \sum_{i=1}^{t-1} w_i \left[\sigma \left((\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i))^\top \boldsymbol{\theta}_1 \right) \right. \\ &\quad \left. - \sigma \left((\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i))^\top \boldsymbol{\theta}_2 \right) \right] (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i)) \\ &= \left[\lambda\kappa\mathbf{I} + \sum_{i=1}^{t-1} w_i \dot{\sigma} \left((\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i))^\top \bar{\boldsymbol{\theta}} \right) \right. \\ &\quad \left. (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i)) (\boldsymbol{\phi}(x_i, a_i) - \boldsymbol{\phi}(x_i, b_i))^\top \right] (\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2). \end{aligned}$$

598 We define $F(\bar{\theta})$ as

$$F(\bar{\theta}) = \lambda\kappa\mathbf{I} + \sum_{i=1}^{t-1} w_i \dot{\sigma} \left((\phi(x_i, a_i) - \phi(x_i, b_i))^\top \bar{\theta} \right) (\phi(x_i, a_i) - \phi(x_i, b_i)) (\phi(x_i, a_i) - \phi(x_i, b_i))^\top \Big].$$

599 Moreover, we can see that $G_t(\theta^*) = \lambda\kappa\theta^*$. Recall $\Sigma_t = \lambda\mathbf{I} + \sum_{i=1}^{t-1} w_i (\phi(x_i, a_i) -$
600 $\phi(x_i, b_i)) (\phi(x_i, a_i) - \phi(x_i, b_i))^\top$. We have

$$\begin{aligned} \|G_t(\theta_t) - G_t(\theta^*)\|_{\Sigma_t^{-1}}^2 &= (\theta_t - \theta^*)^\top F(\bar{\theta}) \Sigma_t^{-1} F(\bar{\theta}) (\theta_t - \theta^*) \\ &\geq \kappa^2 (\theta_t - \theta^*)^\top \Sigma_t (\theta_t - \theta^*) \\ &= \kappa^2 \|\theta_t - \theta^*\|_{\Sigma_t}^2, \end{aligned}$$

601 where the first inequality holds due to $\dot{\mu}(\cdot) \geq \kappa > 0$ and $F(\bar{\theta}) \succeq \kappa \Sigma_t$. Then we have the following
602 estimate of the estimation error:

$$\begin{aligned} \|\theta_t - \theta^*\|_{\Sigma_t} &\leq \frac{1}{\kappa} \|G_t(\theta_t) - G_t(\theta^*)\|_{\Sigma_t^{-1}} \\ &\leq \lambda \|\theta^*\|_{\Sigma_t^{-1}} + \frac{1}{\kappa} \|Z_t\|_{\Sigma_t^{-1}} \\ &\leq \sqrt{\lambda} \|\theta^*\|_2 + \frac{1}{\kappa} \|Z_t\|_{\Sigma_t^{-1}}, \end{aligned}$$

603 where the second inequality holds due to the triangle inequality and $G_t(\theta^*) = \lambda\kappa\theta^*$. The last
604 inequality holds due to $\Sigma_t \succeq \lambda\mathbf{I}$. Finally, we need to bound the $\|Z_t\|_{\Sigma_t^{-1}}$ term. To study the impact
605 of adversarial feedback, we decompose the summation in (A.1) based on the adversarial feedback c_t ,
606 i.e.,

$$Z_t = \sum_{i<t:c_i=0} w_i \gamma_i (\phi(x_i, a_i) - \phi(x_i, b_i)) + \sum_{i<t:c_i=1} w_i \gamma_i (\phi(x_i, a_i) - \phi(x_i, b_i)),$$

607 When $c_i = 1$, i.e. with adversarial feedback, $|\gamma_i - \epsilon_i| = 1$. On the contrary, when $c_i = 0$, $\gamma_i = \epsilon_i$.
608 Therefore,

$$\begin{aligned} \sum_{i<t:c_i=0} w_i \gamma_i (\phi(x_i, a_i) - \phi(x_i, b_i)) &= \sum_{i<t:c_i=0} w_i \epsilon_i (\phi(x_i, a_i) - \phi(x_i, b_i)), \\ \sum_{i<t:c_i=1} w_i \gamma_i (\phi(x_i, a_i) - \phi(x_i, b_i)) &= \sum_{i<t:c_i=1} w_i \epsilon_i (\phi(x_i, a_i) - \phi(x_i, b_i)) \\ &\quad + \sum_{i<t:c_i=1} w_i (\gamma_i - \epsilon_i) (\phi(x_i, a_i) - \phi(x_i, b_i)). \end{aligned}$$

609 Summing up the two equalities, we have

$$Z_t = \sum_{i=1}^{t-1} w_i \epsilon_i (\phi(x_i, a_i) - \phi(x_i, b_i)) + \sum_{i<t:c_i=1} w_i (\gamma_i - \epsilon_i) (\phi(x_i, a_i) - \phi(x_i, b_i)).$$

610 Therefore,

$$\|Z_t\|_{\Sigma_t^{-1}} \leq \underbrace{\left\| \sum_{i=1}^{t-1} w_i \epsilon_i (\phi(x_i, a_i) - \phi(x_i, b_i)) \right\|_{\Sigma_t^{-1}}}_{I_1} + \underbrace{\left\| \sum_{i<t:c_i=1} w_i (\phi(x_i, a_i) - \phi(x_i, b_i)) \right\|_{\Sigma_t^{-1}}}_{I_2}.$$

611 For the term I_1 , with probability at least $1 - \delta$, for all $t \in [T]$, it can be bounded by

$$I_1 \leq \sqrt{2 \log \left(\frac{\det(\Sigma_t)^{1/2} \det(\Sigma_0)^{-1/2}}{\delta} \right)},$$

612 due to Lemma D.2. Using $w_i \leq 1$, we have $\sqrt{w_i} \|\phi(x_i, a_i) - \phi(x_i, b_i)\|_2 \leq 2$. Moreover, we have

$$\det(\Sigma_t) \leq \left(\frac{\text{Tr}(\Sigma_t)}{d} \right)^d$$

$$\begin{aligned}
&= \left(\frac{d\lambda + \sum_{i=1}^{t-1} w_i \|(\phi(x_i, a_i) - \phi(x_i, b_i))\|_2^2}{d} \right)^d \\
&\leq \left(\frac{d\lambda + 2T}{d} \right)^d,
\end{aligned}$$

613 where the first inequality holds because for every matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$, $\det \mathbf{A} \leq (\text{Tr}(\mathbf{A})/d)^d$. The
614 second inequality holds due to $\sqrt{w_i} \|\phi(x_i, a_i) - \phi(x_i, b_i)\|_2 \leq 2$. Easy to see that $\det(\Sigma_0) = \lambda^d$.
615 The term I_1 can be bounded by

$$I_1 \leq \sqrt{d \log((1 + 2T/\lambda)/\delta)}. \quad (\text{C.2})$$

616 For I_2 , with our choice of the weight w_i , we have

$$\begin{aligned}
I_2 &\leq \sum_{i < t: c_i=1} w_i \|(\phi(x_i, a_i) - \phi(x_i, b_i))\|_{\Sigma_t^{-1}} \\
&\leq \sum_{i < t: c_i=1} w_i \|(\phi(x_i, a_i) - \phi(x_i, b_i))\|_{\Sigma_i^{-1}} \\
&\leq \sum_{i < t: c_i=1} \alpha \\
&\leq \alpha C,
\end{aligned} \quad (\text{C.3})$$

617 where the second inequality holds due to $\Sigma_t \succeq \Sigma_i$. The third inequality holds due to $w_i \leq$
618 $\alpha / \|(\phi(x_i, a_i) - \phi(x_i, b_i))\|_{\Sigma_i^{-1}}$. The last inequality holds due to the definition of C . Combining
619 (C.2) and (C.3), we complete the proof of Lemma 5.1. \square

620 C.2 Proof of Lemma B.1

621 *Proof of Lemma B.1.* Let the regret incurred in the t -th round by $r_t = 2r^*(x_t, a_t^*) - r^*(x_t, a_t) -$
622 $r^*(x_t, b_t)$. It can be decomposed as

$$\begin{aligned}
r_t &= 2r^*(x_t, a_t^*) - r^*(x_t, a_t) - r^*(x_t, b_t) \\
&= \langle \phi(x_t, a_t^*) - \phi(x_t, a_t), \theta^* \rangle + \langle \phi(x_t, a_t^*) - \phi(x_t, b_t), \theta^* \rangle \\
&= \langle \phi(x_t, a_t^*) - \phi(x_t, a_t), \theta^* - \theta_t \rangle + \langle \phi(x_t, a_t^*) - \phi(x_t, b_t), \theta^* - \theta_t \rangle \\
&\quad + \langle 2\phi(x_t, a_t^*) - \phi(x_t, a_t) - \phi(x_t, b_t), \theta_t \rangle \\
&\leq \|\phi(x_t, a_t^*) - \phi(x_t, a_t)\|_{\Sigma_t^{-1}} \|\theta^* - \theta_t\|_{\Sigma_t} + \|\phi(x_t, a_t^*) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \|\theta^* - \theta_t\|_{\Sigma_t} \\
&\quad + \langle 2\phi(x_t, a_t^*) - \phi(x_t, a_t) - \phi(x_t, b_t), \theta_t \rangle \\
&\leq \beta \|\phi(x_t, a_t^*) - \phi(x_t, a_t)\|_{\Sigma_t^{-1}} + \beta \|\phi(x_t, a_t^*) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \\
&\quad + \langle 2\phi(x_t, a_t^*) - \phi(x_t, a_t) - \phi(x_t, b_t), \theta_t \rangle,
\end{aligned}$$

623 where the first inequality holds due to the Cauchy-Schwarz inequality. The second inequality holds
624 due to the high probability confidence event \mathcal{E} . Using our action selection rule, we have

$$\begin{aligned}
&\langle \phi(x_t, a_t^*) - \phi(x_t, a_t), \theta_t \rangle + \beta \|\phi(x_t, a_t^*) - \phi(x_t, a_t)\|_{\Sigma_t^{-1}} \\
&\leq \langle \phi(x_t, b_t) - \phi(x_t, a_t), \theta_t \rangle + \beta \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \\
&\langle \phi(x_t, a_t^*) - \phi(x_t, b_t), \theta_t \rangle + \beta \|\phi(x_t, a_t^*) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \\
&\leq \langle \phi(x_t, a_t) - \phi(x_t, b_t), \theta_t \rangle + \beta \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}}.
\end{aligned}$$

625 Adding the above two inequalities, we have

$$\begin{aligned}
&\beta \|\phi(x_t, a_t^*) - \phi(x_t, a_t)\|_{\Sigma_t^{-1}} + \beta \|\phi(x_t, a_t^*) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \\
&\leq \langle \phi(x_t, a_t) + \phi(x_t, b_t) - 2\phi(x_t, a_t^*), \theta_t \rangle + 2\beta \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}}.
\end{aligned}$$

626 Therefore, we prove that the regret in round t can be upper bounded by

$$r_t \leq 2\beta \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}}.$$

627 With a simple observation, we have $r_t \leq 4$. Therefore, the total regret can be upper bounded by

$$\text{Regret}(T) \leq \sum_{t=1}^T \min \left\{ 4, 2\beta \|\phi(x_t, a_t) - \phi(x_t, b_t)\|_{\Sigma_t^{-1}} \right\}.$$

628 \square

629 **D Auxiliary Lemmas**

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

$$\sum_{t=1}^T \eta_t \leq R\sqrt{2T \log 1/\delta}.$$

634 **Lemma D.2** (Lemma 9 Abbasi-Yadkori et al. 2011). Let $\{\epsilon_t\}_{t=1}^T$ be a real-valued stochastic process
 635 with corresponding filtration $\{\mathcal{F}_t\}_{t=0}^T$ such that ϵ_t is \mathcal{F}_t -measurable and ϵ_t is conditionally R -sub-
 636 Gaussian, i.e.

$$\forall \lambda \in \mathbb{R}, \mathbb{E}[e^{\lambda \epsilon_t} | \mathcal{F}_{t-1}] \leq \exp\left(\frac{\lambda^2 R^2}{2}\right).$$

637 Let $\{\mathbf{x}_t\}_{t=1}^T$ be an \mathbb{R}^d -valued stochastic process where \mathbf{x}_t is \mathcal{F}_{t-1} -measurable and for any $t \in [T]$,
 638 we further define $\Sigma_t = \lambda \mathbf{I} + \sum_{i=1}^t \mathbf{x}_i \mathbf{x}_i^\top$. Then with probability at least $1 - \delta$, for all $t \in [T]$, we
 639 have

$$\left\| \sum_{i=1}^T \mathbf{x}_i \eta_i \right\|_{\Sigma_t^{-1}}^2 \leq 2R^2 \log \left(\frac{\det(\Sigma_t)^{1/2} \det(\Sigma_0)^{-1/2}}{\delta} \right).$$

640 **Lemma D.3** (Lemma 11, Abbasi-Yadkori et al. 2011). For any $\lambda > 0$ and sequence $\{\mathbf{x}_t\}_{t=1}^T \subseteq \mathbb{R}^d$
 641 for $t \in [T]$, define $\mathbf{Z}_t = \lambda \mathbf{I} + \sum_{i=1}^{t-1} \mathbf{x}_i \mathbf{x}_i^\top$. Then, provided that $\|\mathbf{x}_t\|_2 \leq L$ holds for all $t \in [T]$, we
 642 have

$$\sum_{t=1}^T \min \left\{ 1, \|\mathbf{x}_t\|_{\mathbf{Z}_t^{-1}}^2 \right\} \leq 2d \log(1 + TL^2/(d\lambda)).$$

643 **NeurIPS Paper Checklist**

644 **1. Claims**

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

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

648 Justification: The primary contribution of this paper is addressing the challenge of adversarial
649 feedback within the dueling bandit model, where feedback is represented as a binary
650 preference label. Our research introduces a new perspective to machine learning. Unlike
651 previous works on corruption-robust bandits, where corruption in each round affects the
652 single-arm exploration and exploitation process. Flipping the preference label potentially
653 impacts the expected reward of both actions chosen in a duel. This interaction can further
654 affect subsequent decisions involving only one of these arms. Compared with previous
655 adversarial dueling bandit work, we study the most direct label-flipping attack, which is
656 aligned with many real-life preference-based learning scenarios. Our uncertainty-weighted
657 maximum likelihood estimation method helps to solve this novel problem, in scenarios with
658 known and unknown adversarial feedback. All the scope has been discussed clearly in our
659 abstract and introduction.

660 Guidelines:

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

670 **2. Limitations**

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

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

673 Justification: We have added a Limitations setting in our main paper. We assume that
674 the reward is linear with respect to some known feature maps. Although this setting is
675 common in the literature, we observe that some recent works on dueling bandits can deal
676 with nonlinear rewards (Li et al., 2024). Therefore, it's possible to extend our results to a
677 more general setting. Another assumption concerns the lower bound of the derivative of
678 the link function. Notably, in the logistic bandit model, which shares similarities with our
679 setting through Bernoulli variables, some work (Abeille et al., 2021; Fauray et al., 2022) can
680 improve the dependency of κ from $1/\kappa$ to $\sqrt{\kappa}$. A similar improvement might be achieved in
681 our setting as well.

682 Guidelines:

- 683 • The answer NA means that the paper has no limitation while the answer No means that
684 the paper has limitations, but those are not discussed in the paper.
- 685 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 686 • The paper should point out any strong assumptions and how robust the results are to
687 violations of these assumptions (e.g., independence assumptions, noiseless settings,
688 model well-specification, asymptotic approximations only holding locally). The authors
689 should reflect on how these assumptions might be violated in practice and what the
690 implications would be.
- 691 • The authors should reflect on the scope of the claims made, e.g., if the approach was
692 only tested on a few datasets or with a few runs. In general, empirical results often
693 depend on implicit assumptions, which should be articulated.
- 694 • The authors should reflect on the factors that influence the performance of the approach.
695 For example, a facial recognition algorithm may perform poorly when image resolution
696 is low or images are taken in low lighting. Or a speech-to-text system might not be

697 used reliably to provide closed captions for online lectures because it fails to handle
698 technical jargon.

- 699 • The authors should discuss the computational efficiency of the proposed algorithms
700 and how they scale with dataset size.
- 701 • If applicable, the authors should discuss possible limitations of their approach to
702 address problems of privacy and fairness.
- 703 • While the authors might fear that complete honesty about limitations might be used by
704 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
705 limitations that aren't acknowledged in the paper. The authors should use their best
706 judgment and recognize that individual actions in favor of transparency play an impor-
707 tant role in developing norms that preserve the integrity of the community. Reviewers
708 will be specifically instructed to not penalize honesty concerning limitations.

709 3. Theory Assumptions and Proofs

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

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

713 Justification: We have clearly stated and proved all the lemmas and theorems used in our
714 theoretical results. To help readers understand the proof without checking all the details, we
715 provide a roadmap of our proof in Appendix A. We also write explanation and clarification
716 for every formula in our paper.

717 Guidelines:

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

728 4. Experimental Result Reproducibility

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

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

733 Justification: Our paper is mainly theoretical but we also do numerical experiments to justify
734 the correctness of our results. We provide all the information to reproduce our results in
735 Section 6.

736 Guidelines:

- 737 • The answer NA means that the paper does not include experiments.
- 738 • If the paper includes experiments, a No answer to this question will not be perceived
739 well by the reviewers: Making the paper reproducible is important, regardless of
740 whether the code and data are provided or not.
- 741 • If the contribution is a dataset and/or model, the authors should describe the steps taken
742 to make their results reproducible or verifiable.
- 743 • Depending on the contribution, reproducibility can be accomplished in various ways.
744 For example, if the contribution is a novel architecture, describing the architecture fully
745 might suffice, or if the contribution is a specific model and empirical evaluation, it may
746 be necessary to either make it possible for others to replicate the model with the same
747 dataset, or provide access to the model. In general, releasing code and data is often
748 one good way to accomplish this, but reproducibility can also be provided via detailed
749 instructions for how to replicate the results, access to a hosted model (e.g., in the case

- 750 of a large language model), releasing of a model checkpoint, or other means that are
751 appropriate to the research performed.
- 752 • While NeurIPS does not require releasing code, the conference does require all submis-
753 sions to provide some reasonable avenue for reproducibility, which may depend on the
754 nature of the contribution. For example
 - 755 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
756 to reproduce that algorithm.
 - 757 (b) If the contribution is primarily a new model architecture, the paper should describe
758 the architecture clearly and fully.
 - 759 (c) If the contribution is a new model (e.g., a large language model), then there should
760 either be a way to access this model for reproducing the results or a way to reproduce
761 the model (e.g., with an open-source dataset or instructions for how to construct
762 the dataset).
 - 763 (d) We recognize that reproducibility may be tricky in some cases, in which case
764 authors are welcome to describe the particular way they provide for reproducibility.
765 In the case of closed-source models, it may be that access to the model is limited in
766 some way (e.g., to registered users), but it should be possible for other researchers
767 to have some path to reproducing or verifying the results.

768 5. Open access to data and code

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

772 Answer: [No]

773 Justification: Our experiments involve synthetic data generated from a generalized linear
774 model, which is quite simple and easy to reproduce. That’s why we do not provide access
775 to our data and code. All the information required to reproduce the results is provided in
776 Section 6. Guidelines:

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

796 6. Experimental Setting/Details

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

800 Answer: [Yes]

801 Justification:

802 Guidelines:

- 803 • The answer NA means that the paper does not include experiments.

- 804 • The experimental setting should be presented in the core of the paper to a level of detail
805 that is necessary to appreciate the results and make sense of them.
806 • The full details can be provided either with the code, in appendix, or as supplemental
807 material.

808 7. Experiment Statistical Significance

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

811 Answer: [No]

812 Justification:

813 Guidelines:

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

834 8. Experiments Compute Resources

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

838 Answer: [Yes]

839 Justification: We only have synthetic experiments and it can be reproduced on CPUs.

840 Guidelines:

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

849 9. Code Of Ethics

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

852 Answer: [Yes]

853 Justification:

854 Guidelines:

- 855
- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 856
- If the authors answer No, they should explain the special circumstances that require a
- 857
- deviation from the Code of Ethics.
- 858
- The authors should make sure to preserve anonymity (e.g., if there is a special consid-
- 859
- eration due to laws or regulations in their jurisdiction).

860 10. Broader Impacts

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

863 Answer: [Yes]

864 Justification: This paper studies contextual dueling bandits with adversarial feedback. Our
865 primary objective is to propel advancements in bandit theory by introducing a more robust
866 algorithm backed by solid theoretical guarantees. The uncertainty-weighted approach
867 we have developed for dueling bandits holds significant potential to address the issue of
868 adversarial feedback in preference-based data, which could be instrumental in enhancing the
869 robustness of generative models against adversarial attacks, thereby contributing positively
870 to the societal impact and reliability of machine learning applications.

871 Guidelines:

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

894 11. Safeguards

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

898 Answer: [NA]

899 Justification:

900 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.
- 901
- 902
- 903
- 904
- 905
- 906
- 907

908 • We recognize that providing effective safeguards is challenging, and many papers do
909 not require this, but we encourage authors to take this into account and make a best
910 faith effort.

911 12. Licenses for existing assets

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

915 Answer: [NA]

916 Justification:

917 Guidelines:

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

933 13. New Assets

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

936 Answer: [NA]

937 Justification:

938 Guidelines:

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

947 14. Crowdsourcing and Research with Human Subjects

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

951 Answer: [NA]

952 Justification:

953 Guidelines:

- 954 • The answer NA means that the paper does not involve crowdsourcing nor research with
955 human subjects.
- 956 • Including this information in the supplemental material is fine, but if the main contribu-
957 tion of the paper involves human subjects, then as much detail as possible should be
958 included in the main paper.

959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980

- 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:

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.