
On Private and Robust Bandits

Yulian Wu*
KAUST
yulian.wu@kaust.edu.sa

Xingyu Zhou*
Wayne State University
xingyu.zhou@wayne.edu

Youming Tao
Shandong University
ym.tao99@mail.sdu.edu.cn

Di Wang
KAUST
di.wang@kaust.edu.sa

Abstract

We study private and robust multi-armed bandits (MABs), where the agent receives Huber’s contaminated heavy-tailed rewards and meanwhile needs to ensure differential privacy. We consider both the finite k -th raw moment and the finite k -th central moment settings for heavy-tailed rewards distributions with $k \geq 2$. We first present its minimax lower bound, characterizing the information-theoretic limit of regret with respect to privacy budget, contamination level, and heavy-tailedness. Then, we propose a meta-algorithm that builds on a private and robust mean estimation sub-routine PRM that essentially relies on reward truncation and the Laplace mechanism. For the above two different heavy-tailed settings, we give corresponding schemes of PRM, which enable us to achieve nearly-optimal regrets. Moreover, our two proposed truncation-based or histogram-based PRM schemes achieve the optimal trade-off between estimation accuracy, privacy and robustness. Finally, we support our theoretical results and show the effectiveness of our algorithms with experimental studies.

1 Introduction

The multi-armed bandit (MAB) [1] problem provides a fundamental framework for sequential decision-making under uncertainty with bandit feedback, which has drawn a wide range of applications in medicine [2], finance [3, 4], recommendation system [5], and online advertising [6], to name a few. Consider a portfolio selection in finance as an example. At each decision round $t \in [T]$, the learning agent selects an action $a_t \in [K]$ (i.e., a choice of assets to user t) and receives a reward r_t (e.g., the corresponding payoff) that is i.i.d. drawn from an unknown probability distribution associated with the portfolio choice. The goal is to learn to maximize its cumulative payoff.

In practice, applying the celebrated MAB formulation to real-life applications (e.g., the above finance example) needs to deal with both robustness and privacy issues. On the one hand, it is known that financial data is often heavy-tailed (rather than sub-Gaussian) [7, 8]. Moreover, the received payoff data in finance often contains outliers [9] due to data contamination. On the other hand, privacy concern in finance is growing [10–12]. For instance, even if the adversary does not have direct access to the dataset, they are still able to reconstruct other customers’ personal information by interacting with the pricing platform and observing its decisions [13].

Motivated by this, a line of work on MABs has focused on designing robust algorithms with respect to heavy-tailed rewards [14], adversary contamination [15, 16], or both [17]. Another line of recent work has studied privacy protection in MABs via different trust models of differential privacy (DP) [18] such as central DP [19, 20], local DP [21, 22] and distributed DP [23, 24]. Moreover, there have also

*Equal contribution.

been recent advances in understanding the close relationship between robustness and privacy for the mean estimation problem (e.g., robustness induces privacy [25, 26] and vice versa [27]). In light of this, a fundamental question we are interested in this paper is:

Is there a simple algorithm that can tackle privacy and robustness in MABs simultaneously?

Our contributions. We give an affirmative answer to it by showing that a unified truncation-based algorithm could achieve a nearly optimal trade-off between regret, privacy, and robustness for MABs under two different heavy-tailed settings. The key intuition is that reward truncation not only helps to reduce outliers (due to both heavy tails and contamination), but also bound its sensitivity (useful for DP guarantees). To make our intuition rigorous, we take the following principled approaches.

(i) We first establish the minimax regret lower bound for private and robust MABs, i.e., heavy-tailed MABs with both privacy constraints and Huber’s contamination (see section 4) where the reward distribution of each action has finite k -th moment for $k \geq 2$. This characterizes the information-theoretic limit of regret with respect to privacy budget, contamination level and heavy-tailedness.

(ii) To match the lower bound, we first propose a meta-algorithm (see section 5), which builds upon the idea of batched successive elimination and relies on a generic private and robust mean estimation sub-routine denoted by PRM. Then, for two different settings of (heavy-tailed) reward distributions (i.e., finite raw or central moments), we propose corresponding schemes for the sub-routine PRM, both of which require truncation and the Laplace mechanism to guarantee robustness and privacy, simultaneously. Armed with these, our meta-algorithm can enjoy nearly matching regret upper bounds (see section 6). Experimental studies also corroborate our theoretical results.

(iii) Along the way, several results could be of independent interest. In particular, our proposed PRM shows that reward truncation is sufficient to help achieve the optimal high-probability concentration for private and robust mean estimation in the one-dimension case. Moreover, without contamination, our regret upper bounds not only match the optimal one for private heavy-tailed MABs with finite raw moments, but also provide the first results for the case with finite central moments, hence a complete study for private bandits as well.

Due to space limit, **experiments (section A in Appendix), technical lemmas, and all proofs** are relegated to Appendix.

2 Related Work

Robust MABs. The studies on robust bandits can be largely categorized into two groups. The first group of work mainly focuses on the setting where the total contamination is bounded, i.e., the cumulative difference between observed reward and true reward is bounded by some constant [15]. The second group considers Huber’s α -contamination model [28] (which is also the focus of our paper) or a similar α -fraction model. In these cases, the reward for each round can be contaminated by an arbitrary distribution with probability $\alpha \in [0, 1/2]$ [29, 16, 30], or at most α -fraction of the rewards are arbitrarily contaminated [31]. The existing work in this group has mainly focused on the light-tailed setting where the true inlier distribution is Gaussian or sub-Gaussian and uses a robust median or trimmed-mean estimator. A very recent work [17] studies the (non-private) setting where the inlier distribution only has finite variance and uses Huber’s estimator to establish problem-dependent bounds. In contrast, we take the perspective of minimax regret, i.e., problem-independent bounds, and also account for privacy protection.

Private MABs. In addition to the above mentioned results on private MABs with light-tailed rewards, [22] studies private heavy-tailed MABs with finite raw moments under both central and local models of DP. However, how to design private algorithms for heavy-tailed distributions with finite *central* moments is left as an important problem. In this paper, we resolve the problem as a byproduct of our main results.

Robust and private mean estimation. Our work is also related to robust and private mean estimation, especially the one-dimensional case. On the robustness side with Huber’s model, a high-probability concentration bound for the median of Gaussian (hence the mean by symmetry) is first established in [32]. Recently, [30] gives a high probability mean concentration via a trimmed-mean estimator for general sub-Gaussian inlier distributions while [33] focuses on the heavy-tailed setting. We refer the reader to the survey paper [34] on high-dimensional robust statistics focusing on robust mean estimation. On the robustness side with heavy-tail, we refer the reader to the survey paper

[35] for mean estimation and regression under heavy-tailed distributions. Besides median-of-means techniques and trimmed mean we mentioned above to handle heavy-tailed data, Catoni’s estimator is a very different estimator for heavy tail [36, 37] and also used in bandits [14]. On the privacy side, one close work is [38], which presents the first high-probability mean concentration for private heavy-tailed distributions with finite central moments (via a medians-of-means approach). It is worth noting that there are recent exciting advances in understanding the close relationship between robustness and privacy (e.g., robustness induces privacy [25, 26] and vice versa [27]). From this aspect, our results imply that for the one-dimensional mean estimation problem, truncation alone suffices to help to achieve both.

Concurrent and independent work. While preparing this submission, we have noticed that [39] also studies private and robust bandits problems under the Huber model. However, there are many differences between our work and theirs. First, for the differential privacy and bandit models, we investigate the multi-armed bandit problem under central differential privacy while they study stochastic linear bandits under the local differential privacy model. For the assumption of rewards, we focus on the heavy-tailed cases where the distribution of reward has finite k -th raw moment and central moment for $k \geq 2$ while they essentially assume that the rewards are bounded².

3 Preliminaries

In this section, we first formally introduce our problem of private and robust MABs and subsequently present its corresponding regret notions.

3.1 Private and Robust MABs

In this paper, we consider the multi-armed bandits problem where the agent interacts with the environment for T rounds. In each round t , the agent chooses an action $a_t \in [K]$ and then standard reward r_t is generated independently from inlier distribution. After contamination, the agent observes contaminated reward x_t .

As mentioned before, by robustness, we aim to handle both reward contamination and possible heavy-tailed inlier distributions. To this end, we first introduce the following two classes of heavy-tailed reward distributions.

Definition 3.1 (Finite k -th raw moment). *A distribution over \mathbb{R} is said to have a finite k -th raw moment if it is within*

$$\mathcal{P}_k = \{P : \mathbb{E}_{X \sim P} [|X|^k] \leq 1\}, \quad k \geq 2, \quad (1)$$

where k is considered fixed but arbitrary.

Definition 3.2 (Finite k -th central moment [40, 38]). *A distribution over \mathbb{R} is said to have a finite k -th central moment if it is within*

$$\mathcal{P}_k^c = \{P : \mathbb{E}_{X \sim P} [|X - \mu|^k] \leq 1\}, \quad k \geq 2, \quad (2)$$

where $\mu := \mathbb{E}_{X \sim P}[X] \in [-D, D]$ and $D \geq 1$.

The relationship between the finite raw moment and the finite central moment has been discussed in [22]. We further consider the celebrated Huber contamination model [28] in heavy-tailed MABs.

Definition 3.3 (Heavy-tailed MABs with Huber contamination). *Given the corruption level $\alpha \in [0, 1/2)$. For each round $t \in [T]$, the observed reward³ x_t for action a_t , is sampled independently from the true distribution $P_{a_t} \in \mathcal{P}_k$ (or $P_{a_t} \in \mathcal{P}_k^c$) with probability $1 - \alpha$; otherwise is sampled from some arbitrary and unknown contamination distribution $G_{a_t} \in \mathcal{G}$.*

In addition to robustness, we also consider the privacy protection in MABs via the lens of DP. In particular, we consider the standard central model of DP for MABs (e.g., [41]), where the learning agent has access to users’ raw data (i.e., rewards) and guarantees that its output (i.e., sequence

²Note that although the general problem formulation in [39] considers sub-Gaussian rewards, their privacy guarantee only holds for bounded rewards.

³Here we use x_t in the contaminated case to distinguish with standard reward r_t .

of actions) are indistinguishable in probability on two neighboring reward sequences. Due to contamination, the reward data accessed by the learning agent at round t could have already been contaminated. More precisely, we let $\mathcal{D}_T = (x_1, \dots, x_T) \in \mathbb{R}^T$ be a reward sequence generated in the learning process and $\mathcal{M}(\mathcal{D}_T) = (a_1, \dots, a_T) \in [K]^T$ to denote the sequence of all actions recommended by a learning algorithm \mathcal{M} . With this setup, we have the following formal definition.

Definition 3.4 (Differential Privacy for MABs). *For any $\epsilon > 0$, a learning algorithm $\mathcal{M} : \mathbb{R}^T \rightarrow [K]^T$ is ϵ -DP if for all sequences $\mathcal{D}_T, \mathcal{D}'_T \in \mathbb{R}^T$ differing only in a single element and for all events $E \subset [K]^T$, we have*

$$\mathbb{P}[\mathcal{M}(\mathcal{D}_T) \in E] \leq e^\epsilon \cdot \mathbb{P}[\mathcal{M}(\mathcal{D}'_T) \in E].$$

In other words, we protect the privacy of any individual user who interacts with the learning agent in the sense that an adversary observing the output of the learning agent (i.e., a sequence of actions) cannot infer too much about whether any particular individual has participated in this process, or the specific reward feedback of this individual. In this paper, we will leverage the well-known Laplace mechanism to guarantee differential privacy.

Definition 3.5 (Laplace Mechanism). *Given a function $f : \mathcal{X}^n \rightarrow \mathbb{R}^d$, the Laplacian mechanism is given by*

$$\mathcal{M}_L(\mathcal{D}_n, f, \epsilon) = f(\mathcal{D}_n) + (Y_1, Y_2, \dots, Y_d),$$

where Y_i is i.i.d. drawn from a Laplacian Distribution⁴ $\text{Lap}(\frac{\Delta_1(f)}{\epsilon})$, where $\Delta_1(f)$ is the ℓ_1 -sensitivity of the function f , i.e., $\Delta_1(f) = \sup_{\mathcal{D}_n, \mathcal{D}'_n} \|f(\mathcal{D}_n) - f(\mathcal{D}'_n)\|_1$. Then, for any $\epsilon > 0$, Laplacian mechanism satisfies ϵ -DP.

In the following sections, for brevity, we will simply use *private and robust MABs* to refer to our setting, i.e., heavy-tailed MABs with Huber contamination and privacy constraints.

3.2 Regrets for Private and Robust MABs

In the contamination case, the standard regret using observed (contaminated) rewards $\{x_t\}_{t \in [T]}$ is ill-defined [31]. Instead, the literature focuses on the *clean regret*, that is, to compete with the best policy in hindsight as measured by the expected true uncontaminated rewards [31, 17, 42]. Hence, let μ_a be the mean of the inlier distribution of arm $a \in [K]$ and $\mu^* = \max_{a \in [K]} \mu_a$. We also let Π^ϵ be the set of all ϵ -DP MAB algorithms and $\mathcal{E}_{\alpha, k}$ be the set of all instances of heavy-tailed MABs (with parameter k) with Huber contamination (of level α).

Definition 3.6 (Clean Regret). *Fix an algorithm $\pi \in \Pi^\epsilon$ and an instance $\nu \in \mathcal{E}_{\alpha, k}$. Then, the clean regret of π under ν is given by $\mathcal{R}_T(\pi, \nu) := \mathbb{E}_{\pi, \nu}[T\mu^* - \sum_{t=1}^T \mu_{a_t}]$.*

Note that here the expectation is taken over the randomness generated by the *contaminated* environment and ϵ -DP MAB algorithm while the means are of the true inlier distributions.

To capture the intrinsic difficulty of the private and robust MAB problem, we are also interested in its minimax regret.

Definition 3.7 (Minimax Regret). *The minimax regret of our private and robust MAB problem is defined as $\mathcal{R}_{\epsilon, \alpha, k}^{\text{minimax}} := \inf_{\pi \in \Pi^\epsilon} \sup_{\nu \in \mathcal{E}_{\alpha, k}} \mathbb{E}_{\pi, \nu}[T\mu^* - \sum_{t=1}^T \mu_{a_t}]$.*

4 Lower Bound

We start with the following lower bound on the minimax regret, which characterizes the fundamental impact of privacy budget (via ϵ), contamination level (via α) and heavy-tailedness of rewards (via k) in the regret. Note that this lower bound is mainly established under Definition 3.1, which in turn also serves as one valid lower bound for the central moment case in Definition 3.2.

Theorem 4.1. *Consider a private and robust MAB problem where inlier distributions have finite k -th raw (or central) moments ($k \geq 2$). Then, its minimax regret satisfies*

$$\mathcal{R}_{\epsilon, \alpha, k}^{\text{minimax}} = \Omega\left(\sqrt{KT} + (K/\epsilon)^{1-\frac{1}{k}} T^{\frac{1}{k}} + T\alpha^{1-\frac{1}{k}}\right).$$

⁴For a parameter λ , the Laplacian distribution has the density function $\text{Lap}(\lambda)(x) = \frac{1}{2\lambda} \exp(-\frac{|x|}{\lambda})$.

Algorithm 1 Private and Robust Arm Elimination

```
1: Input: Number of arms  $K$ , time horizon  $T$ , privacy budget  $\epsilon$ , Huber parameter  $\alpha \in (0, 1/2)$ ,  
   error probability  $\delta \in (0, 1]$ , inlier distribution parameters i.e.,  $k$  and optional  $D$ .  
2: Initialize:  $\tau = 0$ , active set of arms  $\mathcal{S} = \{1, \dots, K\}$ .  
3: for batch  $\tau = 1, 2, \dots$  do  
4:   Set batch size  $B_\tau = 2^\tau$ .  
5:   if  $B_\tau < \mathcal{T}$  then  
6:     Randomly select an action  $a \in [K]$ .  
7:     Play action  $a$  for  $B_\tau$  times.  
8:   else  
9:     for each active arm  $a \in \mathcal{S}$  do  
10:    for  $i$  from 1 to  $B_\tau$  do  
11:      Pull arm  $a$ , observe contaminated reward  $x_i^a$ .  
12:      If total number of pulls reaches  $T$ , exit.  
13:    end for  
14:    Set truncation threshold  $M_\tau$ .  
15:    Set additional parameters  $\Phi$ .  
16:    Compute estimate  $\tilde{\mu}_a = \text{PRM}(\{x_i^a\}_{i=1}^{B_\tau}, M_\tau, \Phi)$ .  
17:    end for  
18:    Set confidence radius  $\beta_\tau$ .  
19:    Let  $\tilde{\mu}_{\max} = \max_{a \in \mathcal{S}} \tilde{\mu}_a$ .  
20:    Remove all arms  $a$  from  $\mathcal{S}$  s.t.  $\tilde{\mu}_{\max} - \tilde{\mu}_a > 2\beta_\tau$ .  
21:  end if  
22: end for
```

This lower bound basically takes a maximum of three terms. The first term comes from the standard regret for Gaussian rewards, the second one captures the additional cost in regret due to privacy and heavy-tailed rewards, and the last term indicates the additional cost in regret due to contamination and heavy-tailed rewards. Note that, for a given k , the impact of privacy and contamination is separable. It would also be helpful to compare our lower bound with the related ones, as discussed below.

Remark 4.2. *First, when $k = \infty$ and $\alpha = 0$, our lower bound recovers the state-of-the-art lower bound for private MABs with sub-Gaussian rewards [20]; Second, when there is no privacy protection, a very recent work [17] establishes a problem-dependent regret lower bound for robust MABs while we are interested in the problem-independent lower bound with further privacy protection.*

5 Our Approach: A Meta-Algorithm

In this section, we first introduce a meta-algorithm for private and robust MABs, which not only allows us to tackle inlier distributions with bounded raw or central moments in a unified way, but also highlights the key component, i.e., a private and robust mean estimation sub-routine building on the main idea of reward truncation.

Our meta-algorithm, at a high level, can be viewed as a batched version of the celebrated successive arm elimination [43] along with a private and robust mean estimation sub-routine PRM (see Algorithm 1). That is, it divides the time horizon T into batches with exponentially increasing size and eliminates sub-optimal arms successively based on the mean estimate via PRM. More specifically, based on the batch size, it consists of two phases. That is, when the batch size is less than a threshold \mathcal{T} , it simply recommends actions randomly (line 5-7) (more on this will be explained soon). Otherwise, for each active arm a in batch τ , it first prescribes a to a batch of $B_\tau = 2^\tau$ fresh *new* users (i.e., “doubling”) and observes possibly contaminated rewards (line 8). Then, it calls the sub-routine PRM to compute a private and robust mean estimate for each active arm a (line 12). In particular, it *only* uses the rewards within the most recent batch (i.e., “forgetting”) along with a proper reward truncation threshold M_τ . Finally, it adopts the classic idea of arm elimination with a proper choice of confidence radius β_τ to remove sub-optimal arms with high confidence (line 18-20).

We now provide more intuitions behind our algorithm design by highlighting how its main components work in concert. *First*, the reason behind the first phase (i.e., $B_\tau \leq \mathcal{T}$), named forced exploration, is necessary as it ensures that concentration is satisfied so that the optimal arm will not be eliminated.

Algorithm 2 PRM for the finite raw moment case

- 1: **Input:** A collection of data $\{x_i\}_{i=1}^n$, truncation parameter M , additional parameters $\Phi = \{\epsilon\}$.
 - 2: **for** $i = 1, 2, \dots, n$ **do**
 - 3: Truncate data $\bar{x}_i = x_i \cdot \mathbb{1}_{\{|x_i| \leq M\}}$.
 - 4: **end for**
 - 5: Return private estimate $\tilde{\mu} = \frac{\sum_{i=1}^n \bar{x}_i}{n} + \text{Lap}(\frac{2M}{n\epsilon})$.
-

This is due to the fact that in the contamination case, a well-behaved concentration only kicks in when the number of samples is larger than a threshold. Thus, one cannot adopt arm elimination in this phase since it might eliminate the optimal arm. Note that, instead of our choice of random selection, one can also use other methods for the first phase (see Remark 5.1 below). *Second*, for the second phase, the idea of doubling batching and forgetting is the key to achieving privacy with a minimal amount of noise (hence better regret). This is because now any single reward feedback only impacts one computation of estimate. This is in sharp contrast to standard arm elimination (e.g., [43]) where each mean estimate is based on all samples so far (as no batching is used), and hence a single reward change could impact $O(T)$ mean estimations⁵. *Third*, the simple idea of reward truncation in PRM turns out to be extremely useful for both robustness and privacy. On the one hand, truncation helps to reduce the impact of outliers (due to both heavy tails and contamination); On the other hand, truncation also helps to bound the sensitivity, which is necessary for privacy. In fact, as we will show later, a well-tuned truncation threshold enables us to achieve a near-optimal trade-off between regret, privacy and robustness. *Finally*, in contrast to the first phase, we can now eliminate sub-optimal arms with high confidence due to the high probability concentration of mean estimate when batch size is larger than \mathcal{T} (more details will be given later for specific choices of PRM and hence the choice of \mathcal{T}).

We choose to use the successive elimination (SE) technique here because it suffices to enable us to achieve optimal order-wise regret later. In fact, once armed with our novel PRM module, one can also use other exploration strategies like UCB in [20]. One difference is that now instead of first pulling each arm once, it needs to pull each arm \mathcal{T} times to ensure that concentration kicks in later. This is in fact not surprising since on the high-level, the analysis of SE and UCB is very similar, i.e., doubling (batching) and forgetting. We also note that instead of UCB, one can also adapt it to the Thompson sampling strategy, e.g., [45], with our PRM module. Again, the key idea is batching and forgetting.

Remark 5.1. *The algorithmic choice of the first phase can be flexible. For example, instead of playing a randomly selected action for the whole batch, one can choose to play a randomly selected action for each round. Moreover, one can also choose to be greedy or probabilistically greedy with respect to the mean estimate by PRM, which also only uses the rewards collected within the last batch for each arm. All of these choices have the same theoretical guarantees, though some will help to improve the empirical performance.*

We then present the following remark that places our meta-algorithm in the existing literature.

Remark 5.2 (Comparison with existing literature). *For private MABs (without contamination), the state-of-the-art also builds upon the idea of batching and forgetting [19, 20, 24] to achieve optimal regret. For robust MABs (without privacy), existing works take different robust mean estimations. For example, both [31, 30] use a trimmed mean estimator for sub-Gaussian inlier distributions while [17] adopts Huber’s estimator to handle inlier distributions with only bounded variance. We are the first to study privacy and robustness simultaneously, via a simple truncation-based estimator, which in turn reveals the close relationship between privacy and robustness in MABs. This complements the recent advances in capturing the connection between these two in (high-dimensional) statistics [25, 27].*

6 Upper Bounds

In this section, we establish the regret upper bounds for two specific instantiations of our meta-algorithm – one for the finite raw moment case and another for the finite central moment case. The results could match our lower bound up to a log factor on T , demonstrating their near-optimality.

⁵One can use the tree-based mechanism [44] to reduce it to $O(\log T)$, but it is still sub-optimal [19].

6.1 Finite Raw Moment Case

In this section, we will focus on private and robust MABs where the inlier distributions have a finite k -th raw moment as given by Definition 3.1. In particular, we first introduce the choice of PRM in this case (see Algorithm 2) and establish its concentration property, which plays a key role in our implementation of meta-algorithm.

The PRM in Algorithm 2 is simply a truncation-based Laplace mechanism. That is, it first truncates all the received data with the threshold M (line 3). Then, Laplace noise is added to the empirical mean to preserve privacy (line 5). We highlight again that truncation here helps with both robustness (via removing outliers) and privacy (via bounding the sensitivity of empirical mean).

As in the standard algorithm design of MABs, the key is to utilize the concentration of the mean estimator. To this end, we first give the following high-probability concentration result for the mean estimate returned by PRM in Algorithm 2.

Theorem 6.1 (Concentration of Mean Estimate). *Given a collection of Huber-contaminated data $\{x_i\}_{i=1}^n$ where the inlier distribution satisfies Definition 3.1 with mean μ , let $\tilde{\mu}$ be the mean estimate by Algorithm 2. Then, for any privacy budget $\epsilon > 0$ and $\delta \in (0, 1)$, the following results hold:*

Uncontaminated case. For $\alpha = 0$, we have $|\tilde{\mu} - \mu| = O\left(\sqrt{\frac{\log(1/\delta)}{n}} + \frac{M \log(1/\delta)}{n\epsilon} + \frac{1}{M^{k-1}}\right)$, with probability at least $1 - \delta$. Thus, choosing the truncation threshold $M = \Theta\left(\left(\frac{n\epsilon}{\log(1/\delta)}\right)^{\frac{1}{k}}\right)$ yields $|\tilde{\mu} - \mu| = O\left(\sqrt{\frac{\log(1/\delta)}{n}} + \left(\frac{\log(1/\delta)}{n\epsilon}\right)^{1-\frac{1}{k}}\right)$.

Contaminated case. For $0 < \alpha \leq \alpha_1 \in (0, 1/2)$ and $n = \Omega\left(\frac{\log(1/\delta)}{\alpha_1}\right)$, we have the following with probability at least $1 - \delta$, $|\tilde{\mu} - \mu| = O\left(\sqrt{\frac{\log(1/\delta)}{n}} + \frac{M \log(1/\delta)}{n\epsilon} + \frac{1}{M^{k-1}} + \alpha_1 M\right)$. Therefore, choosing the truncation threshold $M = \Theta\left(\min\left\{\left(\frac{n\epsilon}{\log(1/\delta)}\right)^{\frac{1}{k}}, \alpha_1^{-\frac{1}{k}}\right\}\right)$, yields $|\tilde{\mu} - \mu| \leq \beta$, where $\beta = O\left(\sqrt{\frac{\log(1/\delta)}{n}} + \left(\frac{\log(1/\delta)}{n\epsilon}\right)^{1-\frac{1}{k}} + \alpha_1^{1-\frac{1}{k}}\right)$.

With the above result, several remarks are in order. *First*, for the uncontaminated case, our concentration result consists of the standard sub-Gaussian term and a new one due to privacy and heavy-tailed data. It can be translated into a sample complexity bound, i.e., to guarantee $|\tilde{\mu} - \mu| \leq \eta$ for any $\eta \in (0, 1)$, it requires the sample size to be $n \geq O\left(\frac{\log(1/\delta)}{\eta^2} + \frac{\log(1/\delta)}{\epsilon\eta^{k/(k-1)}}\right)$, which is optimal since it matches the lower bound for private heavy-tail mean estimation (cf. Theorem 7.2 in [46]). *Second*, for the contaminated case, it has an additional bias term $O(\alpha_1^{1-1/k})$, which is also known to be information-theoretically optimal [47]. Thus, via truncation, the PRM given by Algorithm 2 achieves the *optimal* trade-off between accuracy, privacy and robustness, which in turn shows its potential to be integrated into our meta-algorithm. Note that α_1 here is any upper bound (e.g., estimate) on the true contamination level α .

Now, based on the concentration result, we can set other missing parameters in our meta-algorithm accordingly. In particular, we have the following theorem that states the specific instantiation along with its performance guarantees.

Theorem 6.2 (Performance Guarantees). *Consider a private and robust MAB with inlier distributions satisfying Definition 3.1 and $0 < \alpha \leq \alpha_1 \in (0, 1/2)$. Let Algorithm 1 be instantiated with Algorithm 2 and M_τ, β_τ be given by Theorem 6.1 with n replaced by B_τ . Set $\mathcal{T} = \Omega\left(\frac{\log(1/\delta)}{\alpha_1}\right)$ and $\delta = 1/T$. Then Algorithm 1 is ϵ -DP with its regret upper bound*

$$\mathcal{R}_T = O\left(\sqrt{KT \log T} + \left(\frac{K \log T}{\epsilon}\right)^{\frac{k-1}{k}} T^{\frac{1}{k}} + T\alpha_1^{1-\frac{1}{k}} + \frac{K \log T}{\alpha_1}\right).$$

The above theorem presents the first achievable regret guarantee for private and robust bandits. The first three terms match our lower bound in Theorem 4.1 up to $\log T$ factor. The last additive term is

Algorithm 3 PRM for the finite central moment case

- 1: **Input:** A collection of data $\{x_i\}_{i=1}^{2n}$, truncation parameter M , additional parameters $\Phi = \{\epsilon, D, r\}$, $r \in \mathbb{R}$.
 - 2: // **First step: initial estimate**
 - 3: $B_j = [j, j + r)$, $j \in \mathcal{J} = \{-D, -D + r, \dots, D - r\}$.
 - 4: Compute private histogram using the first fold of data: $\tilde{p}_j = \frac{\sum_{i=1}^n \mathbb{1}_{\{X_i \in B_j\}}}{n} + \text{Lap}\left(\frac{2}{n\epsilon}\right)$.
 - 5: Get the initial estimate $J = \arg \max_{j \in \mathcal{J}} \tilde{p}_j$.
 - 6: // **Second step: final estimate**
 - 7: Get final estimator using the second fold of data: $\tilde{\mu} = J + \frac{1}{n} \sum_{i=n+1}^{2n} (X_i - J) \mathbb{1}_{\{|X_i - J| \leq M\}} + \text{Lap}\left(\frac{2M}{n\epsilon}\right)$.
-

mainly due to the fact that the mean concentration result only holds when the sample size is larger than $\mathcal{T} = \Omega\left(\frac{\log(1/\delta)}{\alpha_1}\right)$. As a result, each sub-optimal has to be played at least $\Omega\left(\frac{\log(1/\delta)}{\alpha_1}\right)$ times. However, for a sufficiently large T and a constant α , the last term is dominated by other terms. The last term is also exactly the reason that we choose an upper bound α_1 on the actual contamination level α and state the upper bound results in terms of α_1 rather than α . That is, for a very small but non-zero α , one can choose a larger α_1 to balance the regret. This subtle issue is also mentioned in one nice related work [42], see the remark after Theorem 7.4 and Remark 5.4 in their work.

Remark 6.3. *For the uncontaminated case with $\alpha = 0$, using the uncontaminated concentration bound in Theorem 6.1 and the same analysis, we achieve a regret upper bound $O(\sqrt{KT \log T} + \left(\frac{K \log T}{\epsilon}\right)^{\frac{k-1}{k}} T^{\frac{1}{k}})$, which also matches the lower bound in Theorem 4.1 up to a factor of $O(\log T)$.*

6.2 Finite Central Moment Case

The setting in the last section for the finite raw moment case may not be entirely satisfactory as it essentially assumes that the mean of arms is bounded within a small range (hence the sub-optimal gaps). Thus, in this section, we turn to private and robust MABs where the inlier distributions have a finite k -th central moment as given by Definition 3.2 for a reward distribution with large mean but small variability around the mean. To this end, we first need a new PRM, since now simply truncating around zero as in Algorithm 2 will not work well (as discussed in detail in Remark 6.13).

Our new PRM is presented in Algorithm 3, which consists of two steps. The intuition is simple: the first step aims to have a rough estimate of the mean, which is necessary since now the mean could be far away from zero. Then, in the second step, it truncates around the initial estimate to return the final result. More specifically, in the first step, we first construct bins over the range $[-D, D]$, which is assumed to contain the true mean by Definition 3.2. Then, we compute the private histogram via the Laplace mechanism. The initial estimate is given by the left endpoint of the bin that has the largest empirical mass. Next, in the second step, it simply truncates around the initial estimate and again adds Laplace noise for privacy.

Remark 6.4. *It is worth noting that a similar idea of two-step estimation has been used in previous work on robust mean estimation in the one-dimensional heavy-tailed case [33, 38, 40]. However, there are several differences in our algorithm design and analysis. In particular, while [33] considers mean estimation under Huber's model without privacy constraints, we further impose differential privacy requirements. As a result, the estimates for both two steps are in different forms in our case compared to [33], though they share the same high-level intuition. On the other hand, while [38] considers mean estimation under differential privacy, there is no consideration of Huber contamination as in our case. Moreover, our second estimate is based on truncation while their method is via medians-of-means. In fact, as will be shown later (see Remark 6.7), when our result reduces to the uncontaminated case, it achieves improvement over the one in [38]. Finally, [40] considers both Huber contamination and local differential privacy, and establishes the corresponding mean square error (MSE). In contrast, we consider the central differential privacy and aim to establish a high-probability tail concentration. To this end, we take a different truncation method (i.e., using an indicator function in Line 7) compared to the one in [40] (Section 3.2 in its first arxiv version).*

As before, we first present the concentration property of our new PRM, which will manifest in the specific instantiation of our meta-algorithm. In particular, we first give the following general theorem and then state two more detailed corollaries.

Theorem 6.5 (Concentration of Mean Estimate). *Given a collection of Huber-contaminated data $\{x_i\}_{i=1}^{2n}$ where the inlier distribution satisfies Definition 3.2 with mean μ , let $\tilde{\mu}$ be the mean estimate by Algorithm 3. For any $\alpha \leq \alpha_1 \in (0, \alpha_{\max})$, $\epsilon \in (0, 1]$ and $\delta \in (0, 1)$, there exist some constants $\mathcal{T}(\alpha_1, \epsilon, \delta)$, r , M and $D \geq 2r$ such that for all $n \geq \mathcal{T}(\alpha_1, \epsilon, \delta)$, with probability at least $1 - \delta$*

$$|\tilde{\mu} - \mu| \leq O \left(\sqrt{\frac{\log(1/\delta)}{n}} + \frac{M \log(1/\delta)}{n\epsilon} + \frac{1}{M^{k-1}} + \alpha_1 M \right),$$

where $\alpha_{\max} < 1/2$ is the breakdown point.

The above theorem follows the same pattern as the one for the raw moment case (Theorem 6.1). The key differences are the threshold value $\mathcal{T}(\alpha, \epsilon, \delta)$ and the breakdown point α_{\max} , which are summarized in the following results.

Corollary 6.6 (Mean Concentration, $\alpha = 0$). *Let the same assumptions in Theorem 6.5 hold. For any $\epsilon \in (0, 1]$, setting $r = 10^{1/k}$ and $M = \Theta \left(\frac{n\epsilon}{\log(1/\delta)} \right)^{1/k}$, then for all $n \geq \Omega(\log(D/\delta)/\epsilon)$ and $D \geq 2r$, we have that for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, $|\tilde{\mu} - \mu| \leq \beta$ where $\beta = O \left(\sqrt{\frac{\log(1/\delta)}{n}} + \left(\frac{\log(1/\delta)}{n\epsilon} \right)^{1-\frac{1}{k}} \right)$. In other words, taking number of samples n such that $n \geq O \left(\frac{\log(1/\delta)}{\eta^2} + \frac{\log(1/\delta)}{\epsilon \eta^{\frac{k}{k-1}}} + \frac{\log(D/\delta)}{\epsilon} \right)$ we have $|\tilde{\mu} - \mu| \leq \eta$ with probability at least $1 - \delta$.*

Remark 6.7. *The above lemma strictly improves the result in [38, Theorem 3.5]⁶. In particular, it uses the method of medians-of-means and achieves $\frac{\log D \cdot \log(1/\delta)}{\epsilon}$ for the third term. In contrast, our third term is additive rather than multiplicative. In fact, our concentration is optimal, which matches the lower bound for the one-dimensional case (see [46, Theorem 7.2]).*

Corollary 6.8 (Mean Concentration, $\alpha > 0$). *Let the same assumptions in Theorem 6.5 hold. For any $\epsilon \in (0, 1]$ and $0 < \alpha \leq \alpha_1 \in (0, 0.133)$, we let $r = \iota^{1/k}$ where $\iota = \frac{1-\alpha}{0.249-\alpha}$ and $M = \Theta(\min\{ \left(\frac{n\epsilon}{\log(1/\delta)} \right)^{1/k}, (\alpha_1)^{-1/k} \})$. Then, there exists constant c_1 , for all n such that $n \geq \mathcal{T} = \Omega(\max\{ \frac{\iota \log(1/\delta)}{\epsilon}, \frac{c_1 \log(D/\delta)}{\epsilon}, \frac{\log(1/\delta)}{\alpha_1^2} \})$ and $D \geq 2r$, we have that for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, $|\tilde{\mu} - \mu| \leq \beta$ with $\beta = O \left(\sqrt{\frac{\log(1/\delta)}{n}} + \left(\frac{\log(1/\delta)}{n\epsilon} \right)^{1-\frac{1}{k}} + \alpha_1^{1-\frac{1}{k}} \right)$.*

Remark 6.9. *The above concentration has the same form as the one in Theorem 6.1. Specifically, for a large sample size n , it has the optimal concentration (for small α). The threshold \mathcal{T} on n depends on both α, ϵ now. We note that even for the sub-Gaussian inlier distributions without privacy protection, the existing concentration also has a threshold $\mathcal{T} = \frac{\log(1/\delta)}{\alpha_1^2}$ (see Lemma 4.1 in [30]). The finite central moment assumption will help us get logarithmic results for sample complexity on D of $n \geq \Omega(\log(D/\delta)/\epsilon)$ which is better than the polynomial results on D from the finite raw moment. Actually, finding a private estimator that achieves an error on the logarithmic of the range of the parameter is significant in the topic of DP estimation (see [48–50] on the importance of the problem).*

Now, we are left to leverage the above two concentration results to design specific instantiations of our meta-algorithm and establish their performance guarantees. Our first instantiation is for the uncontaminated case, i.e., $\alpha = 0$. Therefore, robustness is then only with respect to heavy-tailed rewards while privacy is still preserved.

Theorem 6.10 (Performance Guarantees, $\alpha = 0$). *Consider a private and robust MAB with inlier distributions satisfying Definition 3.2 and $\alpha = 0$. Let Algorithm 1 be instantiated with Algorithm 3, and r, M_τ, β_τ be given by Corollary 6.6 with n replaced by B_τ . Set $\mathcal{T} = \Omega(\frac{\log(D/\delta)}{\epsilon})$ and $\delta = 1/\mathcal{T}$. Then, Algorithm 1 is ϵ -DP with its regret upper bound*

$$\mathcal{R}_T = O \left(\sqrt{KT \log T} + (K \log T / \epsilon)^{\frac{k-1}{k}} T^{\frac{1}{k}} + \gamma \right),$$

where $\gamma := O(KD \log(DT)/\epsilon)$.

⁶We also note that the main focus of [38] is not on achieving the optimal estimate though.

Remark 6.11. *To the best of our knowledge, this is the first result on private and heavy-tailed bandits with the finite central moment assumption. The state-of-the-art result only focuses on the simpler case, i.e., the finite raw moment assumption [22].*

Finally, armed with Corollary 6.8, we have the second instantiation of our meta-algorithm that deals with the contaminated case.

Theorem 6.12 (Performance Guarantees, $\alpha > 0$). *Consider a private and robust MAB with inlier distributions satisfying Definition 3.2 and $\alpha \leq \alpha_1 \in (0, 0.133)$. Let Algorithm 1 be instantiated with Algorithm 3, and $r, \mathcal{T}, M_\tau, \beta_\tau$ be given by Corollary 6.8 with n replaced by B_τ . Set $\delta = 1/T$, then Algorithm 1 is ϵ -DP with its regret upper bound*

$$\mathcal{R}_T = O\left(\sqrt{KT \log T} + (K \log T / \epsilon)^{\frac{k-1}{k}} T^{\frac{1}{k}} + T \alpha_1^{1-\frac{1}{k}} + \hat{\gamma}\right),$$

where $\hat{\gamma} := O\left(\frac{DK \log T}{\alpha_1^2} + \frac{\iota DK \log T}{\epsilon} + \frac{DK \log(DT)}{\epsilon}\right)$ and $\iota = \frac{1-\alpha}{0.249-\alpha}$.

Remark 6.13. *The above upper bound also matches our lower bound up to log factor and $O(\hat{\gamma})$, which is dominated by other terms for a sufficiently large T and small α, D . A comparison between the results in Theorem 6.2 (Remark 6.3) and Theorem 6.12 (Theorem 6.10) for contaminated case (uncontaminated case), respectively, allows us to address a fundamental question: why we do not first transform the central moment condition to the raw moment condition, and then simply apply Algorithms 1 and 2 to get a regret upper bound? This is because the central moment condition implies $\mathbb{E}[|X|^k] = O(D^k)$. By employing a simple scaling technique, Algorithms 1 and 2 would yield a regret bound of $O(D \cdot \text{Reg})$, where Reg corresponds to the bound stated in Theorem 6.2. Thus, a significant contrast arises when comparing this result with Theorem 6.12, where the dependence on D is only an additive linear term rather than a multiplicative factor of $O(D)$. Consequently, the bounds presented in Theorem 6.12 and Theorem 6.10 prove to be substantially tighter when D assumes larger values, hence highlighting the contributions of our new PRM module in Algorithm 3.*

7 Simulations and Conclusion

We refer readers to Appendix A for our simulation results. In this paper, we investigated private and robust multi-armed bandits with heavy-tailed rewards under Huber's contamination model as well as differential privacy constraints. We proposed a meta-algorithm that builds on a private and robust mean estimation sub-routine PRM for different heavy-tailed assumptions of rewards, i.e. finite k -th raw moment and finite k -th central moment with $k \geq 2$. Moreover, we also established regret upper bounds for these algorithms, which nearly match our derived minimax lower bound.

There remain many open questions in this direction. First, the problem-dependent lower bound is unclear for private and robust bandits but some work has been done for corrupted heavy-tailed bandits (see Theorem 1 in [17]). One interesting future direction is to study how to leverage both the insights in [17] and the lower bound under privacy to derive a problem-dependent lower bound for private and robust bandits. Second, how to design algorithms and analyze the theoretical results for private and robust bandits under other corrupted models such as "bandits with total corruption budget" in [15]. Third, throughout the paper T is known in advance, and it remains to see how to handle the case where T is unknown and infinity.

8 Acknowledgments

We thank Mengchu Li for the insightful discussions and for pointing out the first arxiv version of [40]. And we also thank the anonymous NeurIPS reviewers for their feedback. Xingyu Zhou is supported in part by NSF grant CNS-2153220. Di Wang and Yulian Wu were supported in part by the baseline funding BAS/1/1689-01-01, funding from the CRG grand URF/1/4663-01-01, FCC/1/1976-49-01 from CBRC of King Abdullah University of Science and Technology (KAUST). Di Wang was also supported by the funding of the SDAIA-KAUST Center of Excellence in Data Science and Artificial Intelligence (SDAIA-KAUST AI)

References

- [1] Donald A Berry and Bert Fristedt. Bandit problems: sequential allocation of experiments (monographs on statistics and applied probability). *London: Chapman and Hall*, 5(71-87):7–7, 1985.
- [2] Benjamín Gutiérrez, Loïc Peter, Tassilo Klein, and Christian Wachinger. A multi-armed bandit to smartly select a training set from big medical data. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 38–45. Springer, 2017.
- [3] Matthew Hoffman, Eric Brochu, Nando De Freitas, et al. Portfolio allocation for bayesian optimization. In *UAI*, pages 327–336. Citeseer, 2011.
- [4] Weiwei Shen, Jun Wang, Yu-Gang Jiang, and Hongyuan Zha. Portfolio choices with orthogonal bandit learning. In *Twenty-fourth international joint conference on artificial intelligence*, 2015.
- [5] Stéphane Caron and Smriti Bhagat. Mixing bandits: A recipe for improved cold-start recommendations in a social network. In *Proceedings of the 7th Workshop on Social Network Mining and Analysis*, pages 1–9, 2013.
- [6] Eric M Schwartz, Eric T Bradlow, and Peter S Fader. Customer acquisition via display advertising using multi-armed bandit experiments. *Marketing Science*, 36(4):500–522, 2017.
- [7] Svetlozar Todorov Rachev. *Handbook of heavy tailed distributions in finance: Handbooks in finance, Book 1*. Elsevier, 2003.
- [8] John Hull. *Risk management and financial institutions, + Web Site*, volume 733. John Wiley & Sons, 2012.
- [9] John Adams, Darren Hayunga, Sattar Mansi, David Reeb, and Vincenzo Verardi. Identifying and treating outliers in finance. *Financial Management*, 48(2):345–384, 2019.
- [10] Yanzhe Murray Lei, Sentao Miao, and Ruslan Momot. Privacy-preserving personalized revenue management. *HEC Paris Research Paper No. MOSI-2020-1391*, 2020.
- [11] Xi Chen, David Simchi-Levi, and Yining Wang. Privacy-preserving dynamic personalized pricing with demand learning. *Management Science*, 68(7):4878–4898, 2022.
- [12] Xi Chen, Sentao Miao, and Yining Wang. Differential privacy in personalized pricing with nonparametric demand models. *Operations Research*, 2022.
- [13] Matthew Fredrikson, Eric Lantz, Somesh Jha, Simon Lin, David Page, and Thomas Ristenpart. Privacy in pharmacogenetics: An {End-to-End} case study of personalized warfarin dosing. In *23rd USENIX Security Symposium (USENIX Security 14)*, pages 17–32, 2014.
- [14] Sébastien Bubeck, Nicolo Cesa-Bianchi, and Gábor Lugosi. Bandits with heavy tail. *IEEE Transactions on Information Theory*, 59(11):7711–7717, 2013.
- [15] Thodoris Lykouris, Vahab Mirrokni, and Renato Paes Leme. Stochastic bandits robust to adversarial corruptions. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages 114–122, 2018.
- [16] Sayash Kapoor, Kumar Kshitij Patel, and Purushottam Kar. Corruption-tolerant bandit learning. *Machine Learning*, 108(4):687–715, 2019.
- [17] Debabrota Basu, Odalric-Ambrym Maillard, and Timothée Mathieu. Bandits corrupted by nature: Lower bounds on regret and robust optimistic algorithm. *arXiv preprint arXiv:2203.03186*, 2022.
- [18] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.
- [19] Touqir Sajed and Or Sheffet. An optimal private stochastic-mab algorithm based on optimal private stopping rule. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*, pages 5579–5588, 2019.

- [20] Achraf Azize and Debabrota Basu. When privacy meets partial information: A refined analysis of differentially private bandits. *arXiv preprint arXiv:2209.02570*, 2022.
- [21] Wenbo Ren, Xingyu Zhou, Jia Liu, and Ness B Shroff. Multi-armed bandits with local differential privacy. *arXiv preprint arXiv:2007.03121*, 2020.
- [22] Youming Tao, Yulian Wu, Peng Zhao, and Di Wang. Optimal rates of (locally) differentially private heavy-tailed multi-armed bandits. *arXiv preprint arXiv:2106.02575*, 2021.
- [23] Jay Tenenbaum, Haim Kaplan, Yishay Mansour, and Uri Stemmer. Differentially private multi-armed bandits in the shuffle model. *Advances in Neural Information Processing Systems*, 34, 2021.
- [24] Sayak Ray Chowdhury and Xingyu Zhou. Distributed differential privacy in multi-armed bandits. *arXiv preprint arXiv:2206.05772*, 2022.
- [25] Samuel B Hopkins, Gautam Kamath, Mahbod Majid, and Shyam Narayanan. Robustness implies privacy in statistical estimation. *arXiv preprint arXiv:2212.05015*, 2022.
- [26] Hilal Asi, Jonathan Ullman, and Lydia Zakynthinou. From robustness to privacy and back. *arXiv preprint arXiv:2302.01855*, 2023.
- [27] Kristian Georgiev and Samuel B Hopkins. Privacy induces robustness: Information-computation gaps and sparse mean estimation. *arXiv preprint arXiv:2211.00724*, 2022.
- [28] Peter J Huber. Robust estimation of a location parameter. *Ann. Math. Statist.*, 35(4):73–101, 1964.
- [29] Pranjal Awasthi, Sreenivas Gollapudi, Kostas Kollias, and Apaar Sadhwani. Online learning under adversarial corruptions. 2020.
- [30] Arpan Mukherjee, Ali Tajer, Pin-Yu Chen, and Payel Das. Mean-based best arm identification in stochastic bandits under reward contamination. *Advances in Neural Information Processing Systems*, 34:9651–9662, 2021.
- [31] Laura Niss and Ambuj Tewari. What you see may not be what you get: Ucb bandit algorithms robust to ε -contamination. In *Conference on Uncertainty in Artificial Intelligence*, pages 450–459. PMLR, 2020.
- [32] Kevin A Lai, Anup B Rao, and Santosh Vempala. Agnostic estimation of mean and covariance. In *2016 IEEE 57th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 665–674. IEEE, 2016.
- [33] Adarsh Prasad, Sivaraman Balakrishnan, and Pradeep Ravikumar. A unified approach to robust mean estimation. *arXiv preprint arXiv:1907.00927*, 2019.
- [34] Ilias Diakonikolas and Daniel M Kane. Recent advances in algorithmic high-dimensional robust statistics. *arXiv preprint arXiv:1911.05911*, 2019.
- [35] Gábor Lugosi and Shahar Mendelson. Mean estimation and regression under heavy-tailed distributions: A survey. *Foundations of Computational Mathematics*, 19(5):1145–1190, 2019.
- [36] Olivier Catoni and Ilaria Giullini. Dimension-free pac-bayesian bounds for the estimation of the mean of a random vector. *arXiv preprint arXiv:1802.04308*, 2018.
- [37] Sujay Bhatt, Guanhua Fang, Ping Li, and Gennady Samorodnitsky. Nearly optimal catoni’s m-estimator for infinite variance. In *International Conference on Machine Learning*, pages 1925–1944. PMLR, 2022.
- [38] Gautam Kamath, Vikrant Singhal, and Jonathan Ullman. Private mean estimation of heavy-tailed distributions. In *Conference on Learning Theory*, pages 2204–2235. PMLR, 2020.
- [39] Vasileios Charisopoulos, Hossein Esfandiari, and Vahab Mirrokni. Robust and private stochastic linear bandits. 2023.

- [40] Mengchu Li, Thomas B Berrett, and Yi Yu. On robustness and local differential privacy. *arXiv preprint arXiv:2201.00751*, 2022.
- [41] Nikita Mishra and Abhradeep Thakurta. (Nearly) optimal differentially private stochastic multi-arm bandits. In *Proceedings of the 31st Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 592–601, 2015.
- [42] Sitan Chen, Frederic Koehler, Ankur Moitra, and Morris Yau. Online and distribution-free robustness: Regression and contextual bandits with huber contamination. In *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)*, pages 684–695. IEEE, 2022.
- [43] Eyal Even-Dar, Shie Mannor, Yishay Mansour, and Sridhar Mahadevan. Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *Journal of Machine Learning Research*, 7(6), 2006.
- [44] T-H Hubert Chan, Elaine Shi, and Dawn Song. Private and continual release of statistics. *ACM Transactions on Information and System Security (TISSEC)*, 14(3):1–24, 2011.
- [45] Bingshan Hu and Nidhi Hegde. Near-optimal thompson sampling-based algorithms for differentially private stochastic bandits. In *Uncertainty in Artificial Intelligence*, pages 844–852. PMLR, 2022.
- [46] Samuel B Hopkins, Gautam Kamath, and Mahbod Majid. Efficient mean estimation with pure differential privacy via a sum-of-squares exponential mechanism. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1406–1417, 2022.
- [47] Ilias Diakonikolas. Algorithmic high-dimensional robust statistics. *Webpage <http://www.iliasdiakonikolas.org/simons-tutorial-robust.html>*, 2018.
- [48] Gautam Kamath, Jerry Li, Vikrant Singhal, and Jonathan Ullman. Privately learning high-dimensional distributions. In *Conference on Learning Theory*, pages 1853–1902. PMLR, 2019.
- [49] Vishesh Karwa and Salil Vadhan. Finite sample differentially private confidence intervals. *arXiv preprint arXiv:1711.03908*, 2017.
- [50] Gautam Kamath, Xingtu Liu, and Huanyu Zhang. Improved rates for differentially private stochastic convex optimization with heavy-tailed data. In *International Conference on Machine Learning*, pages 10633–10660. PMLR, 2022.
- [51] Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Found. Trends Theor. Comput. Sci.*, 9(3-4):211–407, 2014.
- [52] Joseph P Near and Chiké Abuah. Programming differential privacy. *URL: <https://uvm>*, 2021.
- [53] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- [54] Mengjie Chen, Chao Gao, and Zhao Ren. Robust covariance and scatter matrix estimation under huber’s contamination model. *The Annals of Statistics*, 46(5):1932–1960, 2018.
- [55] Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.

A Experiments

In this section, we empirically evaluate the practical performance of our private and robust arm elimination algorithms, denoted as PRAE-R and PRAE-C. These algorithms employ PRM as a sub-routine, which can be either Algorithm 2 for the finite raw moment case or Algorithm 3 for the finite central moment case. We benchmark our algorithms against DPRSE [22], which attains the optimal regret bound for DP heavy-tailed MAB, and RUCB [16], which is a non-private robust algorithm.

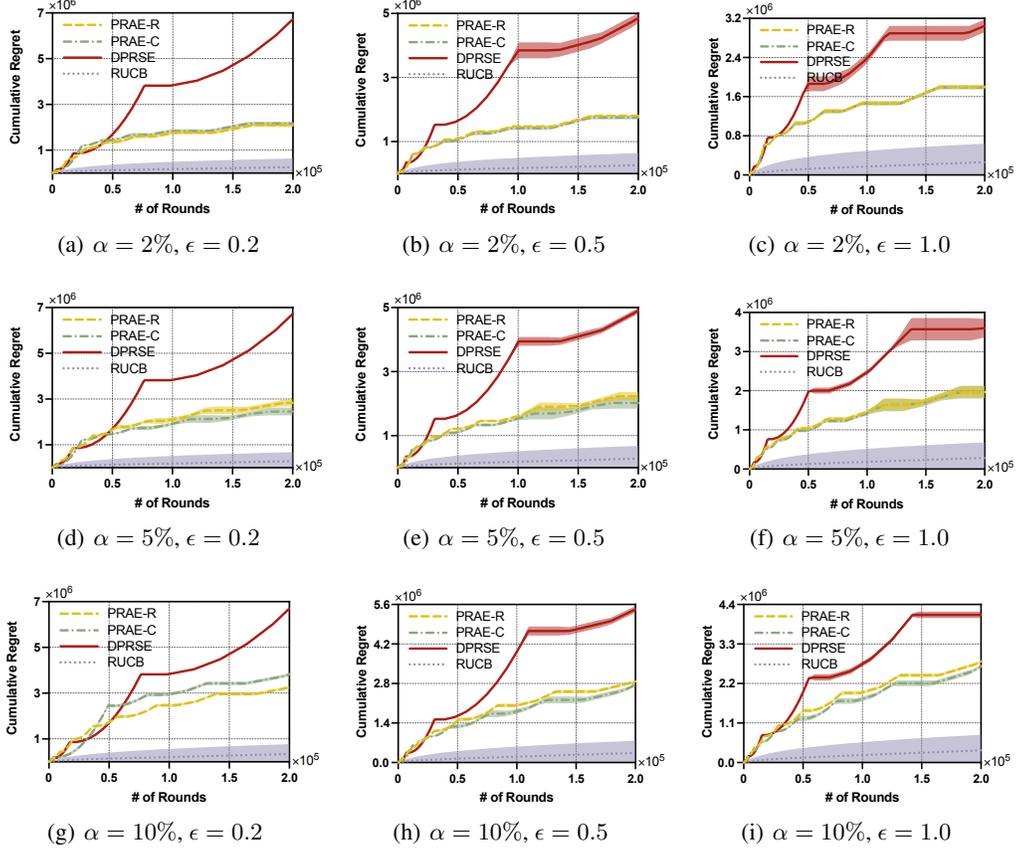


Figure 1: Comparison of cumulative regret for PRAE-R, PRAE-C and DPRSE under Pareto reward.

A.1 Experiment Setup

We consider the case where there are $K = 11$ arms, and the mean of each arm is within the range of $[0, 100]$. We set the means of the optimal arm and the worst arm to be 100 and 0, respectively, and the means of other arms scale linearly with a gap of 10. Specifically, for each arm $a \in [k]$, we have $\mu_a = 100 - \frac{100(a-1)}{K-1}$. We consider the following two types of heavy-tailed distributions for the true inlier reward generation:

- *Pareto distribution:* For each pull of arm a , we generate a reward that is sampled from the distribution $\mu_a + \eta - 2.5$, where $\eta \sim \frac{sx_m^s}{x^s+1} \mathbb{1}_{\{x \geq x_m\}}$ for $x \in \mathbb{R}$ and we set the shape parameter $s = 3$ and the scale parameter $x_m = 40$.
- *Student's t -distribution:* For each pull of arm a we generate a reward that is sampled from the distribution $\mu_a + \eta$, where $\eta \sim \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}}$. Here we set the degree of freedom $\nu = 2.0017$.

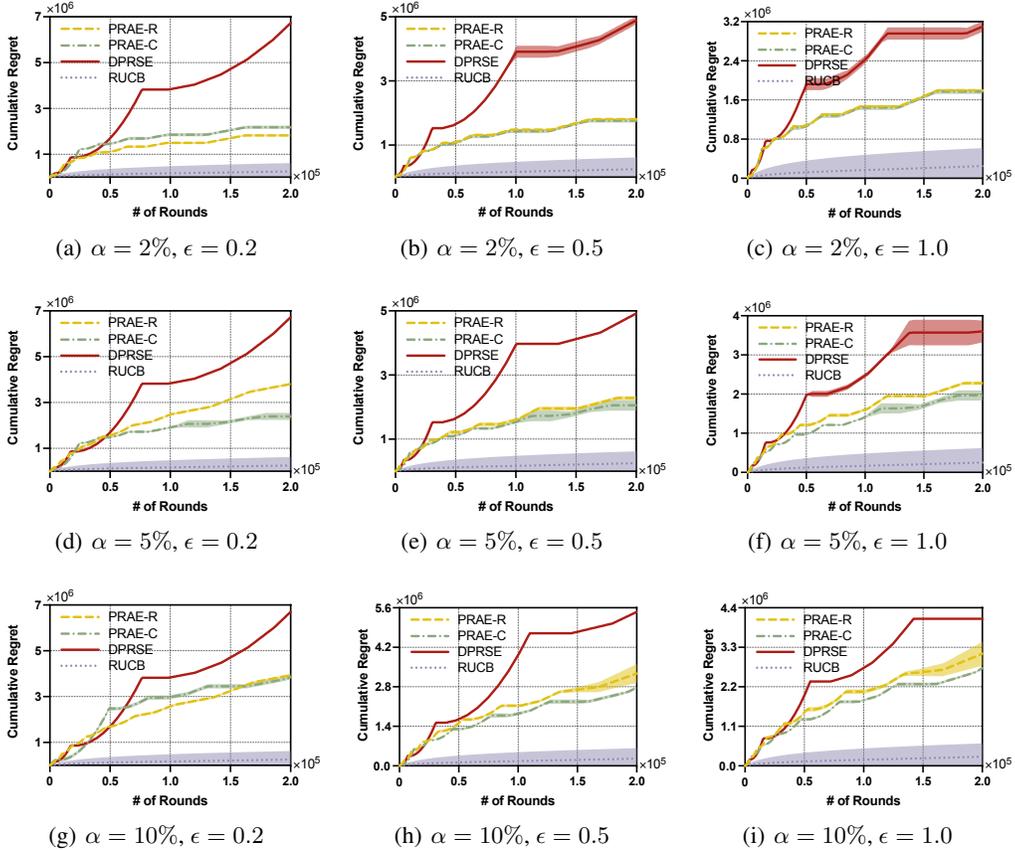


Figure 2: Comparison of cumulative regret for PRAE-R, PRAE-C and DPRSE under Student’s t reward.

For both cases, the stochastic rewards have finite second central moment of 35. The main difference between the above two types of distribution is that the Student’s t -distribution is symmetric while the Pareto distribution is one-sided.

To generate contaminated rewards, we still use Gaussian distribution. For the optimal arm, we set the mean of the Gaussian corruption distribution to be 0. For the other sub-optimal arms, we let the mean of their corrupted reward be 100. This way, the optimal arm is under-evaluated and the non-optimal arms are over-evaluated.

In our experiments, we will vary α and ϵ in $\{2\%, 5\%, 10\%\}$ and $\{0.2, 0.5, 1\}$, respectively. For each case, we repeat 30 times and set the total number of round $T = 10^5$ (thus we set $\delta = 10^{-5}$). We will report the average of cumulative regrets \mathcal{R}_T with respect to the number of rounds.

A.2 Results and Discussions

We present the experimental results under Pareto reward and Student’s t reward in Figure 1 and Figure 2, respectively. From these results, we can make the following observations:

- Firstly, by comparing the cumulative regret, we see that for all cases, PRAE-R and PRAE-C achieve smaller cumulative regret and thus better expected performance than DPRSE. In particular, when the Huber parameter α increases, DPRSE diverges to a larger regret, while PRAE-R and PRAE-C are only slightly affected. This is because DPRSE adopts more aggressive truncation thresholds that incorporate more outliers. In contrast, the truncation thresholds in PRAE-R and PRAE-C are carefully designed and thus provide robustness against contaminated rewards.

- Secondly, by looking at the error bars, we notice that PRAE-R and PRAE-C have smaller variance compared with both the private baseline and robust baseline under both symmetric and one-sided types of heavy-tailed distribution. In particular, although RUCB provides lower bounds on the performance with privacy for all the settings, PRAE-R and PRAE-C are more stable. The main reason could be that their robust mean estimator is based on median, which incurs smaller estimation error than the estimators developed in our paper.
- Thirdly, for both PRAE-R and PRAE-C, we also observe that when ϵ is smaller or α is larger, the regret will increase for both types of distributions, which is consistent with the fact that the regret bound is proportional to $1/\epsilon$ and α when T is large enough. Moreover, compared with PRAE-R, we can see that the regret of PRAE-C is lower in most cases. This is due to the fact that the PRM subroutine for PRAE-C leverages the prior information (i.e., range D) of mean for each arm, which could provide tighter performance bound. In summary, all these results corroborate our theoretical analysis.

B Useful Lemmas

Lemma B.1 (Post-Processing [51]). *Let $\mathcal{M} : \mathcal{X} \rightarrow \mathcal{Y}$ be a randomized algorithm that is (ϵ, δ) -differentially private. Let $f : \mathcal{Y} \rightarrow \mathcal{Z}$ be an arbitrary randomize mapping. Then $f \circ \mathcal{M} : \mathcal{X} \rightarrow \mathcal{Z}$ is (ϵ, δ) -differentially private.*

Lemma B.2 (Parallel Composition [52]). *Suppose there are n ϵ -differentially private mechanisms $\{\mathcal{M}_i\}_{i=1}^n$ and n disjoint datasets denoted by $\{D_i\}_{i=1}^n$. Then for the algorithm which applies each \mathcal{M}_i on the corresponding D_i , it is ϵ -DP.*

Lemma B.3 (Markov's inequality). *If $Y \in \mathbb{R}$ is a random variable and $a > 0$, we have*

$$\mathbb{P}(|Y| \geq a) \leq \frac{\mathbb{E}(|Y|^k)}{a^k}$$

Lemma B.4 (Chebyshev's inequality). *For a real-valued random variable $Y \in \mathbb{R}$, $a > 0$ and $k \in \mathbb{N}$, we have*

$$\mathbb{P}(|Y - \mathbb{E}Y| \geq a) = \mathbb{P}(|Y - \mathbb{E}Y|^k \geq a^k) \leq \frac{\mathbb{E}(|Y - \mathbb{E}Y|^k)}{a^k}$$

Lemma B.5 (Tail Bound of Laplacian Vairable [18]). *If $X \sim \text{Lap}(b)$, then*

$$\mathbb{P}(|X| \geq t \cdot b) = \exp(-t).$$

Lemma B.6 (Hoeffding's inequality). *Let Z_1, \dots, Z_n be independent bounded random variables with $Z_i \in [a, b]$ for all i , where $-\infty < a < b < \infty$. Then*

$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n (Z_i - \mathbb{E}[Z_i])\right| \geq t\right) \leq 2 \exp\left(-\frac{2nt^2}{(b-a)^2}\right)$$

Lemma B.7 (Hölder's Inequality). *Let X, Y be random variables over \mathbb{R} , and let $k > 1$. Then,*

$$\mathbb{E}[|XY|] \leq (\mathbb{E}[|X|^k])^{\frac{1}{k}} \left(\mathbb{E}[|Y|^{\frac{k}{k-1}}]\right)^{\frac{k-1}{k}}$$

Lemma B.8 (Bernstein's Inequality [53]). *Let X_1, \dots, X_n be n independent zero-mean random variables. Suppose $|X_i| \leq M$ and $\mathbb{E}[X_i^2] \leq s$ for all i . Then for any $t > 0$, we have*

$$\mathbb{P}\left\{\left|\frac{1}{n} \sum_{i=1}^n X_i\right| \geq t\right\} \leq 2 \exp\left(-\frac{\frac{1}{2}t^2n}{s + \frac{1}{3}Mt}\right)$$

C Proofs of Section 4

Lemma C.1 (Upper Bound on KL-divergence for Bandits with ϵ -DP [20]). *If π is a mechanism satisfying ϵ -DP, then for two instances $\nu_1 = (r_a : a \in [K])$ and $\nu_2 = (r'_a : a \in [K])$ we have*

$$KL(\mathbb{P}_{\pi, \nu_1}^T \| \mathbb{P}_{\pi, \nu_2}^T) \leq 6\epsilon \mathbb{E}_{\pi, \nu_1} \left[\sum_{t=1}^T TV(r_{a_t} \| r'_{a_t}) \right]$$

where $TV(r_a \| r'_a)$ is the total-variation distance between r_a and r'_a .

Lemma C.2 (Theorem 5.1 in [54]). *Let R_1 and R_2 be two distributions on \mathcal{X} . If for some $\alpha \in [0, 1/2)$, we have that $TV(R_1, R_2) = \frac{\alpha}{1-\alpha}$, then there exists two distributions on the same probability space G_1 and G_2 such that*

$$(1 - \alpha)R_1 + \alpha G_1 = (1 - \alpha)R_2 + \alpha G_2.$$

Proof of Theorem 4.1. Let Π be the set of all policies and Π^ϵ be the set of all ϵ -DP policies. We denote the environment corresponding to the set of K -Gaussian reward distributions with means $\mu \in \mathbb{R}^K$ the same variance σ_k^2 where the value of σ_k is determined by k to make the k -th raw moments of the distributions are bounded by 1 as $\mathcal{E}_N^K(\sigma_k) \triangleq \left\{ (\mathcal{N}(\mu_i, \sigma_k^2))_{i=1}^K : \mu = (\mu_1, \dots, \mu_K) \in \mathbb{R}^K \right\}$. Since $\Pi^\epsilon \subset \Pi$, we can have that

$$\mathcal{R}_T^{\text{minimax}}(\pi, \nu) \geq \inf_{\pi \in \Pi} \sup_{\nu \in \mathcal{E}_N^K(\sigma_k)} \text{Reg}_T(\pi, \nu) \geq \Omega(\sqrt{KT})$$

where the last inequality is due to Theorem 15.2 in [55].

Case 1: Uncontaminated case. By the definition of minimax regret, we know that $\mathcal{R}_{\epsilon, \alpha}^{\text{minimax}} \geq \mathcal{R}_{\epsilon, 0}^{\text{minimax}}$. Therefore, we first derive the lower bound of private bandits without contamination.

We consider two environments. In the first environment ν_1 , the optimal arm (denote by a_1) follows

$$r_{a_1} = \begin{cases} 1/\gamma & \text{with probability of } \frac{1}{2}\gamma^k \\ 0 & \text{with probability of } 1 - \frac{1}{2}\gamma^k \end{cases}$$

where $\gamma \in (0, 1]$. We can verify $\mathbb{E}[r_{a_1}] = \frac{1}{2}\gamma^{k-1}$ and $\mathbb{E}[r_{a_1}^k] = \frac{1}{2} \leq 1$. Any other sub-optimal arm $a \neq a_1$ in ν_1 follows the same reward distribution

$$r_a = \begin{cases} 1/\gamma & \text{with probability of } \frac{3}{10}\gamma^k \\ 0 & \text{with probability of } 1 - \frac{3}{10}\gamma^k \end{cases}$$

We can verify $\mathbb{E}[r_a] = \frac{3}{10}\gamma^{k-1}$ and $\mathbb{E}[r_{a_1}^k] = \frac{3}{10} \leq 1$. Then the gap of means between the optimal arm and sub-optimal arm is $\Delta = \frac{1}{5}\gamma^{k-1}$.

For algorithm π and instance ν_1 , we denote $i = \arg \min_{a \in \{2, \dots, K\}} \mathbb{E}_{\pi, \nu_1}[N_a(T)]$. Thus, $\mathbb{E}_{\pi, \nu_1}[N_i(T)] \leq \frac{T}{K-1}$.

Now, consider another instance ν_2 where r_{a_1}, \dots, r_{a_k} are the same as those in ν_1 except the i -th arm such that

$$r'_i = \begin{cases} 1/\gamma & \text{with probability of } \frac{7}{10}\gamma^k \\ 0 & \text{with probability of } 1 - \frac{7}{10}\gamma^k \end{cases}$$

We can verify $\mathbb{E}[r'_i] = \frac{7}{10}\gamma^{k-1}$ and $\mathbb{E}[(r'_i)^k] = \frac{7}{10} \leq 1$. Then in ν_2 , the arm i is optimal.

Now by the classic regret decomposition, we obtain

$$\mathcal{R}_T(\pi, \nu_1) = (T - \mathbb{E}_{\pi, \nu_1}[N_1(T)])\Delta \geq \mathbb{P}_{\pi, \nu_1}^T \left[N_1(T) \leq \frac{T}{2} \right] \frac{T\Delta}{2}.$$

$$\mathcal{R}_T(\pi, \nu_2) = \Delta \mathbb{E}_{\pi, \nu_2}[N_1(T)] + \sum_{a \notin \{1, i\}} 2\Delta \mathbb{E}_{\pi, \nu_2}[N_a(T)] \geq \mathbb{P}_{\pi, \nu_2}^T \left[N_1(T) \geq \frac{T}{2} \right] \frac{T\Delta}{2}.$$

By applying the Bretagnolle–Huber inequality ([55], Theorem 14.2), we have

$$\begin{aligned} \mathcal{R}_T(\pi, \nu_1) + \mathcal{R}_T(\pi, \nu_2) &\geq \frac{T\Delta}{2} \left(\mathbb{P}_{\pi, \nu_1}^T \left[N_1(T) \leq \frac{T}{2} \right] + \mathbb{P}_{\pi, \nu_2}^T \left[N_1(T) \geq \frac{T}{2} \right] \right) \\ &\geq \frac{T\Delta}{4} \exp(-\text{KL}(\mathbb{P}_{\pi, \nu_1}^T \parallel \mathbb{P}_{\pi, \nu_2}^T)) \end{aligned}$$

Based on Lemma C.1, we can get the upper bound of the KL-Divergence between the marginals.

$$\begin{aligned} \text{KL}(\mathbb{P}_{\pi, \nu_1}^T \parallel \mathbb{P}_{\pi, \nu_2}^T) &\leq 6\epsilon \mathbb{E}_{\pi, \nu_1} \left[\sum_{t=1}^T \text{TV}(r_{a_t} \parallel r'_{a_t}) \right] \\ &\leq 6\epsilon \mathbb{E}_{\pi, \nu_1}[N_i(T)] \text{TV}(r_i \parallel r'_i) \end{aligned}$$

since ν_1 and ν_i only differ in the arm i .

Thus,

$$\begin{aligned}\mathcal{R}_T(\pi, \nu_1) + \mathcal{R}_T(\pi, \nu_2) &\geq \frac{T\Delta}{4} \exp(-6\epsilon\mathbb{E}_{\pi, \nu_1}[N_i(T)] \cdot \frac{2}{5}\gamma^k) \\ &\geq \frac{T\gamma^{k-1}}{20} \exp\left(-\frac{12 \cdot \epsilon T\gamma^k}{5(K-1)}\right).\end{aligned}$$

Taking $\gamma = \left(\frac{K-1}{T\epsilon}\right)^{\frac{1}{k}}$, we get the result

$$\mathcal{R}_T(\pi, \nu_1) \geq \Omega\left(\left(\frac{K}{\epsilon}\right)^{\frac{k-1}{k}} T^{\frac{1}{k}}\right).$$

Case 2: Contaminated case. For $\alpha \neq 0$ and $\alpha \in (0, 1/2)$, we still consider the true distributions of arms are the same in above ν_1 and ν_2 . In the first environment ν_1 , the optimal arm (denote by a_1) follows

$$r_{a_1} = \begin{cases} 1/\gamma & \text{with probability of } \frac{1}{2}\gamma^k \\ 0 & \text{with probability of } 1 - \frac{1}{2}\gamma^k \end{cases}$$

where $\gamma \in (0, 1]$. We can verify $\mathbb{E}[r_{a_1}] = \frac{1}{2}\gamma^{k-1}$ and $\mathbb{E}[r_{a_1}^k] = \frac{1}{2} \leq 1$. Any other sub-optimal arm $a \neq a_1$ in ν_1 follows the same reward distribution

$$r_a = \begin{cases} 1/\gamma & \text{with probability of } \frac{3}{10}\gamma^k \\ 0 & \text{with probability of } 1 - \frac{3}{10}\gamma^k \end{cases}$$

We can verify $\mathbb{E}[r_a] = \frac{3}{10}\gamma^{k-1}$ and $\mathbb{E}[r_a^k] = \frac{3}{10} \leq 1$. Then the gap of means between the optimal arm and sub-optimal arm is $\Delta = \frac{1}{5}\gamma^{k-1}$.

And we denote the contaminated version of ν_1 as $\tilde{\nu}_1$. For algorithm π and instance $\tilde{\nu}_1$, we denote $i = \arg \min_{a \in \{2, \dots, K\}} \mathbb{E}_{\pi, \tilde{\nu}_1}[N_a(T)]$. Thus, $\mathbb{E}_{\pi, \tilde{\nu}_1}[N_i(T)] \leq \frac{T}{K-1}$.

Now, consider another instance ν_2 where r_{a_1}, \dots, r_{a_k} are the same as those in ν_1 except the i -th arm such that

$$r'_i = \begin{cases} 1/\gamma & \text{with probability of } \frac{7}{10}\gamma^k \\ 0 & \text{with probability of } 1 - \frac{7}{10}\gamma^k \end{cases}$$

We can verify $\mathbb{E}[r'_i] = \frac{7}{10}\gamma^{k-1}$ and $\mathbb{E}[(r'_i)^k] = \frac{7}{10} \leq 1$. Then in ν_2 , the arm i is optimal.

Also, we denote the contaminated version of ν_2 as $\tilde{\nu}_2$. Take $\gamma = \alpha^{\frac{1}{k}} \in (0, 1]$, since for any $a \in [K]$, $\text{TV}(r_a \| r'_a) \leq \frac{2}{5}\gamma^k = \frac{2}{5}\alpha \leq \frac{\alpha}{1-\alpha}$, from Lemma C.2, we have for any arm $a \in [K]$, there exists distribution G_a and G'_a such that

$$(1-\alpha)r_a + \alpha G_a = (1-\alpha)r'_a + \alpha G'_a.$$

We consider $\tilde{\nu}_1 = \{x_a = (1-\alpha)r_a + \alpha G_a : a \in [K]\}$ and $\tilde{\nu}_2 = \{x'_a = (1-\alpha)r'_a + \alpha G'_a : a \in [K]\}$.

Now by the classic regret decomposition, we obtain

$$\mathcal{R}_T(\pi, \tilde{\nu}_1) = (T - \mathbb{E}_{\pi, \tilde{\nu}_1}[N_1(T)])\Delta \geq \mathbb{P}_{\pi, \tilde{\nu}_1}^T \left[N_1(T) \leq \frac{T}{2} \right] \frac{T\Delta}{2}.$$

$$\mathcal{R}_T(\pi, \tilde{\nu}_2) = \Delta \mathbb{E}_{\pi, \tilde{\nu}_2}[N_1(T)] + \sum_{a \notin \{1, i\}} 2\Delta \mathbb{E}_{\pi, \tilde{\nu}_2}[N_a(T)] \geq \mathbb{P}_{\pi, \tilde{\nu}_2}^T \left[N_1(T) \geq \frac{T}{2} \right] \frac{T\Delta}{2}.$$

By applying the Bretagnolle–Huber inequality ([55], Theorem 14.2), we have

$$\begin{aligned}\mathcal{R}_T(\pi, \tilde{\nu}_1) + \mathcal{R}_T(\pi, \tilde{\nu}_2) &\geq \frac{T\Delta}{2} \left(\mathbb{P}_{\pi, \tilde{\nu}_1}^T \left[N_1(T) \leq \frac{T}{2} \right] + \mathbb{P}_{\pi, \tilde{\nu}_2}^T \left[N_1(T) \geq \frac{T}{2} \right] \right) \\ &\geq \frac{T\Delta}{4} \exp(-\text{KL}(\mathbb{P}_{\pi, \tilde{\nu}_1}^T \| \mathbb{P}_{\pi, \tilde{\nu}_2}^T))\end{aligned}$$

Based on Lemma C.1, we can get the upper bound of the KL-Divergence between the marginals.

$$\text{KL}(\mathbb{P}_{\pi, \tilde{\nu}_1}^T \| \mathbb{P}_{\pi, \tilde{\nu}_2}^T) \leq 6\epsilon \mathbb{E}_{\pi, \tilde{\nu}_1} \left[\sum_{t=1}^T \text{TV}(x_{a_t} \| x'_{a_t}) \right]$$

Since, $\text{TV}(x_a \| x'_a) = 0$ for $\forall a \in [K]$, $\Delta = \frac{1}{5}\gamma^{k-1}$ and $\gamma = \alpha^{\frac{1}{k}}$. We obtain

$$\mathcal{R}_T(\pi, \tilde{\nu}_1) \geq \Omega(T\alpha^{1-\frac{1}{k}}).$$

Combine Gaussian case, case 1 and case 2, we have

$$\mathcal{R}_T = \Omega \left(\max \left\{ \sqrt{KT}, \left(\frac{K}{\epsilon} \right)^{1-\frac{1}{k}} T^{\frac{1}{k}}, T\alpha^{1-\frac{1}{k}} \right\} \right).$$

□

D Proofs of Section 6.1

Proof of Theorem 6.1. We denote the finite raw moments distribution for rewards by P_k , and denote P_k under α -Huber contamination by $P_{\alpha, k}$. Let $\hat{\mu} = \frac{1}{n} \sum_{\substack{i \in [n] \\ X_i \sim P_{\alpha, k}}} X_i \mathbb{1}_{(|X_i| \leq M)}$ and $\mu = \mathbb{E}_{X_i \sim P_k} [X_i]$.

$$\begin{aligned} & |\tilde{\mu} - \mu| \\ & \leq \left| \text{Lap} \left(\frac{2M}{n\epsilon} \right) \right| + |\hat{\mu} - \mu| \\ & \leq \left| \text{Lap} \left(\frac{2M}{n\epsilon} \right) \right| + \left| \frac{1}{n} \sum_{\substack{i \in [n] \\ X_i \sim P_{\alpha, k}}} X_i \mathbb{1}_{(|X_i| \leq M)} - \mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| \leq M)}] \right| + |\mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| \leq M)}] - \mu| \\ & = \left| \text{Lap} \left(\frac{2M}{n\epsilon} \right) \right| + \left| \frac{1}{n} \sum_{\substack{i \in [n] \\ X_i \sim P_{\alpha, k}}} X_i \mathbb{1}_{(|X_i| \leq M)} - \mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| \leq M)}] \right| + |\mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| > M)}]| \\ & \stackrel{(a)}{\leq} \frac{2M \log(2/\delta)}{n\epsilon} + \left| \frac{1}{n} \sum_{\substack{i \in [n] \\ X_i \sim P_{\alpha, k}}} X_i \mathbb{1}_{(|X_i| \leq M)} - \mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| \leq M)}] \right| + \mathbb{E}_{X_i \sim P_k} [|X_i| \mathbb{1}_{(|X_i| > M)}] \quad \text{w.p. } 1 - \frac{\delta}{2} \\ & \stackrel{(b)}{\leq} \frac{2M \log(2/\delta)}{n\epsilon} + \left| \frac{1}{n} \sum_{\substack{i \in [n] \\ X_i \sim P_{\alpha, k}}} X_i \mathbb{1}_{(|X_i| \leq M)} - \mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| \leq M)}] \right| + (\mathbb{E}_{X_i \sim P_k} [|X_i|^k])^{\frac{1}{k}} (\mathbb{P}_{X_i \sim P_k} (|X_i| > M))^{\frac{k-1}{k}} \\ & \stackrel{(c)}{\leq} \frac{2M \log(2/\delta)}{n\epsilon} + \left| \frac{1}{n} \sum_{\substack{i \in [n] \\ X_i \sim P_{\alpha, k}}} X_i \mathbb{1}_{(|X_i| \leq M)} - \mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| \leq M)}] \right| + \frac{1}{M^{k-1}} \end{aligned}$$

where the inequality (a) follows from Lemma B.5, (b) is from Hölder's Inequality in Lemma B.7 and (c) follows from Markov's inequality in Lemma B.3.

Now we focus on the upper bound of $\left| \frac{1}{n} \sum_{\substack{i \in [n] \\ X_i \sim P_{\alpha, k}}} X_i \mathbb{1}_{(|X_i| \leq M)} - \mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| \leq M)}] \right|$. Let N_G

be the set of indices in n samples distributed according to G , and N_{P_k} be the set of indices in n samples distributed according to P_k . Then

Case 1: uncontaminated case ($\alpha = 0$) Now, the only thing left is to upper bound

$$\left| \frac{1}{n} \sum_{\substack{i \in [n] \\ X_i \sim P_k}} X_i \mathbb{1}_{(|X_i| \leq M)} - \mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| \leq M)}] \right|.$$

For $X_i \sim P_k$, let $Y_i = X_i \mathbb{1}_{(|X_i| \leq M)}$, then $|Y_i| \leq M$ and $\text{Var}(Y_i) = \mathbb{E}[Y_i^2] - (\mathbb{E}[Y_i])^2 \leq \mathbb{E}[Y_i^2] \leq \mathbb{E}_{X_i \sim P_k} [X_i^2] \leq 1$. Then, from Bernstein's inequality in Lemma B.8, we have with probability $1 - \delta/2$

$$\left| \frac{1}{n} \sum_{\substack{i \in [n] \\ X_i \sim P_k}} X_i \mathbb{1}_{(|X_i| \leq M)} - \mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| \leq M)}] \right| \leq \sqrt{\frac{2 \log(4/\delta)}{n}} + \frac{4M \log(4/\delta)}{3n}. \quad (3)$$

Then we get with probability at least $1 - \delta$,

$$|\tilde{\mu} - \mu| \leq \sqrt{\frac{2 \log(4/\delta)}{n}} + \frac{4M \log(4/\delta)}{3n} + \frac{2M \log(2/\delta)}{n\epsilon} + \frac{1}{M^{k-1}}.$$

For $\epsilon > 0$ and $\delta \in (0, 1)$, we have with probability at least $1 - \delta$,

$$|\tilde{\mu} - \mu| \leq \sqrt{\frac{2 \log(4/\delta)}{n}} + \frac{4M \log(4/\delta)}{n\epsilon} + \frac{1}{M^{k-1}}.$$

Taking the truncation threshold $M = \left(\frac{n\epsilon}{4 \log(4/\delta)}\right)^{\frac{1}{k}}$, we have

$$|\tilde{\mu} - \mu| \leq \sqrt{\frac{2 \log(4/\delta)}{n}} + 2 \left(\frac{4 \log(4/\delta)}{n\epsilon}\right)^{1-\frac{1}{k}}.$$

Case 2: contaminated case ($\alpha \in (0, \frac{1}{2})$)

$$\begin{aligned} & \left| \frac{1}{n} \sum_{\substack{i \in [n] \\ X_i \sim P_{\alpha,k}}} X_i \mathbb{1}_{(|X_i| \leq M)} - \mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| \leq M)}] \right| \\ &= \left| \frac{1}{n} \sum_{i \in N_G} X_i \mathbb{1}_{(|X_i| \leq M)} + \frac{1}{n} \sum_{i \in N_{P_k}} X_i \mathbb{1}_{(|X_i| \leq M)} - \mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| \leq M)}] \right| \\ &\leq \underbrace{\left| \frac{1}{n} \sum_{i \in N_G} X_i \mathbb{1}_{(|X_i| \leq M)} \right|}_{T_1} + \underbrace{\left| \frac{1}{n} \sum_{i \in N_{P_k}} X_i \mathbb{1}_{(|X_i| \leq M)} - \mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| \leq M)}] \right|}_{T_2}. \end{aligned}$$

To control T_1 , we can write it as

$$\begin{aligned} T_1 &= \left| \frac{1}{n} \sum_{i \in N_G} X_i \mathbb{1}_{(|X_i| \leq M)} \right| \\ &\leq \frac{1}{n} \sum_{i \in N_G} |X_i| \mathbb{1}_{(|X_i| \leq M)} \\ &\leq \frac{|N_G|}{n} M. \end{aligned}$$

Then $\frac{|N_G|}{n}$ can be treat as a mean estimation of Bernoulli distribution $Ber(\alpha)$. Then based on Bernstein's inequality in Lemma B.8, we get with probability $1 - \delta/4$,

$$\left| \frac{|N_G|}{n} - \alpha \right| \leq \sqrt{\frac{2\alpha(1-\alpha) \log(8/\delta)}{n}} + \frac{2 \log(8/\delta)}{3n} \leq \sqrt{\frac{2\alpha_1(1-\alpha_1) \log(8/\delta)}{n}} + \frac{2 \log(8/\delta)}{3n}$$

for $\alpha \leq \alpha_1 \in (0, 1/2)$.

Thus,

$$T_1 \leq \left(\alpha_1 + \sqrt{\frac{2\alpha_1 \log(8/\delta)}{n}} + \frac{2 \log(8/\delta)}{3n} \right) M. \quad \text{with probability } 1 - \delta/4$$

When $n \geq \frac{\log(8/\delta)}{\alpha_1}$, we have

$$T_1 \leq 4\alpha_1 M.$$

To bound T_2 , we have

$$\begin{aligned} & \left| \frac{1}{n} \sum_{i \in N_{P_k}} X_i \mathbb{1}_{(|X_i| \leq M)} - \mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| \leq M)}] \right| \\ &= \left| \frac{1}{n} \sum_{\substack{i \in N_G \cup N_{P_k} \\ X_i \sim P_k}} X_i \mathbb{1}_{(|X_i| \leq M)} - \frac{1}{n} \sum_{\substack{i \in N_G \\ X_i \sim P_k}} X_i \mathbb{1}_{(|X_i| \leq M)} - \mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| \leq M)}] \right| \\ &\leq \left| \frac{1}{n} \sum_{\substack{i \in N_G \\ X_i \sim P_k}} X_i \mathbb{1}_{(|X_i| \leq M)} \right| + \left| \frac{1}{n} \sum_{\substack{i \in [n] \\ X_i \sim P_k}} X_i \mathbb{1}_{(|X_i| \leq M)} - \mathbb{E}_{X_i \sim P_k} [X_i \mathbb{1}_{(|X_i| \leq M)}] \right| \\ &\leq 4\alpha_1 M + \sqrt{\frac{2 \log(16/\delta)}{n}} + \frac{4M \log(16/\delta)}{3n} \quad \text{w.p. } 1 - \delta/4 \end{aligned}$$

where the last inequality is based on the similar analysis of T_1 and the inequality of (3).

Put everything together, we have with probability at least $1 - \delta$,

$$|\tilde{\mu} - \mu| \leq \sqrt{\frac{2 \log(16/\delta)}{n}} + \frac{4M \log(16/\delta)}{3n} + 8\alpha_1 M + \frac{2M \log(2/\delta)}{n\epsilon} + \frac{1}{M^{k-1}}.$$

Thus, for $\epsilon > 0$, we have

$$|\tilde{\mu} - \mu| \leq \sqrt{\frac{2 \log(16/\delta)}{n}} + \frac{4M \log(16/\delta)}{n\epsilon} + \frac{1}{M^{k-1}} + 8\alpha_1 M.$$

Taking $M = \min \left\{ \left(\frac{n\epsilon}{4 \log(16/\delta)} \right)^{\frac{1}{k}}, (8\alpha_1)^{-\frac{1}{k}} \right\}$, we have

$$|\tilde{\mu} - \mu| \leq \sqrt{\frac{2 \log(16/\delta)}{n}} + 2 \left(\frac{4 \log(16/\delta)}{n\epsilon} \right)^{1-\frac{1}{k}} + 2(8\alpha_1)^{1-\frac{1}{k}}.$$

□

Proof of Theorem 6.2. Let τ_0 be the maximal epoch such that $B_\tau < \frac{\log(16|\mathcal{S}|\tau^2/\delta)}{\alpha_1}$.

For all epoch $\tau \leq \tau_0$, the batch size is less than 2^{τ_0} . Since batch size doubles, until epoch τ_0 , we have the number of pulls for each arm $a \in [K]$ is less than $2 \cdot 2^{\tau_0} \leq 2 \frac{\log(16|\mathcal{S}|\tau_0^2/\delta)}{\alpha_1}$. Then the regret has to suffer $\frac{2 \log(16|\mathcal{S}|\tau_0^2/\delta)}{\alpha_1} \Delta_a$ for each $a \in [K]$.

For $\tau > \tau_0$, $B_\tau \geq \frac{\log(16|\mathcal{S}|\tau^2/\delta)}{\alpha_1}$. For each $a \in \mathcal{S}$, from Theorem 6.1, we have with probability at least $1 - \frac{\delta}{2|\mathcal{S}|\tau^2}$,

$$|\tilde{\mu}_a - \mu_a| \leq \beta_\tau.$$

Given an epoch $\tau > \tau_0$, we denote by \mathcal{E}_τ the event where for all $a \in \mathcal{S}$ it holds that $|\tilde{\mu}_a - \mu_a| \leq \beta_\tau$. and denote $\mathcal{E} = \cup_{\tau > \tau_0} \mathcal{E}_\tau$. By taking union bound, we have

$$\mathbb{P}(\mathcal{E}_\tau) \geq 1 - \frac{\delta}{2\tau^2},$$

and

$$\mathbb{P}(\mathcal{E}) \geq 1 - \frac{\delta}{2} \left(\sum_{\tau > \tau_0} \tau^{-2} \right) \geq 1 - \delta.$$

In the following, we condition on the good event \mathcal{E} . We first show that the optimal arm a^* is never eliminated. For any epoch $\tau > \tau_0$, let $a_\tau = \arg \max_{a \in \mathcal{S}} \tilde{\mu}_a$. Since

$$(\tilde{\mu}_{a_\tau} - \tilde{\mu}_{a^*}) + \Delta_{a_\tau} = |(\tilde{\mu}_{a_\tau} - \tilde{\mu}_{a^*}) + \Delta_{a_\tau}| \leq |\tilde{\mu}_{a_\tau} - \mu_{a_\tau}| + |\tilde{\mu}_{a^*} - \mu_{a^*}| \leq 2\beta_\tau,$$

it is easy to see that the algorithm doesn't eliminate a^* .

Then, we show that at the end of epoch $\tau > \tau_0$, all arms such that $\Delta_a \geq 4\beta_\tau$ will be eliminated. To show this, we have that under good event \mathcal{E} ,

$$\tilde{\mu}_a + \beta_\tau \leq \mu_a + 2\beta_\tau < \mu_{a^*} - 4\beta_\tau + 2\beta_\tau \leq \tilde{\mu}_{a^*} - \beta_\tau \leq \tilde{\mu}_{a_\tau} - \beta_\tau$$

which implies that arm a will be eliminated by the rule. Thus, for each sub-optimal arm a , let $\tau(a)$ be the last epoch that arm a is not eliminated. By the above result, we have

$$\Delta_a \leq 4\beta_{\tau(a)} = 4\sqrt{\frac{2\log(16|\mathcal{S}|\tau(a)^2/\delta)}{B_{\tau(a)}}} + 8 \left(\frac{4\log(16|\mathcal{S}|\tau(a)^2/\delta)}{B_{\tau(a)}\epsilon} \right)^{1-\frac{1}{k}} + 8(8\alpha_1)^{1-\frac{1}{k}}.$$

We divide the arms $a \in [K]$ into two groups: $\mathcal{G}_1 = \{a \in [K] : 16(8\alpha_1)^{1-\frac{1}{k}} \leq \Delta_a\}$ and $\mathcal{G}_2 = \{a \in [K] : 16(8\alpha_1)^{1-\frac{1}{k}} \geq \Delta_a\}$.

Group 1: Now, for all arm $a \in \mathcal{G}_1$, we have

$$\Delta_a \leq 8\sqrt{\frac{2\log(16|\mathcal{S}|\tau(a)^2/\delta)}{B_{\tau(a)}}} + 16 \left(\frac{4\log(16|\mathcal{S}|\tau(a)^2/\delta)}{B_{\tau(a)}\epsilon} \right)^{1-\frac{1}{k}}.$$

Hence, we have

$$B_{\tau(a)} \leq \max \left\{ \frac{128\log(16|\mathcal{S}|\tau(a)^2/\delta)}{\Delta_a^2}, \frac{4\log(16|\mathcal{S}|\tau(a)^2/\delta)}{\epsilon} \left(\frac{16}{\Delta_a} \right)^{\frac{k}{k-1}}, \frac{\log(16|\mathcal{S}|\tau_0^2/\delta)}{\alpha_1} \right\}.$$

Since $|\mathcal{S}| \leq K$ and $2^\tau \leq T$ for any τ . Thus,

$$B_{\tau(a)} \leq \max \left\{ \frac{128\log(16K\log^2 T/\delta)}{\Delta_a^2}, \frac{4\log(16K\log^2 T/\delta)}{\epsilon} \left(\frac{16}{\Delta_a} \right)^{\frac{k}{k-1}}, \frac{\log(16K\log^2 T/\delta)}{\alpha_1} \right\},$$

Since the batch size doubles, we have $N_a(T) \leq 2B_{\tau(a)}$ for each sub-optimal arm a . Therefore, for all arm $a \in \mathcal{G}_1$,

$$\mathcal{R}_T = \sum_{a \in \mathcal{G}_1} N_a(T)\Delta_a \leq 2B_{\tau(a)}\Delta_a.$$

Let η be a number in $(0, 1)$. For all arms $a \in \mathcal{G}_1$ with $\Delta_a \leq \eta$, the regret incurred by pulling these arms is upper bounded by $T\eta$. For any arm $a \in \mathcal{G}_1$ with $\Delta_a > \eta$, choose $\delta = \frac{1}{T}$ and assume $T \geq K$, then the expected regret incurred by pulling arm a is upper bounded by

$$\begin{aligned} \mathbb{E} \left[\sum_{a \in \mathcal{G}_1, \Delta_a > \eta} \Delta_a N_a(T) \right] &\leq \mathbb{P}(\bar{\mathcal{E}}) \cdot T + O \left(\sum_{a \in \mathcal{G}_1, \Delta_a > \eta} \left\{ \frac{\log T}{\Delta_a} + \frac{\log T}{\epsilon} \left(\frac{1}{\Delta_a} \right)^{\frac{1}{k-1}} + \frac{\log T}{\alpha_1} \Delta_a \right\} \right) \\ &\leq O \left(\frac{K \log T}{\eta} + \frac{K \log T}{\epsilon \eta^{\frac{1}{k-1}}} + \frac{K \log T}{\alpha_1} \right) \end{aligned}$$

where the last term in the last inequality is based on following result: from the heavy-tailed assumption for rewards distributions in (1), we have for any $a \in [K]$, $|\mu_a| \leq \mathbb{E}_{r_a \sim P_k} |r_a| \leq \mathbb{E}_{r_a \sim P_k} |r_a|^k \leq 1$, so $\Delta_a = \mu_{a^*} - \mu_a \leq 2$.

Thus the regret from group 1 is at most

$$T\eta + O\left(\frac{K \log T}{\eta} + \frac{K \log T}{\epsilon\eta^{\frac{1}{k-1}}} + \frac{K \log T}{\alpha_1}\right).$$

Taking $\eta = \max\left\{\sqrt{\frac{K \log T}{T}}, \left(\frac{K \log T}{T\epsilon}\right)^{\frac{k-1}{k}}\right\}$, the regret from group 1 is at most

$$O\left(\sqrt{KT \log T} + \left(\frac{K \log T}{\epsilon}\right)^{\frac{k-1}{k}} T^{\frac{1}{k}} + \frac{K \log T}{\alpha_1}\right)$$

Group 2: For all other arms $a \in \mathcal{G}_2$, we have the total regret is at most $O(T\Delta_a) = O(T\alpha_1^{1-\frac{1}{k}})$.

Combine the two groups, choose $\delta = \frac{1}{T}$ and assume $T \geq K$, we have the that the expected regret satisfies,

$$\mathcal{R}_T \leq O\left(\sqrt{KT \log T} + \left(\frac{K \log T}{\epsilon}\right)^{\frac{k-1}{k}} T^{\frac{1}{k}} + \frac{K \log T}{\alpha_1} + T\alpha_1^{1-\frac{1}{k}}\right)$$

We also give privacy guarantee for the algorithm. Based on Laplacian mechanism in Definition 3.5 and Post-processing in Lemma B.1, we can get that Algorithm 1 is ϵ -DP. \square

E Proofs of Section 6.2

Proof of Theorem 6.5. Step 1: we will show that with high probability $1 - \delta/2$, $|J - \mu| \leq 2r$.

To this end, we first study the private histogram. Note $\mathbb{E}[\mathbb{1}(X_i \in B_j)] = P_{\alpha,k}(B_j)$, then

$$\begin{aligned} \mathbb{P}(|\tilde{p}_j - P_{\alpha,k}(B_j)| > t) &= \mathbb{P}\left(\left|\frac{\sum_{i=1}^n \mathbb{1}(X_i \in B_j)}{n} + \text{Lap}\left(\frac{2}{n\epsilon}\right) - P_{\alpha,k}(B_j)\right| > t\right) \\ &\leq \mathbb{P}\left(\left|\frac{\sum_{i=1}^n \mathbb{1}(X_i \in B_j)}{n} - P_{\alpha,k}(B_j)\right| > t/2\right) + \mathbb{P}\left(\left|\text{Lap}\left(\frac{2}{n\epsilon}\right)\right| > t/2\right) \\ &\leq 2 \exp\left(-\frac{nt^2}{2}\right) + \exp\left(-\frac{n\epsilon t}{4}\right), \end{aligned}$$

where the last inequality is from Lemma B.6 and Lemma B.5. By a union bound over j , we further have

$$\begin{aligned} \mathbb{P}\left(\max_{j \in \mathcal{J}} |\tilde{p}_j - P_{\alpha,k}(B_j)| > t\right) &\leq \frac{2D}{r} \left(2 \exp\left(-\frac{nt^2}{2}\right) + \exp\left(-\frac{n\epsilon t}{4}\right)\right) \\ &\leq 2D \left(2 \exp\left(-\frac{nt^2}{2}\right) + \exp\left(-\frac{n\epsilon t}{4}\right)\right) \end{aligned}$$

Thus, we have with probability $1 - \delta/2$,

$$\max_{j \in \mathcal{J}} |\tilde{p}_j - P_{\alpha,k}(B_j)| \leq \max\left\{\sqrt{\frac{2 \ln \frac{16D}{\delta}}{n}}, \frac{4 \ln \frac{16D}{\delta}}{n\epsilon}\right\} := C_1.$$

In the following, we condition on the above event. Next, by Chebyshev's inequality in Lemma B.4 and the assumption of \mathcal{P}_k that k -th central moment is less than 1, we have

$$\begin{aligned} \mathbb{P}_{X \sim \mathcal{P}_{\alpha,k}}(|X - \mu| \geq r) &\leq \alpha \mathbb{P}_{X \sim \mathcal{G}}(|X - \mu| \geq r) + (1 - \alpha) \mathbb{P}_{X \sim \mathcal{P}_k}(|X - \mu| \geq r) \\ &\leq \alpha + (1 - \alpha)(1/r)^k := C_2 \end{aligned} \quad (4)$$

Let j^* is the index of the bin containing the true mean μ and we consider three consecutive intervals $A_{j^*} = B_{j^*-1} \cup B_{j^*} \cup B_{j^*+1}$

$$\begin{aligned} P_{\alpha,k}(A_{j^*}) &= P_{\alpha,k}(B_{j^*-1}) + P_{\alpha,k}(B_{j^*}) + P_{\alpha,k}(B_{j^*+1}) \\ &\geq P_{\alpha,k}((\mu - r, \mu + r)) \\ &\geq 1 - C_2. \end{aligned}$$

where the first inequality is from inequality (4). Now, for any $j \notin \{j^* - 1, j^*, j^* + 1\}$, we have when $D \geq 2r$

$$\tilde{p}_j \leq P_{\alpha,k}(B_j) + C_1 \leq 1 - P_{\alpha,k}(A_{j^*}) + C_1 \leq C_2 + C_1.$$

On the other hand, since $P_{\alpha,k}(A_{j^*}) \geq 1 - C_2$, there must exist some $j \in \{j^* - 1, j^*, j^* + 1\}$ such that $P_{\alpha,k}(B_j) \geq \frac{1-C_2}{3}$. Therefore, for this j , we have

$$\tilde{p}_j \geq P_{\alpha,k}(B_j) - C_1 \geq \frac{1 - C_2}{3} - C_1.$$

Therefore, if n (depending on α, ϵ, r) such that $\frac{1-C_2}{3} - C_1 > C_2 + C_1$, the true mean μ is in the bin chosen by line 3 in Algorithm 3 or it's neighboring bin, which implies that with probability at least $1 - \delta/2$, $|J - \mu| \leq 2r$.

Step 2: Utilizing the above result, we aim to show that truncation can handle heavy-tail, privacy and robustness in the concentration.

$$\begin{aligned} |\tilde{\mu} - \mu| &= \left| J + \frac{1}{n} \sum_{i=n+1}^{2n} (X_i - J) \mathbb{1}(|X_i - J| \leq M) + \text{Lap}\left(\frac{2M}{n\epsilon}\right) - \mu \right| \\ &= \left| \frac{1}{n} \sum_{i=n+1}^{2n} (X_i - J) \mathbb{1}(|X_i - J| \leq M) + \text{Lap}\left(\frac{2M}{n\epsilon}\right) \right. \\ &\quad \left. + \frac{1}{n} \sum_{i=n+1}^{2n} (J - \mu) \{ \mathbb{1}(|X_i - J| \leq M) + \mathbb{1}(|X_i - J| > M) \} \right| \\ &= \left| \frac{1}{n} \sum_{i=n+1}^{2n} (X_i - J + J - \mu) \mathbb{1}(|X_i - J| \leq M) + \text{Lap}\left(\frac{2M}{n\epsilon}\right) + \frac{1}{n} \sum_{i=n+1}^{2n} (J - \mu) \mathbb{1}(|X_i - J| > M) \right| \\ &\leq \left| \frac{1}{n} \sum_{i=n+1}^{2n} (X_i - \mu) \mathbb{1}(|X_i - J| \leq M) \right| + \left| \text{Lap}\left(\frac{2M}{n\epsilon}\right) \right| + \left| \frac{1}{n} \sum_{i=n+1}^{2n} (J - \mu) \mathbb{1}(|X_i - J| > M) \right| \end{aligned}$$

We first focus on the first term in the right hand of the last inequality. Let N_G be the set of indices in n samples distributed according to $G \in \mathcal{G}$, and N_{P_k} be the set of indices in n samples distributed according to $P_k \in \mathcal{P}_k^c$. Then, we have

$$\begin{aligned} &\left| \frac{1}{n} \sum_{i=n+1}^{2n} (X_i - \mu) \mathbb{1}(|X_i - J| \leq M) \right| \\ &\leq \underbrace{\left| \frac{1}{n} \sum_{i \in N_G} (X_i - \mu) \mathbb{1}(|X_i - J| \leq M) \right|}_{T_1} + \underbrace{\left| \frac{1}{n} \sum_{i \in N_{P_k}} (X_i - \mu) \mathbb{1}(|X_i - J| \leq M) \right|}_{T_2}. \end{aligned}$$

To control T_1 , we can write it as

$$\begin{aligned}
T_1 &= \left| \frac{1}{n} \sum_{i \in N_G} (X_i - \mu) \mathbb{1}(|X_i - J| \leq M) \right| \\
&\leq \frac{1}{n} \sum_{i \in N_G} |(X_i - \mu)| \mathbb{1}(|X_i - J| \leq M) \\
&\leq \frac{1}{n} \sum_{i \in N_G} |(X_i - \mu)| \mathbb{1}(|X_i - \mu| \leq M + 2r) \\
&\leq \frac{|N_G|}{n} (M + 2r).
\end{aligned}$$

Then $\frac{|N_G|}{n}$ can be treated as a mean estimation of Bernoulli distribution $Ber(\alpha)$. Then based on Bernstein's inequality in Lemma B.8, we get with probability $1 - \delta/8$,

$$\left| \frac{|N_G|}{n} - \alpha \right| \leq \sqrt{\frac{2\alpha(1-\alpha)\log(16/\delta)}{n}} + \frac{2\log(16/\delta)}{3n} \leq \sqrt{\frac{2\alpha_1(1-\alpha_1)\log(16/\delta)}{n}} + \frac{2\log(16/\delta)}{3n}$$

for $\alpha \leq \alpha_1 \in (0, 1/2)$.

Thus,

$$T_1 \leq \left(\alpha_1 + \sqrt{\frac{2\alpha_1\log(16/\delta)}{n}} + \frac{2\log(16/\delta)}{3n} \right) (M + 2r), \quad \text{with probability } 1 - \delta/8.$$

Thus, if n satisfies $\sqrt{\frac{2\alpha_1\log(16/\delta)}{n}} + \frac{2\log(16/\delta)}{3n} = O(\alpha_1)$, then we have $T_1 = O(\alpha_1(M + 2r))$. Now, we bound T_2 ,

$$\begin{aligned}
T_2 &= \left| \frac{1}{n} \sum_{\substack{i \in N_G \cup N_{P_k} \\ X_i \sim P_k}} (X_i - \mu) \mathbb{1}(|X_i - J| \leq M) - \frac{1}{n} \sum_{\substack{i \in N_G \\ X_i \sim P_k}} (X_i - \mu) \mathbb{1}(|X_i - J| \leq M) \right| \\
&\leq \left| \frac{1}{n} \sum_{\substack{i \in N_G \cup N_{P_k} \\ X_i \sim P_k}} (X_i - \mu) \mathbb{1}(|X_i - J| \leq M) \right| + \left| \frac{1}{n} \sum_{\substack{i \in N_G \\ X_i \sim P_k}} (X_i - \mu) \mathbb{1}(|X_i - J| \leq M) \right| \\
&\leq \left| \frac{1}{n} \sum_{\substack{i \in [n] \\ X_i \sim P_k}} (X_i - \mu) \mathbb{1}(|X_i - J| \leq M) \right| + T_1.
\end{aligned}$$

Now we focus on the upper bound of $\left| \frac{1}{n} \sum_{\substack{i \in [n] \\ X_i \sim P_k}} (X_i - \mu) \mathbb{1}(|X_i - J| \leq M) \right|$. With probability $1 - \delta/8$,

$$\begin{aligned}
& \left| \frac{1}{n} \sum_{\substack{i \in [n] \\ X_i \sim P_k}} (X_i - \mu) \mathbb{1}(|X_i - J| \leq M) \right| \\
& \leq \left| \frac{1}{n} \sum_{\substack{i \in [n] \\ X_i \sim P_k}} (X_i - \mu) \mathbb{1}(|X_i - J| \leq M) - \mathbb{E}[(X_1 - \mu) \mathbb{1}(|X_1 - J| \leq M)] \right| \\
& \quad + |\mathbb{E}[(X_1 - \mu) \mathbb{1}(|X_1 - J| \leq M)] - \mathbb{E}[(X_1 - \mu)]| \\
& \leq \sqrt{\frac{2 \log(16/\delta)}{n}} + \frac{4(M + 2r) \log(16/\delta)}{3n} + |\mathbb{E}[(X_i - \mu) \mathbb{1}(|X_i - J| \geq M)]| \\
& \leq \sqrt{\frac{2 \log(16/\delta)}{n}} + \frac{4(M + 2r) \log(16/\delta)}{3n} + (\mathbb{E}[|X_i - \mu|^k])^{\frac{1}{k}} (\mathbb{P}(|X_i - \mu| \geq M - 2r))^{\frac{k-1}{k}} \\
& \leq \sqrt{\frac{2 \log(16/\delta)}{n}} + \frac{4(M + 2r) \log(16/\delta)}{3n} + \frac{1}{(M - 2r)^{k-1}} \\
& \leq \sqrt{\frac{2 \log(16/\delta)}{n}} + \frac{4(M + 2r) \log(16/\delta)}{3n} + \left(\frac{2}{M}\right)^{k-1}
\end{aligned}$$

where the last inequality follows from $M \geq 4r$, the third inequality follows from Hölder's Inequality in Lemma B.7 and the second inequality follows from Bernstein inequality in Lemma B.8. That is, let

$$Y_i = (X_i - \mu) \mathbb{1}(|X_i - J| \leq M),$$

then

$$\begin{aligned}
|Y_i - \mathbb{E}[Y_i]| & \leq |Y_i| + |\mathbb{E}[Y_i]| \\
& \leq |X_i - \mu| \mathbb{1}(|X_i - \mu| \leq M + 2r) + \mathbb{E}[|X_i - \mu| \mathbb{1}(|X_i - \mu| \leq M + 2r)] \\
& \leq 2(M + 2r)
\end{aligned}$$

and

$$\begin{aligned}
\text{Var}(Y_i - \mathbb{E}[Y_i]) & = \mathbb{E}(Y_i - \mathbb{E}[Y_i])^2 \leq \mathbb{E}[Y_i^2] \\
& \leq \mathbb{E}_{X_i \sim P_k} [(X_i - \mu)^2 \mathbb{1}(|X_i - \mu| \leq M + 2r)] \\
& \leq \mathbb{E}_{X_i \sim P_k} [(X_i - \mu)^2] \leq 1.
\end{aligned}$$

Therefore, with probability $1 - 3\delta/8$,

$$T_2 \leq \sqrt{\frac{2 \log(16/\delta)}{n}} + \frac{4(M + 2r) \log(16/\delta)}{3n} + \left(\frac{2}{M}\right)^{k-1} + T_1.$$

Now, we focus on the upper bound of $T_3 := \left| \frac{1}{n} \sum_{i=n+1}^{2n} (J - \mu) \mathbb{1}(|X_i - J| > M) \right|$.

$$\begin{aligned}
& \left| \frac{1}{n} \sum_{i=n+1}^{2n} (J - \mu) \mathbb{1}(|X_i - J| > M) \right| \\
& \leq \frac{1}{n} \sum_{i=n+1}^{2n} |J - \mu| \mathbb{1}(|X_i - J| > M) \\
& \leq 2r \frac{\sum_{i=n+1}^{2n} \mathbb{1}(|X_i - J| > M)}{n} \\
& \leq 2r \frac{\sum_{i=n+1}^{2n} \mathbb{1}(|X_i - \mu| > M - 2r)}{n}
\end{aligned}$$

where

$$\begin{aligned}\mathbb{E}_{X_i \sim P_{k,\alpha}}[\mathbb{1}(|X_i - \mu| > M - 2r)] &= \mathbb{P}_{X_i \sim P_{k,\alpha}}(|X_i - \mu| > M - 2r) \\ &\leq \alpha + (1 - \alpha)\mathbb{P}_{X_i \sim P_k}(|X_i - \mu| > M - 2r) \\ &\leq \alpha + \frac{1}{(M - 2r)^k} \leq \alpha_1 + \left(\frac{2}{M}\right)^k\end{aligned}$$

By Hoeffding's inequality, we have with probability $1 - \delta/8$,

$$\frac{\sum_{i=n+1}^{2n} \mathbb{1}(|X_i - \mu| > M - 2r)}{n} \leq \mathbb{P}_{X_i \sim P_{k,\alpha}}(|X_i - \mu| > M - 2r) + \sqrt{\frac{\log(16/\delta)}{2n}}.$$

Thus, we have

$$T_3 = \left| \frac{1}{n} \sum_{i=n+1}^{2n} (J - \mu) \mathbb{1}(|X_i - J| > M) \right| \leq 2r \left(\alpha_1 + \left(\frac{2}{M}\right)^k + \sqrt{\frac{\log(16/\delta)}{2n}} \right).$$

Putting everything together, we have

$$\begin{aligned}|\tilde{\mu} - \mu| &= O \left(\left(\alpha_1 + \sqrt{\frac{2\alpha_1 \log(16/\delta)}{n}} + \frac{2 \log(16/\delta)}{3n} \right) (M + 2r) \right) \\ &\quad + O \left(\sqrt{\frac{2 \log(16/\delta)}{n}} + \frac{4(M + 2r) \log(16/\delta)}{3n} + \left(\frac{2}{M}\right)^{k-1} \right) \\ &\quad + O \left(2r \left(\alpha_1 + \left(\frac{2}{M}\right)^k + \sqrt{\frac{\log(16/\delta)}{2n}} \right) \right) \\ &\quad + O \left(\frac{M \log(1/\delta)}{n\epsilon} \right)\end{aligned}$$

Case I: $\alpha = 0$, Uncontaminated concentration. We want to show that our concentration is better than medians-of-mean in [38] (Theorem 3.5). That is, we are additive for their third term therein (i.e., $\log(D) + \log(1/\delta)$), while they are multiplicative.

In this case, our $C_2 = (1/r)^k$, and by our first condition on n , it need to satisfy $6C_1 + 4C_2 < 1$. This implies that $C_2 < 1/4$. Thus, setting $r = 10^{1/k}$ is sufficient. Hence, we have $C_1 < 0.1$, which requires n to satisfy $n \geq 200 \log(16D/\delta)$ and $n \geq 20 \log(16D/\delta)/\epsilon$. We can safely set $n \geq 200 \log(16D/\delta)/\epsilon$.

In the case of $\alpha = 0$, T_1 is not a problem, which only introduces another $O(M \log(1/\delta)/n)$. T_3 is also not a problem which is dominated by $O((2/M)^{k-1} + \sqrt{\log(1/\delta)}/\sqrt{n})$

Let's summarize all the values: when $\alpha = 0$, $r = 10^{1/k}$ and $n \geq 200 \log(16D/\delta)/\epsilon$, we have

$$\begin{aligned}|\tilde{\mu} - \mu| &= O \left(\sqrt{\frac{2 \log(16/\delta)}{n}} + \frac{M \log(16/\delta)}{3n} + \left(\frac{2}{M}\right)^{k-1} \right) \\ &\quad + O \left(\frac{M \log(1/\delta)}{n\epsilon} \right) \\ &= O \left(\sqrt{\frac{2 \log(16/\delta)}{n}} + \frac{M \log(16/\delta)}{\epsilon n} + \left(\frac{2}{M}\right)^{k-1} \right)\end{aligned}$$

Now, we need to choose M to minimize the above while satisfying $M \geq 4r$. By standard choice, we set $M = 4 \left(\frac{n\epsilon}{\log(1/\delta)} \right)^{1/k}$, which satisfies $M \geq 4r$ when $n \geq \frac{10 \log(1/\delta)}{\epsilon}$.

Case II: $\alpha > 0$ and $\alpha \in (0, 1/2)$. Contaminated concentration. We want to minimize the term $T(\alpha, \epsilon)$ while maximizing the possible range of α .

In this case, $C_2 = \alpha + (1 - \alpha)(1/r)^k$ and again we need to satisfy that $6C_1 + 4C_2 < 1$, which first implies that α needs to be $\alpha < 1/4$. Setting $r = \iota^{1/k}$, we have $C_2 = \alpha + \frac{1}{\iota}(1 - \alpha)$, which needs to be less than $1/4$. Let's set $\iota = \frac{1-\alpha}{0.249-\alpha}$ (hence $\alpha < 0.249$), we have there exists an absolute constant c_1 such that when $n \geq c_1 \log(16D/\delta)/\epsilon$, we guarantee $6C_1 + 4C_2 < 1$.

Now, we turn to T_1 . If $n \geq \frac{\log(16/\delta)}{\alpha_1}$ and $M \geq 4r$, we have $T_1 = O(\alpha_1 M)$.

For T_3 , we have

$$T_3 = 2r \left(\alpha_1 + \left(\frac{2}{M} \right)^k + \sqrt{\frac{\log(16/\delta)}{2n}} \right)$$

One simple way is to set $n \geq \log(16/\delta)/\alpha_1^2$. Then, we have $T_3 = O(\alpha_1 M + (1/M)^{k-1})$.

Let's summarize it. For any $\alpha \in (0, 0.249)$, setting $r = \left(\frac{1-\alpha}{0.249-\alpha} \right)^{1/k}$. Then, for all $n \geq \max\{c_1 \log(16D/\delta)/\epsilon, \log(16/\delta)/\alpha_1^2\}$, we have

$$|\tilde{\mu} - \mu| = O \left(\sqrt{\frac{2 \log(16/\delta)}{n}} + \frac{M \log(16/\delta)}{\epsilon n} + \left(\frac{2}{M} \right)^{k-1} + \alpha_1 M \right)$$

Now, we need to choose M to minimize the above while satisfying $M \geq 4r$. By standard choice, we set $M = \min\left\{4 \left(\frac{n\epsilon}{\log(1/\delta)} \right)^{1/k}, 4\alpha_1^{-1/k}\right\}$, which satisfies $M \geq 4r$ when n and α satisfy

$$n \geq \frac{\iota \log(1/\delta)}{\epsilon} \quad \text{and} \quad \frac{1}{\alpha} \geq \iota,$$

where recall that $\iota = \frac{1-\alpha}{0.249-\alpha}$. Hence, we only have a valid concentration for $\alpha \in (0, 0.133)$. □

Proof of Theorem 6.10. Let τ_0 be the maximal epoch such that $B_\tau < \frac{200 \log(16D|\mathcal{S}|\tau^2/\delta)}{\epsilon}$.

For all epoch $\tau \leq \tau_0$, the batch size is less than 2^{τ_0} . Since batch size doubles, until epoch τ_0 , we have the number of pulls for each arm $a \in [K]$ is less than $2 \cdot 2^{\tau_0} \leq 2 \frac{200 \log(16D|\mathcal{S}|\tau_0^2/\delta)}{\epsilon}$. Then the regret has to suffer $\frac{400 \log(16D|\mathcal{S}|\tau_0^2/\delta)}{\epsilon} \Delta_a$ for each $a \in [K]$.

For $\tau > \tau_0$, $B_\tau \geq \frac{200 \log(16D|\mathcal{S}|\tau^2/\delta)}{\epsilon}$. For each $a \in \mathcal{S}$, from Corollary 6.6, we have with probability at least $1 - \frac{\delta}{2|\mathcal{S}|\tau^2}$,

$$|\tilde{\mu}_a - \mu_a| \leq \beta_\tau.$$

Given an epoch $\tau > \tau_0$, we denote by \mathcal{E}_τ the event where for all $a \in \mathcal{S}$ it holds that $|\tilde{\mu}_a - \mu_a| \leq \beta_\tau$. and denote $\mathcal{E} = \cup_{\tau > \tau_0} \mathcal{E}_\tau$. By taking union bound, we have

$$\mathbb{P}(\mathcal{E}_\tau) \geq 1 - \frac{\delta}{2\tau^2},$$

and

$$\mathbb{P}(\mathcal{E}) \geq 1 - \frac{\delta}{2} \left(\sum_{\tau > \tau_0} \tau^{-2} \right) \geq 1 - \delta.$$

In the following, we condition on the good event \mathcal{E} . We first show that the optimal arm a^* is never eliminated. For any epoch $\tau > \tau_0$, let $a_\tau = \arg \max_{a \in \mathcal{S}} \tilde{\mu}_a$. Since

$$(\tilde{\mu}_{a_\tau} - \tilde{\mu}_{a^*}) + \Delta_{a_\tau} = |(\tilde{\mu}_{a_\tau} - \tilde{\mu}_{a^*}) + \Delta_{a_\tau}| \leq |\tilde{\mu}_{a_\tau} - \mu_{a_\tau}| + |\tilde{\mu}_{a^*} - \mu_{a^*}| \leq 2\beta_\tau,$$

it is easy to see that the algorithm doesn't eliminate a^* .

Then, we show that at the end of epoch $\tau > \tau_0$, all arms such that $\Delta_a \geq 4\beta_\tau$ will be eliminated. To show this, we have that under good event \mathcal{E} ,

$$\tilde{\mu}_a + \beta_\tau \leq \mu_a + 2\beta_\tau < \mu_{a^*} - 4\beta_\tau + 2\beta_\tau \leq \tilde{\mu}_{a^*} - \beta_\tau \leq \tilde{\mu}_{a_\tau} - \beta_\tau$$

which implies that arm a will be eliminated by the rule. Thus, for each sub-optimal arm a , let $\tau(a)$ be the last epoch that arm a is not eliminated. By the above result, we have

$$\Delta_a \leq 4\beta_{\tau(a)} = O\left(\sqrt{\frac{\log(|\mathcal{S}|\tau(a)^2/\delta)}{B_{\tau(a)}}} + \left(\frac{\log(|\mathcal{S}|\tau(a)^2/\delta)}{B_{\tau(a)}\epsilon}\right)^{1-\frac{1}{k}}\right).$$

Hence, we have

$$B_{\tau(a)} \leq O\left(\frac{\log(|\mathcal{S}|\tau(a)^2/\delta)}{\Delta_a^2} + \frac{\log(|\mathcal{S}|\tau(a)^2/\delta)}{\epsilon} \left(\frac{1}{\Delta_a}\right)^{\frac{k}{k-1}} + \frac{\log(D|\mathcal{S}|\tau_0^2/\delta)}{\epsilon}\right).$$

Since $|\mathcal{S}| \leq K$ and $2^\tau \leq T$ for any τ . Thus,

$$B_{\tau(a)} \leq O\left(\frac{\log(K \log^2 T/\delta)}{\Delta_a^2}, \frac{\log(K \log^2 T/\delta)}{\epsilon} \left(\frac{1}{\Delta_a}\right)^{\frac{k}{k-1}}, \frac{\log(DK \log^2 T/\delta)}{\epsilon}\right),$$

Since the batch size doubles, we have $N_a(T) \leq 2B_{\tau(a)}$ for each sub-optimal arm a . Therefore, for all arm $a \in [K]$,

$$\mathcal{R}_T = \sum_{a \in [K]} N_a(T) \Delta_a \leq 2B_{\tau(a)} \Delta_a.$$

Let η be a number in $(0, 1)$. For all arms $a \in [K]$ with $\Delta_a \leq \eta$, the regret incurred by pulling these arms is upper bounded by $T\eta$. For any arm $a \in [K]$ with $\Delta_a > \eta$, choose $\delta = \frac{1}{T}$ and assume $T \geq K$, then the expected regret incurred by pulling arm a is upper bounded by

$$\begin{aligned} \mathbb{E} \left[\sum_{a \in [K], \Delta_a > \eta} \Delta_a N_a(T) \right] &\leq \mathbb{P}(\bar{\mathcal{E}}) \cdot T + O\left(\sum_{a \in [K], \Delta_a > \eta} \left\{ \frac{\log T}{\Delta_a} + \frac{\log T}{\epsilon} \left(\frac{1}{\Delta_a}\right)^{\frac{1}{k-1}} + \frac{\log DT}{\epsilon} \Delta_a \right\}\right) \\ &\leq O\left(\frac{K \log T}{\eta} + \frac{K \log T}{\epsilon \eta^{\frac{1}{k-1}}} + \frac{KD \log(DT)}{\epsilon}\right) \end{aligned}$$

where the last term in the last inequality is based on following result: from the heavy-tailed assumption for rewards distributions in Definition 3.2, we have for any $a \in [K]$, $\mu_a \in [-D, D]$, so $\Delta_a = \mu^* - \mu_a \leq 2D$.

Thus the regret is at most

$$T\eta + O\left(\frac{K \log T}{\eta} + \frac{K \log T}{\epsilon \eta^{\frac{1}{k-1}}} + \frac{KD \log(DT)}{\epsilon}\right).$$

Taking $\eta = \max\left\{\sqrt{\frac{K \log T}{T}}, \left(\frac{K \log T}{T\epsilon}\right)^{\frac{k-1}{k}}\right\}$, the regret is at most

$$O\left(\sqrt{KT \log T} + \left(\frac{K \log T}{\epsilon}\right)^{\frac{k-1}{k}} T^{\frac{1}{k}} + \frac{DK \log(DT)}{\epsilon}\right).$$

For privacy guarantee, based on Laplacian mechanism in Definition 3.5, privacy guarantee for histogram learner in [49, Lemma 2.3], parallel composition theorem in Lemma B.2 and Post-processing in Lemma B.1, we can get the result. \square

Proof of Theorem 6.12. The proof of the theorem is similar to the proof of Theorem 6.2, now the requirement for batch size to start to arm elimination becomes $\max\left\{\frac{\epsilon \log(16/\delta)}{\epsilon}, \frac{c_1 \log(16D/\delta)}{\epsilon}, \frac{\log(16/\delta)}{\alpha_1^2}\right\}$ and the upper bound of Δ_a for each $a \in [K]$ is $2D$. Then we can get the result of upper bound for regret.

For privacy guarantee, based on Laplacian mechanism in Definition 3.5, privacy guarantee for histogram learner in [49, Lemma 2.3], parallel composition theorem in Lemma B.2 and Post-processing in Lemma B.1, we can get the result. \square