

DIFFERENTIALLY PRIVATE ALGORITHMS FOR EFFICIENT ONLINE MATROID OPTIMIZATION

Kushagra Chandak

Department of Computing Science, University of Alberta
Edmonton, Canada
kchandak@ualberta.ca

Bingshan Hu

Department of Computing Science, University of Alberta
Alberta Machine Intelligence Institute (Amii)
Edmonton, Canada
bingsha1@ualberta.ca

Nidhi Hegde

Department of Computing Science, University of Alberta
Canada CIFAR AI Chair, Alberta Machine Intelligence Institute (Amii)
Edmonton, Canada
nidhi.hegde@ualberta.ca

ABSTRACT

A matroid bandit is the online version of combinatorial optimization on a matroid, in which the learner chooses K actions from a set of L actions that can form a matroid basis. Many real-world applications such as recommendation systems can be modeled as matroid bandits. In such learning problems, the revealed data may involve sensitive user information. Therefore, privacy considerations are crucial. We propose two simple and practical differentially private algorithms for matroid bandits built on the well-known Upper Confidence Bound algorithms and Thompson Sampling. The key idea behind our first algorithm, Differentially Private Upper Confidence Bound for Matroid Bandits (DPUCB-MAT), is to construct differentially private upper confidence bounds. The second algorithm, Differentially Private Thompson Sampling for Matroid Bandits (DPTS-MAT), is based on the idea of drawing random samples from differentially private posterior distributions. Both algorithms achieve $O(L \ln(n)/\Delta + LK \ln(n) \min\{K, \ln(n)\}/\varepsilon)$ regret bounds, where Δ denotes the mean reward gap and ε is the required privacy parameter. Our derived regret bounds rely on novel technical arguments that deeply explore the special structure of matroids. We show a novel way to construct ordered pairs between the played actions and the optimal actions, which contributes to decomposing a matroid bandit problem into K stochastic multi-armed bandit problems. Finally, we conduct experiments to demonstrate the empirical performance of our proposed learning algorithms on both a synthetic dataset and a real-world movie-recommendation dataset.

1 INTRODUCTION

We study the learning problem of stochastic matroid bandits first proposed by [Kveton et al. \(2014\)](#) with differential privacy. In a stochastic matroid bandit problem, we have a matroid (E, \mathcal{I}) and a stochastic environment $\nu(E)$, where E is a set of L *base arms* and \mathcal{I} is a set of *independent sets*. Each base arm $e \in E$ is associated with a weight that is independently drawn from a fixed but unknown probability distribution P_e collected by $\nu(E) = (P_e : e \in E)$. As matroids generalize the notion of linear independence in vector spaces, we are interested in a *basis* of a matroid which is a maximal independent set $A \in \mathcal{I}$. All the bases of a matroid have the same size denoted by K . Each basis can be viewed as a *super arm* which consists of exactly K base arms. An important problem for matroid bandits is to learn the basis with the maximum total weight by interacting with environment $\nu(E)$ sequentially. In each round $t \in [n]$, the learner chooses a basis A^t and simultaneously, the environment generates a random weight vector $w_t := (w_t(e_1), w_t(e_2), \dots, w_t(e_L))$ for all the base arms $e_i \in E$. At the end of round t , the learner observes each individual weight $w_t(e)$ for all $e \in A^t$ and obtains as *return* the sum of the weights associated with the played basis. The goal of the learner is to choose bases sequentially to maximize the total return over a finite number of n rounds.

The application of recommending a set of diverse movies to users can be framed as a matroid bandit learning problem ([Kveton et al., 2014](#)), where the learning algorithm recommends a set of carefully chosen movies to the users. In this example, each movie can be characterized by a feature vector denoting the genres of that movie. In each round, the set of recommended movies can be viewed as a super arm. Since we would like the recommended movies to

represent diverse genres, the matroid structure ensures that the feature vectors of the recommended movies are linearly independent, which further indicates that the recommended movies do not include movies with similar genres. At the end of the round, based on the feedback from the users, the learning algorithm adjusts its future recommendations.

Since the learner only observes the weights associated with the played super arm, the central challenge that the learner faces is the *exploitation-versus-exploration* dilemma. In each round, the learner needs to decide whether to choose a super arm with the highest empirical return based on the past information (exploitation) or choose an under-explored super arm to gain information about the unknown environment (exploration). Upper Confidence Bound (UCB)-based (Auer et al., 2002; Garivier & Cappé, 2011) and Thompson Sampling-based learning algorithms (Kaufmann et al., 2012; Agrawal & Goyal, 2017) are both successful in balancing exploitation-versus-exploration. In bandit problems, the UCB-based algorithms are motivated by the principle of *optimism in the face of uncertainty* and can be justified by well-known concentration inequalities. In the UCB-based algorithms, each arm maintains a UCB index, which is an upper bound of the constructed confidence interval. Then, the learning algorithm makes a decision based on UCB indices for all the arms. Different from the UCB-based algorithms that rely on the UCB indices to tackle the exploitation-versus-exploration dilemma, Thompson Sampling-based algorithms maintain a posterior distribution for each arm. For each arm, a random sample is drawn from the posterior distribution. Then, the learning algorithm makes a decision based on the random samples for all the arms. If the learning algorithm knows the mean rewards of all the arms, pulling the arm with the highest mean reward is always the best strategy to achieve the highest expected reward. The performance of bandit algorithms is commonly measured by *Regret*, which represents the expected cumulative performance gap between always choosing the arm with the highest expected reward (fixed but unknown) and the learning algorithm’s actual pulled arms.

The example of recommending diverse movies highlights the necessity for preserving the privacy of users. The feedback information from users (e.g. movie ratings) also reveals their watch history or preferences towards the recommended movies. Being users, we may wish to keep this information private. Take the Netflix Prize dataset for example. As shown in Narayanan & Shmatikov (2008), just anonymizing the ratings is not enough for preserving privacy. In this paper, we focus on the learning problem of matroid bandits with differential privacy. Differential privacy is the most commonly-used notion of privacy for machine learning algorithms (Dwork et al., 2014). If a learning algorithm is implemented in a differentially private manner, the information associated with an individual has almost no impact on the output of the learning. In other words, private learning algorithms are not sensitive to information from a single individual. In this work, we focus on ϵ -differential privacy, where ϵ is the privacy parameter. The parameter ϵ can be viewed as the privacy budget that can be distributed among different components of the learning algorithm.

Now, we summarize the key contributions of this work.

1. In this work, we propose two *sample efficient* and *computationally fast* differentially private algorithms for matroid bandits. Our first algorithm, DPUCB-MAT, is built upon the well-established UCB1 policy of Auer et al. (2002). The regret bound of DPUCB-MAT is $O(L \ln(n)/\Delta + \min\{K, \log(n)\} LK \ln(n)/\epsilon)$, where Δ is the mean reward gap and ϵ is the required privacy parameter. Our second algorithm, DPST-MAT, is built on Thompson Sampling which has demonstrated competitive practical performances (Chapelle & Li, 2011). The regret bound of DPST-MAT is $O(L \ln(n)/\Delta + \min\{K, \log(n)\} LK \ln(n)/\epsilon) + O(K \ln(n)/\Delta)$.
2. Regarding the regret analysis, we propose a unified approach to decompose the regret of DPUCB-MAT and DPST-MAT. Our novel regret decomposition relies on the introduction of a round-dependent permutation on the order of the base arms in the optimal super arm. The permutation contributes to decomposing the regret of matroid bandits into K different stochastic bandit problems.
3. We conduct experiments to evaluate the empirical performance of our proposed algorithms by using both synthetic and real-world movie-recommendation datasets. The experimental results demonstrate that our proposed differentially private algorithms are efficient.

2 DIFFERENTIALLY PRIVATE MATROID BANDITS

In this section, we present our learning problem formally. We first introduce stochastic matroid bandits. Then, we introduce the definition of differential privacy for matroid bandit learning algorithms.

2.1 MATROID BANDITS

A stochastic matroid bandit problem can be specified by $((E, \mathcal{I}); \nu(E))$, where (E, \mathcal{I}) defines a matroid and $\nu(E)$ defines an environment that a learner interacts with. In a matroid bandit problem, E is a set of L items, also called the *set of base arms*, and \mathcal{I} is a family of subsets of E , also called the family of *independent sets*, defined by the

The privacy parameter ε can be thought of as a privacy budget. It can be allocated among different components within the algorithm. An important property of DP is that it is immune to post-processing, which states that the privacy guarantee cannot get worse if we release a function of the output instead of the output itself (Smith & Ullman, 2021).

Fact 1. (Post-processing, (Dwork et al., 2014)). Let $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$ be arbitrary sets and n be a positive integer. Let $\mathcal{B} : \mathcal{X}^n \rightarrow \mathcal{Y}$ and $\mathcal{B}' : \mathcal{Y} \rightarrow \mathcal{Z}$ be randomized algorithms. If \mathcal{B} is ε -DP, then the composition $\mathcal{B}'(\mathcal{B}(\cdot))$ is also ε -DP.

3 RELATED WORKS

Kveton et al. (2014) initiated the study of learning the maximum weight basis for matroid bandits. They proposed a UCB-based algorithm, Optimistic Matroid Maximization (OMM), that achieves the optimal $O(L \ln(n)/\Delta)$ regret bound. The key idea in OMM is to construct upper confidence bounds based on UCB1 of Auer et al. (2002) for all base arms. As OMM is built upon UCB1, OMM cannot be asymptotically optimal. Later, Talebi & Proutiere (2016) proposed another UCB-based algorithm, Efficient Sampling for Matroids (KL-OSM), that achieves the asymptotically optimal $(1 + \epsilon)L \ln(n)\Delta/d_{\text{KL}}(\mu_* - \Delta, \mu_*) + o(\ln(n))$ regret bound, where $\epsilon > 0$ can be any value, μ_* denotes the mean reward of an optimal base arm and $d_{\text{KL}}(a, b)$ denotes the KL-divergence between two Bernoulli distributions with parameters $a, b \in (0, 1)$. The key idea to achieve asymptotic optimality is to construct upper confidence bounds based on KL-UCB of Garivier & Cappé (2011). Other than the aforementioned UCB-based algorithms, Wang & Chen (2018) proposed Combinatorial Thompson Sampling (CTS), a Thompson Sampling-based algorithm with Beta priors, for combinatorial bandits. Combinatorial bandits generalize the setting of matroid bandits with the following key features. In combinatorial bandits, the set of super arms does not have any special structure that we can utilize. Also, the size of a super arm is at most K instead of exactly K . Since matroid bandits are special cases of combinatorial bandits, by a refined regret analysis, CTS achieves an $O(L \ln(n)/\Delta) + \tilde{O}(L/\Delta^4)$ regret for matroid bandits.

Mishra & Thakurta (2015) initiated the study of stochastic multi-armed bandits with differential privacy and proposed both the UCB-based and Thompson Sampling-based learning algorithms. Their proposed algorithms rely on the post-processing property of differential privacy (Fact 1). More specifically, they first guarantee that the internal learning algorithm computing the empirical means is ε -differentially private. Then, from post-processing, they immediately conclude that the proposed bandit algorithms are ε -DP. However, their proposed learning algorithms are very sub-optimal due to the usage of the Tree-based Mechanism (Chan et al., 2011; Dwork et al., 2010) to inject noise to preserve privacy. Later, Chen et al. (2020) introduced differential privacy in combinatorial bandits and proposed a UCB-based algorithm, Differentially Private Combinatorial UCB (CUCB-DP), for differentially private combinatorial bandits. CUCB-DP achieves an $O(LK \ln^2(n)/\Delta) + \tilde{O}(LK \ln^3(n)/\varepsilon)$ regret upper bound and an $\Omega(LK \ln(n)/\Delta + LK \ln(n)/\varepsilon)$ regret lower bound.²

Recently, Hu et al. (2021); Azize & Basu (2022) devised optimal UCB-based algorithms and Hu & Hegde (2022) devised an optimal Thompson Sampling-based algorithm with Beta priors for differentially private stochastic bandits. These algorithms all achieve the optimal $O(L \ln(n)/\Delta + L \ln(n)/\varepsilon)$ regret bound. The key ideas to achieve optimality are the usage of “laziness” and “forgetfulness” along with the Laplace Mechanism (Dwork et al., 2014) to inject noise to mask the true empirical means instead of using the Tree-based Mechanism. The idea of laziness is to update the differentially private empirical mean of an arm in a delayed manner. We only update the DP empirical mean of an arm when a certain number of observations are available from that arm. The idea of forgetfulness is to use fresh observations to update the DP empirical mean. Once observations have been used, we abandon them. Since the change of one reward vector only impacts the aggregated reward by one at most, from Laplace Mechanism, we can add a noise drawn from $\text{Lap}(1/\varepsilon)$ to the aggregated reward of each arm when updating the DP empirical mean.³

Our proposed UCB-based algorithm DPUCB-MAT (Algorithm 2) can be viewed as a differentially private version of OMM. When setting $\varepsilon \rightarrow \infty$, i.e., in the non-private matroid bandit setting, the regret bound of DPUCB-MAT (Theorem 2) recovers the regret bound of OMM. When setting a suitable privacy parameter ε , our regret bound removes an extra $\log(n)$ factor as compared to the regret bound shown in Chen et al. (2020). Although our proposed Thompson sampling-based algorithm DPTS-MAT (Algorithm 3) could be seen as a differentially private version of CTS for matroid bandits, we use different proof techniques. Section 4.3 presents more detail. Table 1 summarizes the regret bounds for matroid bandits in both non-private and private settings.

²The $\tilde{O}(\cdot)$ notation hides an extra $\log \log T$ factor.

³The probability density function of a Laplace distribution $\text{Lap}(b)$ centered at 0 with b as the scale parameter is $h_b(y) = (1/2b)e^{-|y|/b}$ with $b > 0$.

Table 1: Regret upper bounds for UCB and TS-based algorithms for matroid bandits.

Algorithms	Non-private results	Private results (ours)
OMM (UCB1-based)	$O\left(\frac{L \ln(n)}{\Delta}\right)$ (Kveton et al., 2014)	$O\left(\frac{L \ln(n)}{\Delta} + \frac{\min\{K, \log(n)\} L K \ln(n)}{\varepsilon}\right)$
CTS (TS-based)	$\tilde{O}\left(\frac{L \ln(n)}{\Delta} + \frac{L}{\Delta^4}\right)$ (Wang & Chen, 2018)	$O\left(\frac{L \ln(n)}{\Delta} + \frac{\min\{K, \log(n)\} L K \ln(n)}{\varepsilon}\right)$

4 ALGORITHMS AND ANALYSIS

In this section, we present our proposed learning algorithms and their theoretical guarantees. First, we introduce some notation used by Algorithm 2 and Algorithm 3. We denote by $T_e(t-1)$ the “effective” number of observations used to compute the empirical mean of a base arm $e \in E$ by the end of round $t-1$ and denote by $\hat{w}_{e, T_e(t-1)}(t-1)$ the empirical mean of a base arm e among these $T_e(t-1)$ observations. To have a differentially private version of $\hat{w}_{e, T_e(t-1)}(t-1)$ with a small amount of noise, we still use the ideas of “laziness” and “forgetfulness” introduced in Hu et al. (2021); Azize & Basu (2022); Hu & Hegde (2022). As the total privacy budget is ε and at the end of each round t there are exactly K observations revealed, from the basic composition property (Dwork et al., 2014), we know that each base arm has an ε/K privacy budget. Let $\varepsilon_0 := \varepsilon/K$. Let $\tilde{w}_{e, T_e(t-1)}(t-1) = \hat{w}_{e, T_e(t-1)}(t-1) + \frac{\text{Lap}(1/\varepsilon_0)}{T_e(t-1)}$ denote the DP mean estimate of a base arm e . For each base arm e , we store the reward observations in the array \mathcal{T}_e , and $\Sigma \mathcal{T}_e$ denotes the sum of all the observations in \mathcal{T}_e .

4.1 DIFFERENTIALLY PRIVATE UCB FOR MATROID BANDITS

The UCB-based differentially private algorithm for matroid bandits DPUCB-MAT is shown in Algorithm 2. The general idea is to construct *differentially private upper confidence bound* $U_t(e)$ (Line 5) for each base arm $e \in E$. We construct $U_t(e)$ as

$$U_t(e) := \tilde{w}_{e, T_e(t-1)}(t-1) + \sqrt{\frac{3 \ln(Kt)}{T_e(t-1)}} + \frac{3 \ln(Kt)}{\varepsilon_0 \cdot T_e(t-1)}. \quad (3)$$

With all these private upper confidence bounds $U_t(e)$ in hand, DPUCB-MAT selects the best super arm A^t in a greedy way, i.e., invoking Algorithm 1 with all $U_t(e)$ as input and A^t as output. In other words, DPUCB-MAT plays $A^t = \arg \max_{A \in \mathcal{I}} \sum_{e \in A} U_t(e)$. Then, the rewards $w_t(e)$ for $e \in A^t$ are revealed. For each individual $w_t(e)$, we add it to the corresponding \mathcal{T}_e (Line 9). If the number of observations of any base arm e hits 2^{s_e+1} , then, it is the right time to update the DP mean estimate $\tilde{w}_{e, T_e(t-1)}(t-1)$ using the 2^{s_e+1} observations stored in \mathcal{T}_e (Line 12). Since now all the observations in \mathcal{T}_e have been processed, we increment the counter s_e by 1 and reset \mathcal{T}_e (Line 13).

Algorithm 2 DPUCB-MAT

- 1: **Input:** Matroid (E, \mathcal{I}) and privacy budget ε
 - 2: Observe $w_0(e) \sim P_e$, set $T_e \leftarrow 1$, $s_e \leftarrow 0$, $\tilde{w}_{e, T_e} \leftarrow w_0(e) + \text{Lap}\left(\frac{1}{\varepsilon_0}\right)$, $\mathcal{T}_e \leftarrow \{\}$, where $\varepsilon_0 = \frac{\varepsilon}{K} \triangleright$ Initialization
 - 3: **for** $t = 1, 2, \dots$ **do**
 - 4: **for** $e \in E$ **do**
 - 5: $U_t(e) := \tilde{w}_{e, T_e} + \sqrt{\frac{3 \ln(Kt)}{T_e}} + \frac{3 \ln(Kt)}{\varepsilon_0 \cdot T_e}$
 - 6: **end for**
 - 7: Invoke Algorithm 1 with $U_t(e)$ for all $e \in E$ as input and A^t as output
 - 8: Play super arm A^t
 - 9: Observe $w_t(e) \sim P_e$ and add $w_t(e)$ to \mathcal{T}_e , $\forall e \in A^t$
 - 10: Find $B^t = \{e \in A^t : |\mathcal{T}_e| = 2^{s_e+1}\} \triangleright$ Find all the base arms with the number of observations hitting 2^{s_e+1}
 - 11: **for** $e \in B^t$ **do**
 - 12: $T_e \leftarrow 2^{s_e+1}$, $\tilde{w}_{e, T_e} = \frac{\Sigma \mathcal{T}_e + \text{Lap}(1/\varepsilon_0)}{T_e} \triangleright \Sigma \mathcal{T}_e$ denotes the sum of all the observations in \mathcal{T}_e
 - 13: $s_e \leftarrow s_e + 1$, $\mathcal{T}_e \leftarrow \{\}$ \triangleright Doubling the “effective” number of observations and reset \mathcal{T}_e
 - 14: **end for**
 - 15: **end for**
-

We now present theoretical guarantees for Algorithm 2.

Theorem 1. Algorithm 2 is ε -differentially private.

Proof. Since the learner makes decisions based on the differentially private upper confidence bounds and constructing differentially private upper confidence bounds (Lines 4 to 6) can be thought of as post-processing of the DP mean estimates, it suffices to show that as long as the empirical mean computation (lines 9 to 14) is ε -DP, then the matroid bandit learning algorithm itself is ε -DP. Now, we consider two neighboring reward vector sequences \mathbf{W} and \mathbf{W}' which differ in a reward vector in some round s , i.e., $w_s \neq w'_s$. Since the learner only observes the rewards associated with the base arms in the played super arm, changing from w_s to w'_s may only affect the DP mean estimates for all the base arms in A^s . In other words, the DP mean estimates of at most K base arms may be impacted. Since the learning algorithm uses fresh observations to update the DP mean estimates and for each base arm $e \in A^s$ the aggregated reward $\Sigma \mathcal{T}_e$ changes by at most one when working over \mathbf{W} and \mathbf{W}' , from Laplace Mechanism, we know that adding a noise random variable sampled from $\text{Lap}(1/\varepsilon_0)$ is enough to make the function computing the mean estimate of an individual base arm ε_0 -DP. Since at most K base arms may be impacted, using the basic composition property of DP, we know that the learning algorithm to compute the DP mean estimates for all the base arms is ε -DP. \square

Theorem 2. The regret of Algorithm 2 is

$$R_n \leq \sum_{e \in \bar{A}^*: \Delta_{e, \min} > 0} O\left(\frac{\ln(Kn)}{\Delta_{e, \min}} + \frac{\min\{K, \log(Kn)\} \cdot \ln(Kn)}{\varepsilon/K}\right), \quad (4)$$

where $\Delta_{e, \min} = \min_{k \in [K]: \Delta_{e, k} > 0} \Delta_{e, k}$.

Discussion Algorithm 2 can be viewed as a differentially private version of OMM of Kveton et al. (2014). When setting $\varepsilon \rightarrow \infty$, the differentially private matroid bandit problem boils down to the non-private matroid bandits. In this special setting, the regret bound shown in Theorem 2 matches both the regret upper bound and regret lower bound presented in Kveton et al. (2014). With a suitable privacy parameter ε , our regret bound improves the state-of-the-art regret bound presented in Chen et al. (2020) by removing an extra $\log(n)$ factor.⁴

Algorithm 3 DPTS-MAT

- 1: **Input:** Matroid $M = (E, \mathcal{I})$ and privacy parameter ε
 - 2: Observe $w_0(e) \sim P_e$, set $T_e \leftarrow 1, s_e \leftarrow 0, \tilde{w}_e \leftarrow w_0(e) + \text{Lap}\left(\frac{1}{\varepsilon_0}\right), \mathcal{T}_e \leftarrow \{\}$, where $\varepsilon_0 = \frac{\varepsilon}{K}$ \triangleright Initialization
 - 3: **for** $t = 1, 2, \dots$ **do**
 - 4: **for** $e \in E$ **do**
 - 5: Set $w'_{e, T_e}(t) := \tilde{w}_{e, T_e} + \frac{3 \ln(Kt)}{\varepsilon_0 \cdot T_e}$ \triangleright Boost the parameters of the posterior distributions
 - 6: Sample $\theta_e(t) \sim \mathcal{N}\left(w'_{e, T_e}(t), \frac{1}{T_e}\right)$ \triangleright Draw random samples from a Gaussian distribution
 - 7: **end for**
 - 8: Invoke Algorithm 1 with $\theta_e(t)$ for all $e \in E$ as input and A^t as output
 - 9: Play super arm A^t
 - 10: Observe $w_t(e) \sim P_e$ and add $w_t(e)$ to \mathcal{T}_e , $\forall e \in A^t$
 - 11: Find $B^t = \{e \in A^t : |\mathcal{T}_e| = 2^{s_e+1}\}$ \triangleright Find all the base arms with the number of observations hitting 2^{s_e+1}
 - 12: **for** $e \in B^t$ **do**
 - 13: $T_e \leftarrow 2^{s_e+1}, \tilde{w}_{e, T_e} = \frac{\Sigma \mathcal{T}_e + \text{Lap}(1/\varepsilon_0)}{T_e}$ $\triangleright \Sigma \mathcal{T}_e$ denotes the sum of all the observations in \mathcal{T}_e
 - 14: $s_e \leftarrow s_e + 1, \mathcal{T}_e \leftarrow \{\}$ \triangleright Doubling the “effective” number of observations and reset \mathcal{T}_e
 - 15: **end for**
 - 16: **end for**
-

4.2 DIFFERENTIALLY PRIVATE THOMPSON SAMPLING FOR MATROID BANDITS

Different from the UCB-based algorithms where the exploration-exploitation trade-off is done by constructing confidence intervals centered at the mean rewards, Thompson sampling-based algorithms maintain posterior distributions that model the mean rewards of each base arm. As our goal is to design learning algorithms that have good regret guarantees with a finite time horizon, we can use Gaussian priors.⁵

Our Thompson Sampling-based differentially private algorithm for matroid bandits DPTS-MAT is shown in Algorithm 3. The general idea is to boost the parameter of the posterior distribution from $\tilde{w}_{e, T_e(t-1)}(t-1)$ to $w'_{e, T_e(t-1)}(t)$, where

⁴For the discussion, we already translate their regret bound (Theorem 8) from combinatorial bandits to matroid bandits.

⁵As proved in Agrawal & Goyal (2017), Thompson Sampling with Beta priors can be asymptotically optimal while Thompson Sampling with Gaussian priors may not be asymptotically optimal.

$w'_{e, T_e(t-1)}(t) = \tilde{w}_{e, T_e(t-1)}(t-1) + 3 \ln(Kt) / (\varepsilon_0 \cdot T_e(t-1))$ for each base arm $e \in E$. Then, DPTS-MAT draws a random sample $\theta_e(t) \sim \mathcal{N}(w'_{e, T_e(t-1)}(t), 1/T_e(t-1))$ for all $e \in E$. With all these differentially private posterior samples $\theta_e(t)$ in hand, DPTS-MAT selects the best super arm A^t in a greedy way, i.e., invoking Algorithm 1 with all $\theta_e(t)$ as input and A^t as output. That is also to say, DPTS-MAT plays $A^t = \arg \max_{A \in \mathcal{I}} \sum_{e \in A} \theta_e(t)$. DPTS-MAT uses the same way as Algorithm 2 to process the revealed observations, i.e., only update the DP mean estimate of a base arm $e \in A^t$ if the number of observations in \mathcal{T}_e hits 2^{s_e+1} .

We now present theoretical guarantees for Algorithm 3.

Theorem 3. Algorithm 3 is ε -differentially private.

Proof. The proof is similar to the DP proof for Algorithm 2. It suffices to show that as long as the algorithm to compute the empirical mean of each base arm is ε_0 -DP, then the learning algorithm is ε -DP. Note that lines 4 to 9 can be thought of as post-processing. Since Algorithm 3 and Algorithm 2 use the same way to process the obtained observations, by using the same arguments as in the proof of Theorem 1, we can conclude the proof. \square

Theorem 4. The regret of Algorithm 3 is

$$R_n \leq \sum_{e \in \bar{A}^*: \Delta_{e, \min} > 0} O\left(\frac{\ln(Kn)}{\Delta_{e, \min}} + \frac{\min\{K, \log(Kn)\} \cdot \ln(Kn)}{\varepsilon/K}\right) + \sum_{k \in [K]: \Delta_{\min, k} > 0} O\left(\frac{\ln(Kn)}{\Delta_{\min, k}}\right), \quad (5)$$

where $\Delta_{e, \min} = \min_{k \in [K]: \Delta_{e, k} > 0} \Delta_{e, k}$ and $\Delta_{\min, k} = \min_{e \in \bar{A}^*: \Delta_{e, k} > 0} \Delta_{e, k}$.

Discussion Different from Theorem 2 where the regret bound only has one term which is linear in the size of the sub-optimal base arms, there are two terms in Theorem 4. The first term is the same as the regret bound shown in Theorem 2 and it captures the regret for introducing differential privacy. This term characterizes the regret in all the rounds when the posterior distributions of the sub-optimal base arms are not concentrated. It is not surprising why we have the second non-private term which is linear in the size of the optimal base arm set. The second term upper bounds the regret among all the rounds when the posterior distributions of the optimal base arms are not concentrated. As will be shown in Section 4.3, the core of our regret decomposition is to decompose a matroid bandit problem into K stochastic bandit problems. For each individual bandit problem, by slightly modifying the regret analysis in Agrawal & Goyal (2017), when using Gaussian priors, an additive $O(\ln(n)/\Delta_{\min})$ regret occurs among all the rounds when the Gaussian posterior distributions of the optimal arm are not concentrated. In contrast, when using Beta priors, the additive term can be $\tilde{O}(1/\Delta_{\min}^4)$, where \tilde{O} hides problem-dependent constants. Since the regret for a matroid bandit is composed of K different stochastic bandit problems, we have the second term in Theorem 4.

4.3 REGRET DECOMPOSITION

In this section, we present a unified approach to decompose the regret of Algorithm 2 and Algorithm 3. The core is to introduce a round-dependent permutation π_t over the ordered set A^* . Then, we explore the special structures that matroids have, which are all the bases of a matroid have the same size and matroids have certain properties described in Section 2.1. After introducing π_t , the regret for a matroid bandit problem can be decomposed to upper bound the total regret of K different stochastic bandit problems. The construction of π_t is inspired by Lemma 1 in Kveton et al. (2014). Recall $A^t = \{a_1^t, a_2^t, \dots, a_K^t\}$ is a descending ordered set based on the differentially private upper confidence bounds in Algorithm 2 (or the differentially private posterior samples in Algorithm 3) and $A^* = \{a_1^*, a_2^*, \dots, a_K^*\}$ is a descending ordered set based on the mean rewards. The purpose of introducing permutation $\pi_t : \{1, \dots, K\} \rightarrow \{1, \dots, K\}$ over the ordered set A^* is to construct K ordered pairs between A^t and A^* with certain properties. We construct π_t in a backward order as follows.

Fix $B_K = \{a_1^t, \dots, a_{K-1}^t\}$. If $a_K^t \in A^*$, i.e., $a_K^t = a_i^*$ for some $i \in [K]$, we set $\pi_t(K) = i$, i.e., we pair a_K^t to itself. If $a_K^t \notin A^*$, due to the augmentation property of matroids (Property (3) in Section 2.1), we can set $\pi_t(K) = \min\{i : a_i^* \in A^* \setminus B_K, B_K \cup \{a_i^*\} \in \mathcal{I}\}$, i.e., we can pair a_K^t with the optimal base arm with the smallest index that can be added to B_K to form a matroid basis. Now, fix $B_{K-1} = \{a_1^t, \dots, a_{K-2}^t, a_{\pi_t(K)}^t\}$. If $a_{K-1}^t = a_i^*$ for some $i \in [K]$, we set $\pi_t(K-1) = i$. If $a_{K-1}^t \notin A^*$, we can set $\pi_t(K-1) = \min\{i : a_i^* \in A^* \setminus B_{K-1}, B_{K-1} \cup \{a_i^*\} \in \mathcal{I}\}$. The same idea is applied to all the remaining base arms a_{K-2}^t, \dots, a_1^t in A^t .

After applying π_t , all the optimal base arms in A^* will be ordered as $A_{\pi_t}^* = \{a_{\pi_t(1)}^*, \dots, a_{\pi_t(k)}^*, \dots, a_{\pi_t(K)}^*\}$. Still, the order of all the base arms in A^t is $\{a_1^t, \dots, a_k^t, \dots, a_K^t\}$. Now, we can construct the following ordered pairs

$(a_1^t, a_{\pi_t(1)}^*), \dots, (a_k^t, a_{\pi_t(k)}^*), \dots, (a_K^t, a_{\pi_t(K)}^*)$. It is not hard to verify that each $(a_k^t, a_{\pi_t(k)}^*)$ has the following two properties. First, if $a_k^t \in A^*$, its pair is itself, i.e., $\Delta_{a_k^t, \pi_t(k)} = 0$ (recall $\Delta_{e,k} = \bar{w}(a_k^*) - \bar{w}(e)$ for any two base arms a_k^* and e). If $a_k^t \notin A^*$, the differentially private upper confidence bound in Algorithm 2 (or the differentially private posterior sample in Algorithm 3) of a_k^t is no smaller than that of $a_{\pi_t(k)}^*$. Second, given $\{a_1^t, a_2^t, \dots, a_{k-1}^t\}$, both a_k^t and its pair $a_{\pi_t(k)}^*$ can be added to the solution set without breaking the independence of the solution set. In other words, both $\{a_1^t, \dots, a_{k-1}^t\} \cup \{a_k^t\} \in \mathcal{I}$ and $\{a_1^t, \dots, a_{k-1}^t\} \cup \{a_{\pi_t(k)}^*\} \in \mathcal{I}$ hold.

To permute $A_{\pi_t}^*$ back to the original order $\{a_1^*, \dots, a_K^*\}$, we define $\pi_t^{-1} : \{1, \dots, K\} \rightarrow \{1, \dots, K\}$, the inverse permutation of π_t . That is, for each $k \in [K]$, we have $\pi_t^{-1}(\pi_t(k)) = k$. We apply π_t^{-1} to permute the ordered sets $A_{\pi_t}^*$ and A^t separately. Applying π_t^{-1} over $A_{\pi_t}^*$ gives us the original order, which is $\{a_1^*, \dots, a_K^*\}$. Applying π_t^{-1} over A^t gives $\{a_{\pi_t^{-1}(1)}^t, \dots, a_{\pi_t^{-1}(k)}^t, \dots, a_{\pi_t^{-1}(K)}^t\}$. Based on the new order, we construct the following ordered pairs

$$(a_{\pi_t^{-1}(1)}^t, a_1^*), \dots, (a_{\pi_t^{-1}(k)}^t, a_k^*), \dots, (a_{\pi_t^{-1}(K)}^t, a_K^*). \quad (6)$$

Since a single π_t^{-1} is used to permute both the ordered sets, π_t^{-1} can be viewed as a permutation that permutes

$$(a_1^t, a_{\pi_t(1)}^*), \dots, (a_k^t, a_{\pi_t(k)}^*), \dots, (a_K^t, a_{\pi_t(K)}^*) \text{ to } (a_{\pi_t^{-1}(1)}^t, a_1^*), \dots, (a_{\pi_t^{-1}(k)}^t, a_k^*), \dots, (a_{\pi_t^{-1}(K)}^t, a_K^*).$$

Now, we are ready to decompose the regret. Define $\{a \leftrightarrow a^*\}$ as an event that base arms a and a^* form a pair. We have

$$\begin{aligned} R_n &= \sum_{t=1}^n \mathbb{E} \left[\sum_{e \in A^*} \bar{w}(e) - \sum_{e \in A^t} \bar{w}(e) \right] \\ &\stackrel{(a)}{=} \sum_{t=1}^n \mathbb{E} \left[\sum_{k=1}^K (\bar{w}(a_k^*) - \bar{w}(a_{\pi_t^{-1}(k)}^t)) \cdot \mathbb{1} \{a_{\pi_t^{-1}(k)}^t \leftrightarrow a_k^*\} \right] \\ &\leq \sum_{k=1}^K \sum_{t=1}^n \mathbb{E} \left[\Delta_{a_{\pi_t^{-1}(k)}^t, k} \cdot \mathbb{1} \{ \Delta_{a_{\pi_t^{-1}(k)}^t, k} > 0 \} \cdot \mathbb{1} \{a_{\pi_t^{-1}(k)}^t \leftrightarrow a_k^*\} \right] \\ &\leq \sum_{k=1}^K \sum_{e \in \bar{A}^*} \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \{a_{\pi_t^{-1}(k)}^t = e\} \cdot \mathbb{1} \{ \Delta_{e,k} > 0 \} \cdot \mathbb{1} \{e \leftrightarrow a_k^*\} \cdot \Delta_{e,k} \right] \\ &= \underbrace{\sum_{k=1}^K \sum_{e \in \bar{A}^* : \Delta_{e,k} > 0} \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \{a_{\pi_t^{-1}(k)}^t = e\} \right]}_{I_k} \cdot \mathbb{1} \{e \leftrightarrow a_k^*\} \cdot \Delta_{e,k} . \end{aligned} \quad (7)$$

Equality (a) uses the ordered pairs shown in (6). Note that I_k can be viewed as the regret of a stochastic bandit problem with a_k^* as the optimal arm and $\{e \in \bar{A}^* : \Delta_{e,k} > 0\}$ as the set of sub-optimal arms. The regret bounds for both DPUCB-MAT (Theorem 2) and DPIS-MAT (Theorem 4) can be derived based on the regret decomposition shown in (7). We defer the details of the proof to the appendix.

5 EXPERIMENTS

We perform experiments in two different settings. In Section 5.1, we evaluate our algorithms on a synthetic dataset and in Section 5.2, we use a real-world movie-rating dataset with the purpose of recommending diverse and popular movies to users in a differentially private manner. To measure the performance of our algorithms we use the *expected per-round return* as suggested in Kveton et al. (2014). The expected per-round return in round s is computed as $\frac{1}{s} \sum_{t=1}^s \sum_{e \in A^t} \bar{w}(e)$. For DPUCB-MAT, we compare its empirical performance with two baselines. The first baseline is the optimal return $f(A^*, \bar{w}) = \sum_{e \in A^*} \bar{w}(e)$ denoted as Optimal Policy in the plots. The second baseline is OMM in Kveton et al. (2014), which is the non-private UCB1-based algorithm for matroid bandits. Similarly, for DPIS-MAT, we compare its empirical performance with the optimal return and the non-private CTS in Wang & Chen (2018). To have a fair empirical performance comparison, we have already adapted CTS to the matroid bandit setting. We also show the performance of our algorithms with different values of $\epsilon \in \{10^5, 2, 10^{-4}\}$.

Table 2: Synthetic dataset.

Base arm e	Mean reward $\bar{w}(e)$
(1, 0, 0)	0.80
(0, 1, 0)	0.75
(0, 0, 1)	0.60
(1, 0, 1)	0.20
(0, 1, 1)	0.30
(2, 0, 0)	0.40
(0, 0, 0)	0.70

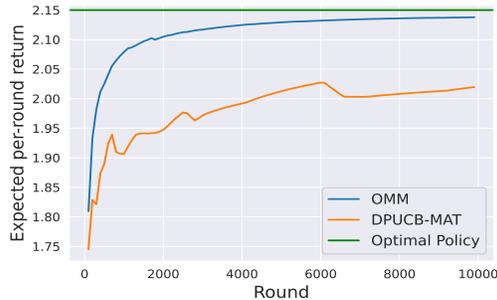
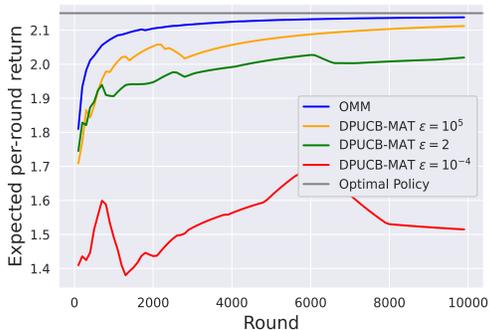
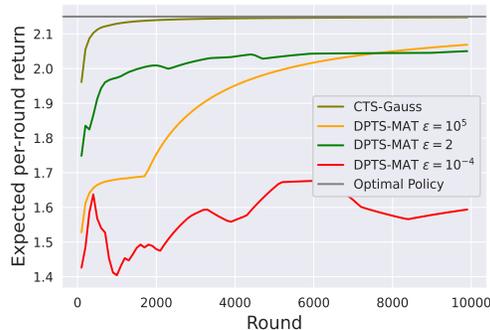


Figure 1: DPUCB-MAT on synthetic dataset.



(a) DPUCB-MAT



(b) DPTS-MAT

Figure 2: **Synthetic dataset.** Figure 2a shows the performance of DPUCB-MAT (Algorithm 2) for different values of ϵ on the synthetic dataset. We observe that the performance of DPUCB-MAT decreases as the value of ϵ decreases. We also observe that as ϵ increases, we can be in the non-private regime and the performance of DPUCB-MAT does not deteriorate by much. Figure 2b shows similar trends for DPTS-MAT (Algorithm 3).

5.1 SYNTHETIC DATASET

In this section, we report the experimental results of our proposed algorithms on a set of 3-dimensional vectors taken from Neel & Neudauer (2009) with $\epsilon = 2$. The base arm set is $E = \{e_1, \dots, e_7\}$, where each base arm is a 3-dimensional vector in the Euclidean space. Table 2 shows the mean rewards of all the base arms. As matroid independence is defined by the linear independence of the vectors, we have $|A^*| = |A^t| = 3$ and $A^* = \{e_1, e_2, e_3\}$. Since the total privacy budget is ϵ , the privacy budget for each $e \in A^t$ is $\epsilon_0 = \epsilon/3 = 2/3$. The reward $w_t(e)$ for each $e \in E$ is generated from a Bernoulli distribution with mean $\bar{w}(e)$. The total number of rounds is $n = 10,000$. The experimental results for DPUCB-MAT, OMM, and Optimal Policy are presented in Figure 1. From the results, we can see that DPUCB-MAT and OMM have a similar growth rate in terms of the return and approach the optimal return. The results for DPTS-MAT, CTS, and Optimal Policy also show similar trends. We defer this set of plots to the appendix.

Figure 2 shows the expected per-round return for different values of the privacy parameter $\epsilon \in \{10^5, 2, 10^{-4}\}$. In Figure 2a, we show the performance of DPUCB-MAT (Algorithm 2). We observe that when ϵ decreases, the performance of DPUCB-MAT deteriorates, and when ϵ increases, the performance of DPUCB-MAT becomes better. This is expected as a good differentially private learning algorithm should balance the privacy and regret guarantees. We also observe that when the privacy parameter ϵ is large, i.e., in the non-private regime, the performance of DPUCB-MAT approaches that of the non-private OMM. This is also expected as in the non-private regime, we do not pay any price for preserving privacy. Similar trends can also be seen for DPTS-MAT (Algorithm 3) in Figure 2b.

Table 3: Movies recommended by DPUCB-MAT that overlap with movies in A^* after 20k rounds.

Movie e	$\bar{w}(e)$	Genres
American Beauty	0.568	Comedy, Drama
Star Wars: Episode IV	0.496	Action, Adventure, Fantasy, Sci-Fi
Star Wars: Episode VI	0.478	Action, Adventure, Romance, Sci-Fi, War
Saving Private Ryan	0.440	Action, Drama, War
Men in Black	0.420	Action, Adventure, Comedy, Sci-Fi
L.A. Confidential	0.379	Crime, Film-Noir, Mystery, Thriller
Ghostbusters	0.361	Comedy, Horror
The Wizard of Oz	0.285	Animation, Children’s, Comedy, Musical

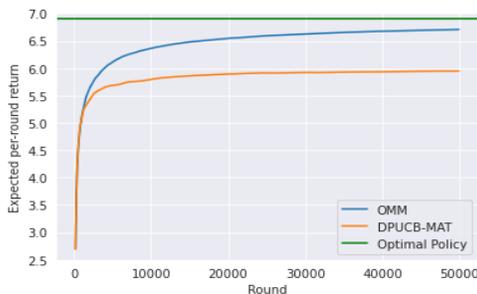


Figure 3: DPUCB-MAT on movie recommendation.

5.2 MOVIE RECOMMENDATION DATASET

In this experiment, we learn to recommend a set of *diverse* and *popular* movies from the MovieLens dataset (Harper & Konstan, 2015) with differential privacy. The experiment design is adopted from Kveton et al. (2014). Still, we report experimental results with the total privacy budget $\varepsilon = 2$. The total number of rounds is $n = 20,000$.

The entire dataset contains 1 million ratings from 6040 users. The total number of movies is 3883 from 18 different genres. To recommend popular movies, we select 100 movies that have received the most ratings. These 100 movies constitute the base arm set E of a matroid. To construct the independent sets \mathcal{I} of the matroid, we select a one-hot feature vector u_e for each movie $e \in E$, which denotes the genres of that movie. If the feature vectors u_e for all $e \in A$ are linearly independent, all the movies collected in A form an independent set, which further indicates that these movies are *diverse*. The expected reward $\bar{w}(e)$ for each movie e is given by the total number of ratings for e divided by the total number of users in the dataset. The optimal solution A^* is computed greedily with respect to \bar{w} using Algorithm 1. In each round, we recommend 17 movies, i.e., $|A^*| = |A^t| = 17$. The privacy budget for each $e \in A^t$ is $\varepsilon_0 = \varepsilon/17 = 2/17$. The randomness comes from the fact that in each round t , a random user is selected. For each movie $e \in A^t$, if e is rated by that selected random user, then the reward $w_t(e)$ is set to 1.

Figure 3 shows the results for DPUCB-MAT, OMM, and Optimal Policy. We observe that the expected per-round return of DPUCB-MAT is comparable to that of the non-private baseline (OMM) and approaches the optimal policy. Table 3 lists the overlapping movies learned by DPUCB-MAT and Optimal Policy at the end of learning. We can see that the movie genres of those movies appear to be diverse. Experimental results for DPTS-MAT are deferred to Appendix A.

6 CONCLUSIONS AND FUTURE DIRECTIONS

In this work, we have shown that learning the maximum weight basis for matroid bandits can be done efficiently in a differentially private manner. We propose two simple differentially private algorithms DPUCB-MAT and DPTS-MAT and conduct experiments to evaluate their practical performance on both synthetic and real-world datasets. There are still some problems remaining for this work. So far, we only have an $\Omega(LK \ln(n)/\Delta + LK \ln(n)/\varepsilon)$ regret lower bound for differentially private combinatorial bandits (Chen et al., 2020). By modifying their proof for Theorem 9, we conjecture that an improved $\Omega(L \ln(n)/\Delta + LK \ln(n)/\varepsilon)$ regret lower bound for differentially private matroid bandits is achievable. However, our Theorem 2 is still an extra $\min\{K, \ln(n)\}$ factor far from this conjectured regret lower bound. We are not sure yet whether our derived regret bound is not tight or whether a better regret lower bound exists for differentially private matroid bandits.

ACKNOWLEDGEMENTS

This work was supported by the Alberta Machine Intelligence Institute (Amii).

REFERENCES

- Shipra Agrawal and Navin Goyal. Near-optimal regret bounds for Thompson Sampling. *Journal of the ACM (JACM)*, 64(5):1–24, 2017.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.

- Achraf Azize and Debabrota Basu. When privacy meets partial information: A refined analysis of differentially private bandits. *arXiv preprint arXiv:2209.02570*, 2022.
- 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.
- Olivier Chapelle and Lihong Li. An empirical evaluation of thompson sampling. *Advances in neural information processing systems*, 24, 2011.
- Xiaoyu Chen, Kai Zheng, Zixin Zhou, Yunchang Yang, Wei Chen, and Liwei Wang. (locally) differentially private combinatorial semi-bandits. In *International Conference on Machine Learning*, pp. 1757–1767. PMLR, 2020.
- Cynthia Dwork, Moni Naor, Toniann Pitassi, and Guy N Rothblum. Differential privacy under continual observation. In *Proceedings of the forty-second ACM symposium on Theory of computing*, pp. 715–724, 2010.
- Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- Aurélien Garivier and Olivier Cappé. The KL-UCB algorithm for bounded stochastic bandits and beyond. In *Proceedings of the 24th annual conference on learning theory*, pp. 359–376, 2011.
- F. Maxwell Harper and Joseph A. Konstan. The movielens datasets: History and context, 2015.
- Bingshan Hu and Nidhi Hegde. Near-optimal Thompson Sampling-based algorithms for differentially private stochastic bandits. In *Uncertainty in Artificial Intelligence*, pp. 844–852. PMLR, 2022.
- Bingshan Hu, Zhiming Huang, and Nishant A Mehta. Optimal algorithms for private online learning in a stochastic environment. *arXiv preprint arXiv:2102.07929*, 2021.
- Emilie Kaufmann, Nathaniel Korda, and Rémi Munos. Thompson Sampling: An asymptotically optimal finite-time analysis. In *International Conference on Algorithmic Learning Theory*, pp. 199–213. Springer, 2012.
- Branislav Kveton, Zheng Wen, Azin Ashkan, Hoda Eydgahi, and Brian Eriksson. Matroid bandits: Fast combinatorial optimization with learning. *arXiv preprint arXiv:1403.5045*, 2014.
- T.L Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1):4–22, 1985. ISSN 0196-8858. doi: [https://doi.org/10.1016/0196-8858\(85\)90002-8](https://doi.org/10.1016/0196-8858(85)90002-8). URL <https://www.sciencedirect.com/science/article/pii/0196885885900028>.
- Nikita Mishra and Abhradeep Thakurta. (nearly) optimal differentially private stochastic multi-arm bandits. In *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence*, pp. 592–601, 2015.
- Arvind Narayanan and Vitaly Shmatikov. Robust de-anonymization of large sparse datasets. In *2008 IEEE Symposium on Security and Privacy (sp 2008)*, pp. 111–125. IEEE, 2008.
- David L Neel and Nancy Ann Neudauer. Matroids you have known. *Mathematics magazine*, 82(1):26–41, 2009.
- Adam Smith and Jonathan Ullman. Privacy in statistics and machine learning course. <https://dpcourse.github.io/2021-spring/schedule.html>, 2021.
- Mohammad Sadegh Talebi and Alexandre Proutiere. An optimal algorithm for stochastic matroid bandit optimization. In *Proceedings of the 2016 International Conference on Autonomous Agents & Multiagent Systems*, pp. 548–556, 2016.
- Siwei Wang and Wei Chen. Thompson Sampling for combinatorial semi-bandits. In *International Conference on Machine Learning*, pp. 5114–5122. PMLR, 2018.

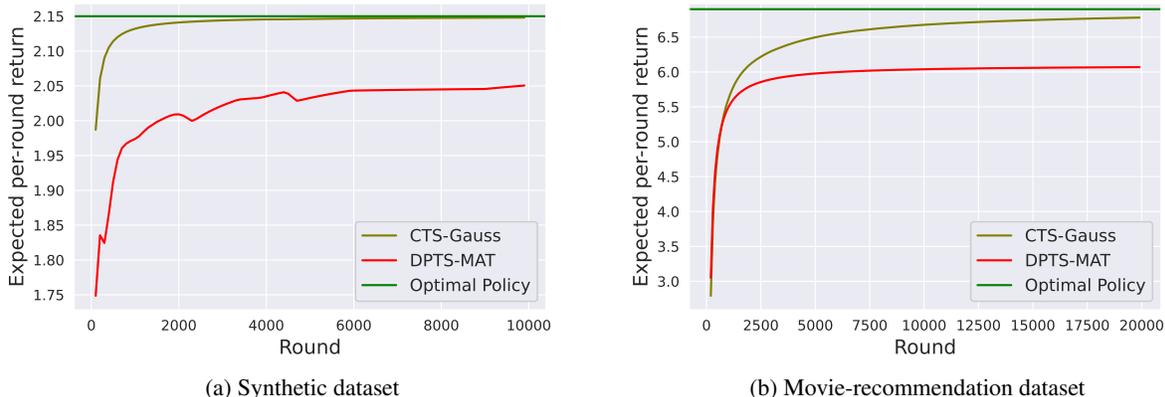


Figure 4: **Performance of DPTS-MAT on synthetic and movie-recommendation datasets.** Figure 4a shows the performances of DPTS-MAT (Algorithm 3), non-private Combinatorial Thompson Sampling (CTS) with Gaussian priors, and Optimal Policy on the synthetic dataset. Figure 4b shows the same algorithms on the movie-recommendation dataset. For both the plots, we again set $\varepsilon = 2$.

Appendices

A Additional Experiments	12
B Regret analysis of Algorithm 2	13
C Regret analysis of Algorithm 3	17

A ADDITIONAL EXPERIMENTS

In this section, we provide more experimental results. Figure 4 shows the performance of DPTS-MAT (Algorithm 3) on both the synthetic dataset (Figure 4a) and the movie-recommendation dataset (Figure 4b) with $\varepsilon = 2$. We compare DPTS-MAT against the non-private Combinatorial Thompson Sampling (CTS) with Gaussian priors and the Optimal Policy. We observe that the expected per-round return of DPTS-MAT is comparable to that of the non-private baseline and approaches the optimal policy for both datasets.

Figure 5a shows the accumulated regret R_n till round n over $\ln(n)$ for $n = 1, 2, \dots, 10^6$ rounds for both our private algorithms DPUCB-MAT and DPTS-MAT on the synthetic dataset. The quantity $\lim_{n \rightarrow \infty} R_n / \ln(n)$ characterizes the asymptotic rate of growth of the regret (Lai & Robbins, 1985). If a learning algorithm has a smaller rate, it suffers less regret. We observe that both DPUCB-MAT and DPTS-MAT converge but with different rates. In addition, we also observe that the Thompson Sampling-based algorithm, DPTS-MAT, empirically outperforms the UCB-based algorithm, DPUCB-MAT. Figure 5b studies the relationship between the accumulated regret R_n and $1/\varepsilon$ for DPUCB-MAT and DPTS-MAT on the synthetic dataset. The values of ε are 50 evenly spaced numbers between 0.5 and 50. We run our algorithms for $n = 10,000$ rounds for each value of ε and plot the accumulated regret. From the experimental results, we can see that for both the algorithms the regret trend is linear in $1/\varepsilon$ and the Thompson Sampling-based algorithm empirically outperforms the UCB-based one.

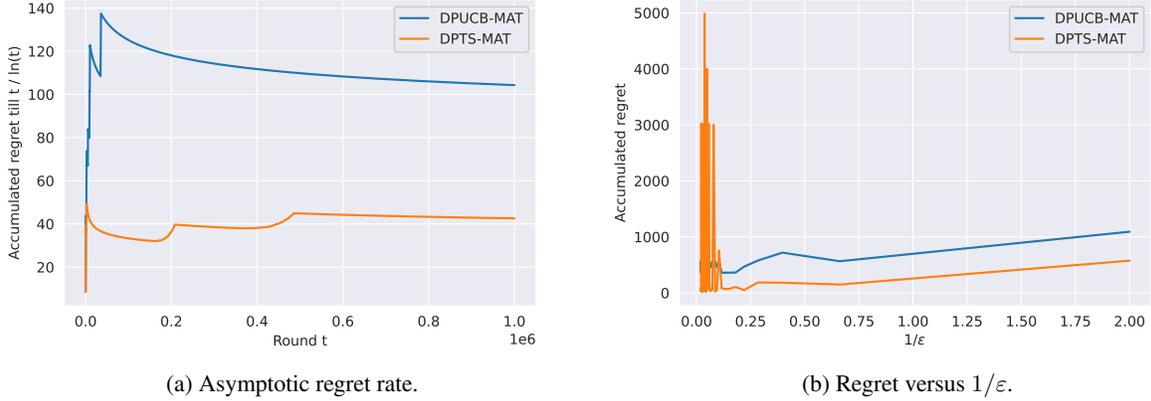


Figure 5: **Regret of DPUCB-MAT and DPTS-MAT on the synthetic dataset.** Figure 5a shows the accumulated regret R_n till round n divided by $\ln(n)$ for $n = 1, 2, \dots, 10^6$ rounds for both DPUCB-MAT and DPTS-MAT. The expression $\lim_{n \rightarrow \infty} R_n / \ln(n)$ characterizes the asymptotic growth rate of the regret. Figure 5b shows the accumulated regret versus $1/\varepsilon$ for both DPUCB-MAT and DPTS-MAT.

B REGRET ANALYSIS OF ALGORITHM 2

Recall that the regret can be expressed as

$$R_n = \sum_{k=1}^K \sum_{e \in \bar{A}^*: \Delta_{e,k} > 0} \underbrace{\sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, e \leftrightarrow a_k^* \right\} \right]}_{I_{e,k}} \cdot \Delta_{e,k} . \quad (8)$$

The indicator function $\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, e \leftrightarrow a_k^* \right\}$ will be 1 when the learner selects a sub-optimal base arm e instead of the optimal base arm a_k^* . This further implies the differentially private upper confidence bound of e is no smaller than that of a_k^* in that round.

Proof. Proof of Theorem 2: We first show the regret is at most

$$\sum_{e \in \bar{A}^*: \Delta_{e,\min} > 0} O \left(\frac{\ln(Kn)}{\Delta_{e,\min}} + \frac{K \ln(Kn)}{\varepsilon/K} \right) . \quad (9)$$

Then, we show the regret is also upper bounded by

$$\sum_{e \in \bar{A}^*: \Delta_{e,\min} > 0} O \left(\frac{\ln(Kn)}{\Delta_{e,\min}} + \frac{\ln^2(Kn)}{\varepsilon/K} \right) . \quad (10)$$

Combining these two claims concludes the proof.

Recall $\varepsilon_0 = \varepsilon/K$ and let $l_{e,k} = O \left(\frac{\ln(Kn)}{\Delta_{e,k} \cdot \min\{\Delta_{e,k}, \varepsilon_0\}} \right)$, where the big-O notation only hides a universal constant.

We decompose $I_{e,k}$ in (8) as

$$I_{e,k} = \underbrace{\sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, e \leftrightarrow a_k^*, T_e(t-1) \leq l_{e,k} \right\} \right]}_{\Gamma_{1,k,e}} \Delta_{e,k} + \underbrace{\sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, e \leftrightarrow a_k^*, T_e(t-1) > l_{e,k} \right\} \right]}_{\Gamma_{2,k,e}} \Delta_{e,k} . \quad (11)$$

From Hu et al. (2021); Azize & Basu (2022), we know that $\Gamma_{2,k,e} = O(1/K^2)$. Now, we use similar arguments as the one shown in Kveton et al. (2014) to complete the proof. As for a fixed $e \in \bar{A}^*$, we maintain a counter $T_e(t-1)$

during the learning. The counter counts the number of observations that are used to compute the differentially private empirical mean and the values of the counter double each time. For all the rounds when the values of the counter are in the range of $[0, l_{e,1}]$, the total regret among all these rounds is at most $\Delta_{e,1} \cdot O(l_{e,1})$. For all the rounds when the values of the counter are in the range of $[l_{e,1} + 1, l_{e,2}]$, the total regret among all these rounds is upper bounded by $\Delta_{e,2} \cdot O(l_{e,2} - l_{e,1}) + \Gamma_{2,1,e} \leq \Delta_{e,2} \cdot O(l_{e,2} - l_{e,1}) + O(1/K)$. Finally, for all the rounds when the values of the counter are in the range of $[l_{e,K-1} + 1, l_{e,K}]$, the total regret among all these rounds is at most $\Delta_{e,K} \cdot O(l_{e,K} - l_{e,K-1}) + \sum_{k=1}^{K-1} \Gamma_{2,k,e} \leq \Delta_{e,K} \cdot O(l_{e,K} - l_{e,K-1}) + O(1/K)$. For all the rounds after the counter hits $l_{e,K}$, the total regret is at most $\sum_{k=1}^K \Gamma_{2,k,e} \leq O(1/K)$.

By using Lemma 1, we have

$$\Delta_{e,1} \cdot O(l_{e,1}) + \sum_{k=2}^K \Delta_{e,k} \cdot O(l_{e,k} - l_{e,k-1}) + \sum_{k=2}^{K+1} \sum_{q=1}^{k-1} \Gamma_{2,q,e} \leq O\left(\frac{\ln(Kn)}{\Delta_{e,\min}} + \frac{K \ln(Kn)}{\varepsilon_0}\right), \quad (12)$$

which yields the result in the first claim, i.e.,

$$\sum_{e \in \bar{A}^* : \Delta_{e,\min} > 0} O\left(\frac{\ln(Kn)}{\Delta_{e,\min}} + \frac{K \ln(Kn)}{\varepsilon/K}\right).$$

We can also use the well-known ‘‘doubling-trick’’ in the bandit literature to have an upper bound.

Let $r_{\max,e} := \min\left\{\log\left(\frac{1}{\Delta_{e,\min}}\right), \log(Kn)\right\}$.

For any $1 \leq r \leq r_{\max,e}$, let $\Phi_r := \{k \in [K] : \Delta_{e,k} \in [0.5^r, 0.5^{r-1}]\}$ and $l_{e,r} = O\left(\frac{\ln(Kn)}{0.5^r \cdot \min\{0.5^r, \varepsilon_0\}}\right)$, where the big-O notation only hides a universal constant.

Then, we take the following regret decomposition. We have

$$\begin{aligned} R_n &= \sum_{k=1}^K \sum_{e \in \bar{A}^* : \Delta_{e,k} > 0} \sum_{t=1}^n \mathbb{E}\left[\mathbb{1}\left\{a_{\pi_t^{-1}(k)}^t = e, e \leftrightarrow a_k^*\right\}\right] \cdot \Delta_{e,k} \\ &\leq 1 + \sum_{e \in \bar{A}^*} \sum_{k \in [K] : \Delta_{e,k} \geq \frac{1}{Kn}} \sum_{t=1}^n \mathbb{E}\left[\mathbb{1}\left\{a_{\pi_t^{-1}(k)}^t = e, e \leftrightarrow a_k^*\right\}\right] \cdot \Delta_{e,k} \\ &\leq 1 + \sum_{e \in \bar{A}^*} \sum_{r=1}^{r_{\max,e}} \sum_{k \in \Phi_r} \sum_{t=1}^n \mathbb{E}\left[\mathbb{1}\left\{a_{\pi_t^{-1}(k)}^t = e, e \leftrightarrow a_k^*\right\}\right] \cdot \Delta_{e,k} \\ &\leq 1 + \sum_{e \in \bar{A}^*} \sum_{r=1}^{r_{\max,e}} \sum_{k \in \Phi_r} \sum_{t=1}^n \mathbb{E}\left[\mathbb{1}\left\{a_{\pi_t^{-1}(k)}^t = e, e \leftrightarrow a_k^*\right\}\right] \cdot 0.5^r \\ &\leq 1 + \sum_{e \in \bar{A}^*} \sum_{r=1}^{r_{\max,e}} \underbrace{\sum_{t=1}^n \mathbb{E}\left[\mathbb{1}\left\{\exists k \in \Phi_r : a_{\pi_t^{-1}(k)}^t = e, e \leftrightarrow a_k^*\right\}\right]}_{I_{e,r}} \cdot 0.5^r. \end{aligned} \quad (13)$$

Let $\left\{a_{\pi_t^{-1}(r)}^t = e, e \leftrightarrow r\right\}$ denote the event that $\left\{\exists k \in \Phi_r : a_{\pi_t^{-1}(k)}^t = e, e \leftrightarrow a_k^*\right\}$.

Now, we decompose $I_{e,r}$ as

$$\begin{aligned} I_{e,r} &= \underbrace{\sum_{t=1}^n \mathbb{E}\left[\mathbb{1}\left\{a_{\pi_t^{-1}(r)}^t = e, e \leftrightarrow r, T_e(t-1) \leq l_{e,r}\right\}\right]}_{\Gamma_{1,r,e}} \cdot 0.5^r \\ &\quad + \underbrace{\sum_{t=1}^n \mathbb{E}\left[\mathbb{1}\left\{a_{\pi_t^{-1}(r)}^t = e, e \leftrightarrow r, T_e(t-1) > l_{e,r}\right\}\right]}_{\Gamma_{2,r,e}} \cdot 0.5^r. \end{aligned} \quad (14)$$

Similarly, we have $\Gamma_{2,r,e} \leq O(1/K^2)$ and the total regret is at most

$$\sum_{e \in \bar{A}^*} \left(0.5^0 \cdot O(l_{e,1}) + \sum_{r=2}^{r_{\max,e}} 0.5^{r-1} \cdot O(l_{e,r} - l_{e,r-1}) + O(1) \right) \leq \sum_{e \in \bar{A}^*} O\left(\frac{\ln(Kn)}{\Delta_{e,\min}} + \frac{\ln^2(Kn)}{\varepsilon_0}\right), \quad (15)$$

where the inequality uses Lemma 2. \square

Lemma 1. Let $\Delta_1 \geq \dots \geq \Delta_K$ be a sequence of reals in $(0, 1]$. For any $\varepsilon_0 > 0$, we have

$$\Delta_1 \frac{1}{\Delta_1 \cdot \min\{\Delta_1, \varepsilon_0\}} + \sum_{k=2}^K \Delta_k \left(\frac{1}{\Delta_k \cdot \min\{\Delta_k, \varepsilon_0\}} - \frac{1}{\Delta_{k-1} \cdot \min\{\Delta_{k-1}, \varepsilon_0\}} \right) \leq \frac{2}{\Delta_K} + \frac{K}{\varepsilon_0}. \quad (16)$$

Lemma 2. For any $\varepsilon_0 > 0$, we have

$$0.5^0 \cdot l_{e,1} + \sum_{r=2}^{r_{\max,e}} 0.5^{r-1} \cdot (l_{e,r} - l_{e,r-1}) \leq O(\ln(Kn)) \cdot \left(\frac{2}{\Delta_{e,\min}} + \frac{\log(Kn)}{\varepsilon_0} \right). \quad (17)$$

Lemma 3. (Kveton et al., 2014, Lemma 3). Let $\Delta_1 \geq \dots \geq \Delta_K$ be a sequence of reals in $(0, 1]$. Then we have

$$\Delta_1 \frac{1}{\Delta_1^2} + \sum_{k=2}^K \Delta_k \left(\frac{1}{\Delta_k^2} - \frac{1}{\Delta_{k-1}^2} \right) \leq \frac{2}{\Delta_K}. \quad (18)$$

Proof. Proof of Lemma 1:

Case 1: When $\varepsilon_0 \geq \Delta_1$, we have

$$\text{LHS of (16)} = \Delta_1 \frac{1}{\Delta_1^2} + \sum_{k=2}^K \Delta_k \left(\frac{1}{\Delta_k^2} - \frac{1}{\Delta_{k-1}^2} \right) \leq \frac{2}{\Delta_K} \leq \frac{2}{\Delta_K} + \frac{K}{\varepsilon_0}, \quad (19)$$

where the first inequality uses Lemma 3.

Case 2: When $\Delta_K \geq \varepsilon_0$, we have

$$\text{LHS of (16)} = \frac{1}{\varepsilon_0} + \frac{1}{\varepsilon_0} \sum_{k=2}^K \left(1 - \frac{\Delta_k}{\Delta_{k-1}} \right) \leq \frac{1}{\varepsilon_0} + \frac{1}{\varepsilon_0} \cdot (K-1) \leq \frac{K}{\varepsilon_0} + \frac{2}{\Delta_K}.$$

Case 3: When $\Delta_1 \geq \dots \geq \Delta_j \geq \varepsilon_0 \geq \Delta_{j+1} \geq \dots \geq \Delta_K$ for some $j \in [K-1]$. We rewrite the LHS of (16) as

$$\begin{aligned} \sum_{k=1}^{K-1} \frac{\Delta_k - \Delta_{k+1}}{\Delta_k \cdot \min\{\Delta_k, \varepsilon_0\}} + \frac{1}{\min\{\Delta_K, \varepsilon_0\}} &= \sum_{k=1}^j \frac{\Delta_k - \Delta_{k+1}}{\Delta_k \cdot \varepsilon_0} + \sum_{k=j+1}^{K-1} \frac{\Delta_k - \Delta_{k+1}}{\Delta_k^2} + \frac{1}{\Delta_K} \\ &\leq \sum_{k=1}^j \left[\frac{1}{\varepsilon_0} - \frac{\Delta_{k+1}}{\varepsilon_0 \Delta_k} \right] + \sum_{k=j+1}^{K-1} \frac{\Delta_k - \Delta_{k+1}}{\Delta_k \cdot \Delta_{k+1}} + \frac{1}{\Delta_K} \\ &= \left[\frac{j}{\varepsilon_0} - \sum_{k=1}^j \frac{\Delta_{k+1}}{\varepsilon_0 \Delta_k} \right] + \left[\frac{1}{\Delta_K} - \frac{1}{\Delta_{j+1}} \right] + \frac{1}{\Delta_K} \\ &\leq \frac{j}{\varepsilon_0} + \frac{2}{\Delta_K} \\ &\leq \frac{K}{\varepsilon_0} + \frac{2}{\Delta_K}, \end{aligned} \quad (20)$$

which concludes the proof. \square

Proof of Lemma 2. We have the LHS in (17) is

$$\begin{aligned}
& 0.5^0 \cdot l_{e,1} + \sum_{r=2}^{r_{\max,e}} 0.5^{r-1} \cdot (l_{e,r} - l_{e,r-1}) \\
= & O(\ln(Kn)) \cdot \underbrace{\left(0.5^0 \frac{1}{0.5^0 \cdot \min\{0.5^0, \varepsilon_0\}} + \sum_{r=2}^{r_{\max,e}} 0.5^{r-1} \left(\frac{1}{0.5^{r-1} \cdot \min\{0.5^{r-1}, \varepsilon_0\}} - \frac{1}{0.5^r \cdot \min\{0.5^r, \varepsilon_0\}} \right) \right)}_I \\
\leq & O(\ln(Kn)) \cdot \left(\frac{2}{\Delta_{e,\min}} + \frac{\log(Kn)}{\varepsilon_0} \right) , \tag{21}
\end{aligned}$$

where the last step uses Lemma 1.

Note that $0.5^0 \geq 0.5^1 \geq \dots \geq 0.5^{r_{\max,e}}$ and from $r_{\max,e} = \min\{\log(1/\Delta_{e,\min}), \log(Kn)\}$, we have $2^{r_{\max,e}} \leq \Delta_{e,\min}$ and $r_{\max,e} \leq \log(Kn)$. Then, we have $I \leq \frac{2}{0.5^{r_{\max,e}}} + \frac{r_{\max,e}}{\varepsilon_0}$, which concludes the proof. \square

C REGRET ANALYSIS OF ALGORITHM 3

Still, recall the regret can be expressed as

$$R_n = \sum_{k=1}^K \sum_{e \in \bar{A}^* : \Delta_{e,k} > 0} \underbrace{\sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, e \leftrightarrow a_k^* \right\} \right]}_{I_{e,k}} \cdot \Delta_{e,k} . \quad (22)$$

The indicator function $\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, e \leftrightarrow a_k^* \right\}$ will be 1 when the learner selects a sub-optimal base arm e instead of the optimal base arm a_k^* . This implies that the posterior sample of e is no smaller than the posterior sample of a_k^* in that round. For a fixed k and $e \in \bar{A}^*$, we define event $\mathcal{E}_{e,k}^\theta(t) := \{\theta_e(t) \leq y_{e,k}\}$, where $y_{e,k} := \bar{w}(a_k^*) - \frac{1}{3}\Delta_{e,k}$, to decompose the regret. By introducing $\mathcal{E}_{e,k}^\theta(t)$, the regret can be decomposed as

$$\begin{aligned} R_n &= \sum_{k=1}^K \sum_{e \in \bar{A}^*} \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, e \leftrightarrow a_k^* \right\} \right] \cdot \Delta_{e,k} \\ &= \underbrace{\sum_{k=1}^K \sum_{e \in \bar{A}^*} \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \overline{\mathcal{E}_{e,k}^\theta(t)}, e \leftrightarrow a_k^* \right\} \right]}_V \cdot \Delta_{e,k} + \underbrace{\sum_{k=1}^K \sum_{e \in \bar{A}^*} \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t), e \leftrightarrow a_k^* \right\} \right]}_U \cdot \Delta_{e,k} . \end{aligned} \quad (23)$$

Bounding $V_{e,k}$ and V : The analysis of this term is very similar to the regret analysis of Algorithm 2. We first define two high probability events. Let $C_e(t) := \left\{ |\hat{w}_{e, T_e(t-1)}(t-1) - \bar{w}(e)| \leq \sqrt{\frac{3 \ln(Kt)}{T_e(t-1)}} \right\}$ be the event that the mean reward of $e \in E$ is within the confidence interval in round t and $G_e(t) := \left\{ |\hat{w}_{e, T_e(t-1)}(t-1) - \hat{w}_{e, T_e(t-1)}(t-1)| \leq \frac{3 \ln(Kt)}{\varepsilon_0 T_e(t-1)} \right\}$ be the event that the noise added is not too much in round t . Let $\overline{C_e(t)}$ and $\overline{G_e(t)}$ denote the complements of events $C_e(t)$ and $G_e(t)$, respectively.

Let $l_{e,k} := \frac{72 \ln(nK)}{\min\{\Delta_{e,k}^2, \varepsilon_0 \cdot \Delta_{e,k}\}}$. Then, $V_{e,k}$ can be further decomposed as

$$\begin{aligned} V_{e,k} &= \underbrace{\sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \overline{\mathcal{E}_{e,k}^\theta(t)}, T_e(t-1) \leq l_{e,k}, e \leftrightarrow a_k^* \right\} \right]}_{\Gamma_{1,k,e}} \cdot \Delta_{e,k} \\ &+ \underbrace{\sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \overline{\mathcal{E}_{e,k}^\theta(t)}, T_e(t-1) > l_{e,k}, e \leftrightarrow a_k^* \right\} \right]}_{\Gamma_{2,k,e}} \cdot \Delta_{e,k} . \end{aligned} \quad (24)$$

We further decompose $\Gamma_{2,k,e}$ as

$$\begin{aligned} \Gamma_{2,k,e} &= \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \overline{\mathcal{E}_{e,k}^\theta(t)}, T_e(t-1) > l_{e,k}, e \leftrightarrow a_k^* \right\} \right] \cdot \Delta_{e,k} \\ &\leq \underbrace{\sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \overline{\mathcal{E}_{e,k}^\theta(t)}, C_e(t), G_e(t), T_e(t-1) > l_{e,k}, e \leftrightarrow a_k^* \right\} \right]}_\gamma \cdot \Delta_{e,k} \\ &+ \underbrace{\sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ \overline{C_e(t)} \right\} \right]}_{\leq O\left(\frac{1}{K^2 n^2}\right), \text{ Lemma 4}} + \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ \overline{G_e(t)} \right\} \right] . \end{aligned} \quad (25)$$

Let $d_{e,k} := \lceil \log(l_{e,k}) \rceil$. Let τ_s denote the round by the end of which there are exactly 2^s fresh observations that will be used to compute the differentially private empirical mean of a sub-optimal base arm e . Now, we upper bound γ .

We have

$$\begin{aligned}
\gamma &= \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \overline{\mathcal{E}_{e,k}^\theta(t)}, C_e(t), G_e(t), e \leftrightarrow a_k^*, T_e(t-1) > l_{e,k} \right\} \right] \cdot \Delta_{e,k} \\
&\leq \sum_{s=d_e}^{\log n} \mathbb{E} \left[\sum_{t=\tau_s+1}^{\tau_{s+1}} \mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \overline{\mathcal{E}_{e,k}^\theta(t)}, C_e(t), G_e(t), e \leftrightarrow a_k^*, T_e(t-1) > l_{e,k} \right\} \right] \cdot \Delta_{e,k} \\
&\leq \sum_{s=d_e}^{\log n} \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \overline{\mathcal{E}_{e,k}^\theta(t)}, C_e(t), G_e(t), e \leftrightarrow a_k^*, T_e(t-1) = 2^s \right\} \right] \cdot \Delta_{e,k} \\
&= \sum_{s=d_e}^{\log n} \sum_{t=1}^n \mathbb{E} \left[\mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \overline{\mathcal{E}_{e,k}^\theta(t)}, C_e(t), G_e(t), e \leftrightarrow a_k^*, T_e(t-1) = 2^s \right\} \mid \mathcal{F}_{t-1} \right] \right] \cdot \Delta_{e,k} \\
&\leq \sum_{s=d_e}^{\log n} \sum_{t=1}^n \mathbb{E} \left[\underbrace{\mathbb{1} \{C_e(t), G_e(t), T_e(t-1) = 2^s\}}_{\Lambda} \underbrace{\mathbb{E} \left[\mathbb{1} \left\{ \overline{\mathcal{E}_{e,k}^\theta(t)} \right\} \mid \mathcal{F}_{t-1} \right]}_{\lambda} \right] \cdot \Delta_{e,k} .
\end{aligned} \tag{26}$$

To upper bound Λ , we divide all the instantiations F_{t-1} of \mathcal{F}_{t-1} into two groups based on whether $\mathbb{1} \{C_e(t), G_e(t), T_e(t) = 2^s\}$ is 1 or 0.

Case 1: For the instantiation F_{t-1} such that $\mathbb{1} \{C_e(t), G_e(t), T_e(t-1) = 2^s\} = 0$, we have $\Lambda = 0$.

Case 2: For the instantiation F_{t-1} such that $\mathbb{1} \{C_e(t), G_e(t), T_e(t-1) = 2^s\} = 1$, we construct an upper bound for λ . To do so, recall that $\theta_e(t) \sim \mathcal{N} \left(w'_{e, T_e(t-1)}(t), \frac{1}{T_e(t-1)} \right)$, where $w'_{e, T_e(t-1)}(t) = \tilde{w}_{e, T_e(t-1)}(t-1) + \frac{3 \ln(Kt)}{\varepsilon_0 \cdot T_e(t-1)}$ and $\mathcal{N}(\mu, \sigma^2)$ is a normal distribution with mean μ and variance σ^2 . We have

$$\begin{aligned}
w'_{e, T_e(t-1)}(t) &= \tilde{w}_{e, T_e(t-1)}(t-1) + \frac{3 \ln(Kt)}{\varepsilon_0 \cdot T_e(t-1)} \\
&\leq \hat{w}_{e, T_e(t-1)}(t-1) + \frac{6 \ln(Kt)}{\varepsilon_0 \cdot T_e(t-1)} \quad (\text{since } \mathbb{1} \{G_e(t-1)\} = 1) \\
&\leq \bar{w}(e) + \frac{6 \ln(Kt)}{\varepsilon_0 \cdot T_e(t-1)} + \sqrt{\frac{3 \ln(Kt)}{T_e(t-1)}} \quad (\text{since } \mathbb{1} \{C_e(t-1)\} = 1) \\
&= \bar{w}(e) + \frac{6 \ln(Kt)}{\varepsilon_0 \cdot 2^s} + \sqrt{\frac{3 \ln(Kt)}{2^s}} \quad (\text{since } \mathbb{1} \{T_e(t-1) = 2^s\} = 1) \\
&\leq \bar{w}(e) + \frac{\Delta_{e,k}}{12} + \frac{\Delta_{e,k}}{\sqrt{18}} \quad \left(\text{for } s \geq \left\lceil \log \frac{72 \ln(nK)}{\min \{\Delta_{e,k}^2, \varepsilon_0 \cdot \Delta_{e,k}\}} \right\rceil \right) \\
&\leq \bar{w}(e) + \frac{1}{3} \Delta_{e,k} .
\end{aligned}$$

From first-order stochastic dominance, we know a $\mathcal{N} \left(w'_{e, T_e(t-1)}(t), \frac{1}{T_e(t-1)} \right)$ distributed random variable is stochastically dominated by a $\mathcal{N} \left(\bar{w}(e) + \frac{1}{3} \Delta_{e,k}, \frac{1}{T_e(t-1)} \right)$ distributed random variable. Next, we slightly abuse the notation and use $\mathbb{P} \{ \mathcal{N}(\mu, \sigma^2) > y_{e,k} \}$ to denote the probability that a $\mathcal{N}(\mu, \sigma^2)$ distributed random variable is drawn greater than $y_{e,k}$. Now, we construct an upper bound for $\mathbb{E} \left[\mathbb{1} \left\{ \overline{\mathcal{E}_{e,k}^\theta(t)} \right\} \mid \mathcal{F}_{t-1} = F_{t-1} \right]$ by using the first-order stochastic

dominance between two Gaussian distributed random variables. We have

$$\begin{aligned}
\mathbb{E} \left[\mathbb{1} \left\{ \overline{\mathcal{E}}_{e,k}^\theta(t) \right\} \mid \mathcal{F}_{t-1} = F_{t-1} \right] &= \mathbb{P} \left\{ \theta_e(t) > \bar{w}(e) + \frac{2}{3} \Delta_{e,k} \right\} \\
&= \mathbb{P} \left\{ \mathcal{N} \left(w'_{e, T_e(t-1)}(t), \frac{1}{T_e(t-1)} \right) > \bar{w}(e) + \frac{2}{3} \Delta_{e,k} \right\} \\
&\leq \mathbb{P} \left\{ \mathcal{N} \left(\bar{w}(e) + \frac{1}{3} \Delta_{e,k}, \frac{1}{T_e(t-1)} \right) > \bar{w}(e) + \frac{2}{3} \Delta_{e,k} \right\} \\
&= \mathbb{P} \left\{ \mathcal{N} \left(\bar{w}(e), \frac{1}{T_e(t-1)} \right) > \bar{w}(e) + \frac{1}{3} \Delta_{e,k} \right\} \\
&\stackrel{(\clubsuit)}{\leq} \exp \left(-\frac{1}{18} \cdot \Delta_{e,k}^2 \cdot \frac{72 \ln(nK)}{\min \{ \Delta_{e,k}^2, \varepsilon_0 \cdot \Delta_{e,k} \}} \right) \\
&\leq O \left(\frac{1}{(nK)^4} \right),
\end{aligned} \tag{27}$$

where (\clubsuit) uses the concentration bounds for a normally distributed random variable (Fact 5) and the fact that $T_e(t-1) = 2^s \geq \frac{72 \ln(nK)}{\min \{ \Delta_{e,k}^2, \varepsilon_0 \cdot \Delta_{e,k} \}}$.

Putting all the pieces together, we have $\Gamma_{2,k,e} \leq O\left(\frac{1}{K^2}\right)$. Now, we use similar arguments that have been used for the proofs of Theorem 2 to upper bound term V in (23). Recall that $T_e(t-1)$ is a counter that counts the number of fresh observations that has been used to compute the differentially private empirical mean of a sub-optimal base arm e . For all the rounds when the counter is in the range of $[0, l_{e,1}]$, the total regret among all these rounds is at most $\Delta_{e,1} \cdot O(l_{e,1})$. For all the rounds when the values of the counter are in the range of $[l_{e,1} + 1, l_{e,2}]$, the total regret among all these rounds is upper bounded by $\Delta_{e,2} \cdot O(l_{e,2} - l_{e,1}) + \Gamma_{2,1,e} \leq \Delta_{e,2} \cdot O(l_{e,2} - l_{e,1}) + O(1/K)$. Finally, for all the rounds when the values of the counter are in the range of $[l_{e,K-1} + 1, l_{e,K}]$, the total regret among all these rounds is at most $\Delta_{e,K} \cdot O(l_{e,K} - l_{e,K-1}) + \sum_{k=1}^{K-1} \Gamma_{2,k,e} \leq \Delta_{e,K} \cdot O(l_{e,K} - l_{e,K-1}) + O(1/K)$. For all the rounds after the counter hits $l_{e,K}$, the total regret is at most $\sum_{k=1}^K \Gamma_{2,k,e} \leq O(1/K)$.

By using Lemma 1, we have

$$V \leq \sum_{e \in \bar{A}^*} \left(\Delta_{e,1} \cdot O(l_{e,1}) + \sum_{k=2}^K \Delta_{e,k} \cdot O(l_{e,k} - l_{e,k-1}) + \sum_{k=2}^{K+1} \sum_{q=1}^{k-1} \Gamma_{2,q,e} \right) \leq \sum_{e \in \bar{A}^*} O \left(\frac{\ln(Kn)}{\Delta_{e,\min}} + \frac{K \ln(Kn)}{\varepsilon_0} \right). \tag{28}$$

Similarly, by using Lemma 2, we have

$$V \leq \sum_{e \in \bar{A}^*} O(\ln(Kn)) \cdot \left(\frac{2}{\Delta_{e,K}} + \frac{\log(Kn)}{\varepsilon_0} \right) = \sum_{e \in \bar{A}^*} O \left(\frac{\ln(Kn)}{\Delta_{e,\min}} + \frac{\ln^2(Kn)}{\varepsilon_0} \right). \tag{29}$$

Combining (28) and (29), we have

$$V \leq \sum_{e \in \bar{A}^*} O(\ln(Kn)) \cdot \left(\frac{2}{\Delta_{e,K}} + \frac{\log(Kn)}{\varepsilon_0} \right) = \sum_{e \in \bar{A}^*} O \left(\frac{\ln(Kn)}{\Delta_{e,\min}} + \frac{\min \{ K, \ln(Kn) \} \cdot \ln(Kn)}{\varepsilon_0} \right). \tag{30}$$

Bounding $U_{e,k}$ and U : Recall

$$U = \sum_{k=1}^K \sum_{e \in \bar{A}^*} \underbrace{\sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t), e \leftrightarrow a_k^* \right\} \right]}_{U_{e,k}} \cdot \Delta_{e,k}. \tag{31}$$

Let $l_{e,k}^* := \left\lceil \frac{288 \ln(L^2(n+e^{32}))}{\Delta_{e,k}^2} \right\rceil$ and $d_{e,k}^* = \log(l_{e,k}^*)$. We first decompose $U_{e,k}$ as

$$\begin{aligned}
U_{e,k} &= \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t), e \leftrightarrow a_k^* \right\} \right] \cdot \Delta_{e,k} \\
&\leq \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t), G_{a_k^*}(t), e \leftrightarrow a_k^* \right\} \right] + \underbrace{\sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ \overline{G_{a_k^*}(t)} \right\} \right]}_{O\left(\frac{1}{K^2}\right), \text{Lemma 4}} \\
&= \underbrace{\sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t), G_{a_k^*}(t), e \leftrightarrow a_k^*, T_{a_k^*}(t-1) \leq l_{e,k}^* \right\} \right]}_{\Gamma_{1,k,e}} \cdot \Delta_{e,k} \\
&\quad + \underbrace{\sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t), G_{a_k^*}(t), e \leftrightarrow a_k^*, T_{a_k^*}(t-1) > l_{e,k}^* \right\} \right]}_{\Gamma_{2,k,e}} \cdot \Delta_{e,k} + O\left(\frac{1}{K^2}\right). \tag{32}
\end{aligned}$$

Let $Y_{e,k}^\theta(t) := \mathbb{P} \left\{ \theta_{a_k^*}(t) > y_{e,k} \mid \mathcal{F}_{t-1} \right\}$. Then, we have

$$\begin{aligned}
\Gamma_{1,k,e} &= \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t), G_{a_k^*}(t), e \leftrightarrow a_k^*, T_{a_k^*}(t-1) \leq l_{e,k}^* \right\} \right] \\
&= \sum_{t=1}^n \mathbb{E} \left[\mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t), G_{a_k^*}(t), e \leftrightarrow a_k^*, T_{a_k^*}(t-1) \leq l_{e,k}^* \right\} \mid \mathcal{F}_{t-1} \right] \right] \\
&= \sum_{t=1}^n \mathbb{E} \left[\mathbb{E} \left[\mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t), G_{a_k^*}(t), e \leftrightarrow a_k^* \right\} \cdot \mathbb{1} \left\{ T_{a_k^*}(t-1) \leq l_{e,k}^* \right\} \mid \mathcal{F}_{t-1} \right] \right] \\
&= \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ T_{a_k^*}(t-1) \leq l_{e,k}^*, e \leftrightarrow a_k^* \right\} \cdot \mathbb{P} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t), G_{a_k^*}(t) \mid \mathcal{F}_{t-1} \right\} \right] \\
&\stackrel{(a)}{\leq} \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ T_{a_k^*}(t-1) \leq l_{e,k}^*, e \leftrightarrow a_k^* \right\} \cdot \frac{1 - Y_{e,k}^\theta(t)}{Y_{e,k}^\theta(t)} \mathbb{P} \left\{ a_{\pi_t^{-1}(k)}^t = a_k^*, \mathcal{E}_{e,k}^\theta(t), G_{a_k^*}(t) \mid \mathcal{F}_{t-1} \right\} \right] \\
&\leq \sum_{t=1}^n \mathbb{E} \left[\underbrace{\mathbb{1} \left\{ T_{a_k^*}(t-1) \leq l_{e,k}^*, e \leftrightarrow a_k^* \right\}}_{\eta} \cdot \frac{1 - Y_{e,k}^\theta(t)}{Y_{e,k}^\theta(t)} \mathbb{1} \left\{ a_{\pi_t^{-1}(k)}^t = a_k^*, G_{a_k^*}(t) \right\} \right], \tag{33}
\end{aligned}$$

where inequality (a) uses Lemma 5.

We now reduce the proof to the non-private setting. First, we divide all the instantiations F_{t-1} of \mathcal{F}_{t-1} into two groups depending on whether $\mathbb{1} \left\{ G_{a_k^*}(t), T_{a_k^*}(t-1) \leq l_{e,k}^* \right\}$ is 1 or 0.

Case 1: For F_{t-1} such that $\mathbb{1} \left\{ G_{a_k^*}(t), T_{a_k^*}(t-1) \leq l_{e,k}^* \right\} = 0$, we have $\eta = 0$.

Case 2: For F_{t-1} such that $\mathbb{1} \left\{ G_{a_k^*}(t), T_{a_k^*}(t-1) \leq l_{e,k}^* \right\} = 1$, we have $w'_{a_k^*, T_{a_k^*}(t-1)}(t) = \tilde{w}_{a_k^*, T_{a_k^*}(t-1)}(t-1) + \frac{3 \ln(Kt)}{\varepsilon_0 \cdot T_{a_k^*}(t-1)} \geq \hat{w}_{a_k^*, T_{a_k^*}(t-1)}(t-1)$. Since a random variable drawn from $\mathcal{N}(\mu, \sigma^2)$ first-order stochastically dominates a random variable drawn from $\mathcal{N}(\mu', \sigma^2)$ if $\mu \geq \mu'$, we have

$$\frac{1 - Y_{e,k}^\theta(t)}{Y_{e,k}^\theta(t)} = \frac{\mathbb{P} \left\{ \theta_{a_k^*}(t) \leq y_{e,k} \mid \mathcal{F}_{t-1} = F_{t-1} \right\}}{\mathbb{P} \left\{ \theta_{a_k^*}(t) > y_{e,k} \mid \mathcal{F}_{t-1} = F_{t-1} \right\}} \leq \frac{\mathbb{P} \left\{ \hat{\theta}_{a_k^*}(t) \leq y_{e,k} \mid \mathcal{F}_{t-1} = F_{t-1} \right\}}{\mathbb{P} \left\{ \hat{\theta}_{a_k^*}(t) > y_{e,k} \mid \mathcal{F}_{t-1} = F_{t-1} \right\}},$$

where $\hat{\theta}_{a_k^*}(t) \sim \mathcal{N}\left(\hat{w}_{a_k^*, T_{a_k^*}(t-1)}(t-1), \frac{1}{T_{a_k^*}(t-1)}\right)$.

From these two cases, for any F_{t-1} , we have

$$\eta \leq \frac{\mathbb{P}\left\{\hat{\theta}_{a_k^*}(t) \leq y_{e,k} \mid \mathcal{F}_{t-1} = F_{t-1}\right\}}{\mathbb{P}\left\{\hat{\theta}_{a_k^*}(t) > y_{e,k} \mid \mathcal{F}_{t-1} = F_{t-1}\right\}} \mathbb{1}\left\{a_{\pi_t^{-1}(k)}^t = a_k^*, G_{a_k^*}(t)\right\}. \quad (34)$$

Now, the proof is reduced to the non-private setting. We divide all the n rounds depending on when the empirical mean of a_k^* changes, i.e., whether $\hat{w}_{a_k^*, T_{a_k^*}(t-1)}(t-1)$ changes. Let τ_s denote the round by the end of which we use fresh 2^s observations to update the empirical mean of a_k^* . Note that τ_s is random. Then, we have

$$\begin{aligned} (33) &\leq \sum_{s=0}^{d_{e,k}^*} \mathbb{E} \left[\sum_{t=\tau_s+1}^{\tau_{s+1}} \frac{\mathbb{P}\left\{\hat{\theta}_{a_k^*}(t) \leq y_{e,k} \mid \mathcal{F}_{t-1} = F_{t-1}\right\}}{\mathbb{P}\left\{\hat{\theta}_{a_k^*}(t) > y_{e,k} \mid \mathcal{F}_{t-1} = F_{t-1}\right\}} \cdot \mathbb{1}\left\{a_{\pi_t^{-1}(k)}^t = a_k^*\right\} \cdot \mathbb{1}\left\{G_{a_k^*}(t)\right\} \right] \\ &\leq \sum_{s=0}^{d_{e,k}^*} \mathbb{E} \left[\sum_{t=\tau_s+1}^{\tau_{s+1}} \frac{\mathbb{P}\left\{\hat{\theta}_{a_k^*}(t) \leq y_{e,k} \mid \mathcal{F}_{t-1} = F_{t-1}\right\}}{\mathbb{P}\left\{\hat{\theta}_{a_k^*}(t) > y_{e,k} \mid \mathcal{F}_{t-1} = F_{t-1}\right\}} \cdot \mathbb{1}\left\{a_{\pi_t^{-1}(k)}^t = a_k^*\right\} \right] \\ &\leq \sum_{s=0}^{d_{e,k}^*} \mathbb{E} \left[2^{s+1} \cdot \frac{\mathbb{P}\left\{\hat{\theta}_{a_k^*}(\tau_s+1) \leq y_{e,k} \mid \mathcal{F}_{\tau_s} = F_{\tau_s}\right\}}{\mathbb{P}\left\{\hat{\theta}_{a_k^*}(\tau_s+1) > y_{e,k} \mid \mathcal{F}_{\tau_s} = F_{\tau_s}\right\}} \right] \\ &\leq O\left(\frac{\ln(Kn)}{\Delta_{e,k}^2}\right), \end{aligned} \quad (35)$$

where the last inequality uses Fact 4.

Similarly, we will have

$$\begin{aligned} \Gamma_{2,k,e} &\leq \sum_{s=d_{e,k}^*+1}^{\log(n)} \mathbb{E} \left[2^{s+1} \cdot \frac{\mathbb{P}\left\{\hat{\theta}_{a_k^*}(\tau_s+1) \leq y_{e,k} \mid \mathcal{F}_{\tau_s} = F_{\tau_s}\right\}}{\mathbb{P}\left\{\hat{\theta}_{a_k^*}(\tau_s+1) > y_{e,k} \mid \mathcal{F}_{\tau_s} = F_{\tau_s}\right\}} \right] \cdot \Delta_{e,k} \\ &\leq O\left(\frac{1}{L^2}\right). \end{aligned} \quad (36)$$

We use the following arguments to complete the proof. Now, as we are tracking the number of observations of the optimal base arm a_k^* , i.e., we are tracking $T_{a_k^*}(t-1)$. For a fixed k , we arrange all the mean reward gaps $\Delta_{e,k}$ for all $e \in \bar{A}^*$ in a descending order $\Delta_{e_1,k} \geq \Delta_{e_2,k} \geq \dots \geq \Delta_{e_{L-K},k} =: \Delta_{e_{\min},k}$.

For all the rounds when the counter is in the range of $[0, l_{e_1,k}^*]$, the total regret among all these rounds is at most $\Delta_{e_1,k} \cdot O(l_{e_1,k}^*)$. When the counter is in the range of $[l_{e_1,k}^* + 1, l_{e_2,k}^*]$, the total regret is at most $\Delta_{e_2,k} \cdot O(l_{e_2,k}^* - l_{e_1,k}^*) + \Gamma_{2,k,e_1} \leq \Delta_{e_2,k} \cdot O(l_{e_2,k}^* - l_{e_1,k}^*) + O(1/K)$. Finally, when the counter is in the range of $[l_{e_{L-K-1},k}^* + 1, l_{e_{L-K},k}^*]$, the total regret among all these rounds is at most $\Delta_{e_{L-K-1},k} \cdot O(l_{e_{L-K-1},k}^* - l_{e_{L-K},k}^*) + \sum_{q=1}^{L-K-1} \Gamma_{2,k,e_q} \leq \Delta_{e_{L-K-1},k} \cdot O(l_{e_{L-K-1},k}^* - l_{e_{L-K},k}^*) + O(1/K)$. For all the rounds after the counter hits $l_{e_{L-K},k}^*$, the total regret is at most $\sum_{k=1}^K O(1/L^2) \leq O(1/K)$.

By combining all these pieces together and using Lemma 3, we have

$$U \leq \sum_{k=1}^K O\left(\frac{\ln(Kn)}{\Delta_{k,\min}}\right). \quad (37)$$

Lemma 4. We have

$$\sum_{t=1}^n \mathbb{E} \left[\mathbb{1}\left\{\overline{G_e(t)}\right\} \right] \leq O\left(\frac{1}{K^2 n^2}\right),$$

and

$$\sum_{t=1}^n \mathbb{E} \left[\mathbb{1}\left\{\overline{C_e(t)}\right\} \right] \leq O\left(\frac{1}{K^2 n^2}\right).$$

Proof. Proof of Lemma 4: The proofs use concentration bounds of Laplace distributions and Hoeffding's inequality. We have

$$\begin{aligned}
\sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ \overline{G_e(t)} \right\} \right] &= \sum_{t=1}^n \mathbb{P} \left\{ \left| \tilde{w}_{e, T_e(t-1)}(t-1) - \hat{w}_{e, T_e(t-1)}(t-1) \right| > \frac{3 \ln(Kt)}{\varepsilon_0 \cdot T_e(t-1)} \right\} \\
&\leq \sum_{t=1}^n \sum_{s=0}^{\log(t-1)} \mathbb{P} \left\{ \left| \tilde{w}_{e, 2^s}(t-1) - \hat{w}_{e, 2^s}(t-1) \right| > \frac{3 \ln(Kt)}{\varepsilon_0 \cdot 2^s} \right\} \\
&\leq \sum_{t=1}^n \sum_{s=0}^{\log(t-1)} \mathbb{P} \left\{ \left| 2^s \cdot \tilde{w}_{e, 2^s}(t-1) - 2^s \cdot \hat{w}_{e, 2^s}(t-1) \right| > \frac{3 \ln(Kt)}{\varepsilon_0} \right\} \\
&= \sum_{t=1}^n \sum_{s=0}^{\log(t-1)} e^{-3 \ln(Kt)} \\
&\leq O \left(\frac{1}{K^2 n^2} \right),
\end{aligned}$$

where the second Inequality used the concentration bound of a Laplace random variable (Fact 3).

Similarly, we have

$$\begin{aligned}
\sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ \overline{C_e(t)} \right\} \right] &= \sum_{t=1}^n \mathbb{P} \left\{ \left| \hat{w}_{e, T_e(t-1)}(t-1) - \bar{w}(e) \right| > \sqrt{\frac{3 \ln(Kt)}{T_e(t-1)}} \right\} \\
&\leq \sum_{t=1}^n \sum_{s=0}^{\log(t-1)} \mathbb{P} \left\{ \left| \hat{w}_{e, 2^s}(t-1) - \bar{w}(e) \right| > \sqrt{\frac{3 \ln(Kt)}{2^s}} \right\} \\
&\leq \sum_{t=1}^n \sum_{s=0}^{\log(t-1)} 2e^{-2 \cdot 2^s \cdot \frac{3 \ln(Kt)}{2^s}} \\
&\leq O \left(\frac{1}{K^2 n^2} \right),
\end{aligned}$$

where the second inequality uses Hoeffding's inequality (Fact 2). \square

Lemma 5. For all t and for any instantiation F_{t-1} of \mathcal{F}_{t-1} , we have

$$\mathbb{P} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t), G_{a_k^*}(t) \mid \mathcal{F}_{t-1} = F_{t-1} \right\} \leq \frac{1 - Y_{e,k}^\theta(t)}{Y_{e,k}^\theta(t)} \mathbb{P} \left\{ a_{\pi_t^{-1}(k)}^t = a_k^*, \mathcal{E}_{e,k}^\theta(t), G_{a_k^*}(t) \mid \mathcal{F}_{t-1} = F_{t-1} \right\}. \quad (38)$$

Proof of Lemma 5. We start by noting that the event $G_{a_k^*}(t)$ is determined by the history \mathcal{F}_{t-1} .

Case 1: If F_{t-1} is the one such that $G_{a_k^*}(t)$ is false, then both sides of the inequality shown in (38) are zero, and the inequality trivially holds.

Case 2: If F_{t-1} is the one such that $G_{a_k^*}(t)$ is true, we can omit $G_{a_k^*}$ in both sides in (38). Let $\pi_t^{-1}(k) = i$ and $A_{i-1}^t = \{a_1^t, \dots, a_{i-1}^t\}$ be the set of the first $i-1$ base arms selected greedily by Algorithm 3. To complete the proof, it suffices to show

$$\mathbb{P} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t) \mid \mathcal{F}_{t-1} = F_{t-1}, A_{i-1}^t \right\} \leq \frac{1 - Y_{e,k}^\theta(t)}{Y_{e,k}^\theta(t)} \mathbb{P} \left\{ a_{\pi_t^{-1}(k)}^t = a_k^*, \mathcal{E}_{e,k}^\theta(t) \mid \mathcal{F}_{t-1} = F_{t-1}, A_{i-1}^t \right\}. \quad (39)$$

Let $\sigma(A_{i-1}^t) = \{e : e \in E \setminus A_{i-1}^t, A_{i-1}^t \cup \{e\} \in \mathcal{I}\}$ be the set of base arms that can be added to the current solution set A_{i-1}^t . Note that $a_k^* \in \sigma(A_{i-1}^t)$ and $a_k^* \notin A_{i-1}^t$.

We first construct an upper bound for the LHS of (39). We have

$$\begin{aligned}
& \mathbb{P} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t) \mid \mathcal{F}_{t-1} = F_{t-1}, A_{i-1}^t \right\} \\
& \leq \mathbb{P} \left\{ \theta_j(t) \leq y_{e,k}, \forall j \in \sigma(A_{i-1}^t) \mid \mathcal{F}_{t-1} = F_{t-1}, A_{i-1}^t \right\} \\
& \stackrel{(\spadesuit)}{=} \mathbb{P} \left\{ \theta_{a_k^*}(t) \leq y_{e,k} \mid A_{i-1}^t, \mathcal{F}_{t-1} = F_{t-1} \right\} \cdot \mathbb{P} \left\{ \theta_j(t) \leq y_{e,k}, \forall j \in \sigma(A_{i-1}^t) \setminus \{a_k^*\} \mid \mathcal{F}_{t-1} = F_{t-1}, A_{i-1}^t \right\} \quad (40) \\
& = \mathbb{P} \left\{ \theta_{a_k^*}(t) \leq y_{e,k} \mid \mathcal{F}_{t-1} = F_{t-1} \right\} \cdot \mathbb{P} \left\{ \theta_j(t) \leq y_{e,k}, \forall j \in \sigma(A_{i-1}^t) \setminus \{a_k^*\} \mid \mathcal{F}_{t-1} = F_{t-1}, A_{i-1}^t \right\} \\
& = (1 - Y_{e,k}^\theta(t)) \cdot \mathbb{P} \left\{ \theta_j(t) \leq y_{e,k}, \forall j \in \sigma(A_{i-1}^t) \setminus \{a_k^*\} \mid \mathcal{F}_{t-1} = F_{t-1}, A_{i-1}^t \right\},
\end{aligned}$$

where (\spadesuit) uses the fact that $\theta_{a_k^*}(t)$ and all base arms in $\sigma(A_{i-1}^t)$ are independent.

Similarly, the RHS of (39) is lower bounded by

$$\begin{aligned}
& \mathbb{P} \left\{ a_{\pi_t^{-1}(k)}^t = a_k^*, \mathcal{E}_{e,k}^\theta(t) \mid \mathcal{F}_{t-1} = F_{t-1}, A_{i-1}^t \right\} \\
& \geq \mathbb{P} \left\{ \theta_{a_k^*}(t) > y_{e,k} \geq \theta_j(t), \forall j \in \sigma(A_{i-1}^t) \setminus \{a_k^*\} \mid \mathcal{F}_{t-1} = F_{t-1}, A_{i-1}^t \right\} \quad (41) \\
& = \mathbb{P} \left\{ \theta_{a_k^*}(t) > y_{e,k} \mid \mathcal{F}_{t-1} = F_{t-1}, A_{i-1}^t \right\} \cdot \mathbb{P} \left\{ \theta_j(t) \leq y_{e,k}, \forall j \in \sigma(A_{i-1}^t) \setminus \{a_k^*\} \mid \mathcal{F}_{t-1} = F_{t-1}, A_{i-1}^t \right\} \\
& = Y_{e,k}^\theta(t) \cdot \mathbb{P} \left\{ \theta_j(t) \leq y_{e,k}, \forall j \in \sigma(A_{i-1}^t) \setminus \{a_k^*\} \mid \mathcal{F}_{t-1} = F_{t-1}, A_{i-1}^t \right\}.
\end{aligned}$$

Combining (40) and (41) gives

$$\mathbb{P} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t) \mid \mathcal{F}_{t-1} = F_{t-1}, A_{i-1}^t \right\} \leq \frac{1 - Y_{e,k}^\theta(t)}{Y_{e,k}^\theta(t)} \mathbb{P} \left\{ a_{\pi_t^{-1}(k)}^t = a_k^*, \mathcal{E}_{e,k}^\theta(t) \mid \mathcal{F}_{t-1} = F_{t-1}, A_{i-1}^t \right\}. \quad (42)$$

To get the stated result, we use law of total expectation and the fact that $Y_{e,k}^\theta$ is determined by \mathcal{F}_{t-1} . We have

$$\begin{aligned}
& \mathbb{P} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t) \mid \mathcal{F}_{t-1} = F_{t-1} \right\} \\
& = \mathbb{E} \left[\mathbf{1} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t) \right\} \mid \mathcal{F}_{t-1} = F_{t-1} \right] \\
& = \mathbb{E} \left[\mathbb{P} \left\{ a_{\pi_t^{-1}(k)}^t = e, \mathcal{E}_{e,k}^\theta(t) \mid A_{i-1}^t, \mathcal{F}_{t-1} = F_{t-1} \right\} \mid \mathcal{F}_{t-1} = F_{t-1} \right] \quad (43) \\
& \leq \mathbb{E} \left[\frac{1 - Y_{e,k}^\theta(t)}{Y_{e,k}^\theta(t)} \mathbb{P} \left\{ a_{\pi_t^{-1}(k)}^t = a_k^*, \mathcal{E}_{e,k}^\theta(t) \mid A_{i-1}^t, \mathcal{F}_{t-1} = F_{t-1} \right\} \mid \mathcal{F}_{t-1} = F_{t-1} \right] \\
& = \frac{1 - Y_{e,k}^\theta(t)}{Y_{e,k}^\theta(t)} \mathbb{P} \left\{ a_{\pi_t^{-1}(k)}^t = a_k^*, \mathcal{E}_{e,k}^\theta(t) \mid \mathcal{F}_{t-1} = F_{t-1} \right\},
\end{aligned}$$

where the only inequality uses (42). \square

Fact 2. (Hoeffding's inequality). Let X_1, \dots, X_n be independent random variables with each $X_i \in [a_i, b_i]$. Then, for any $\epsilon > 0$, we have

$$\mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n (X_i - \mathbb{E}X_i) \right| \geq \epsilon \right) \leq 2 \exp \left(\frac{-2n^2 \epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2} \right). \quad (44)$$

Fact 3. ((Dwork et al., 2014, Fact 3.6); tail probability for Laplace distribution). If $Y \sim \text{Lap}(b)$, for any $0 < \delta \leq 1$, we have

$$\mathbb{P} \{ |Y| \geq b \ln(1/\delta) \} = \delta. \quad (45)$$

Fact 4. (Agrawal & Goyal, 2017, Lemma 2.13). Let τ_s be the round by the end of which we use fresh 2^s observations to update the empirical mean of an optimal base arm a_k^* . Then, we have

$$\mathbb{E} \left[\frac{1 - Y_{r,k}^{\hat{\theta}}(\tau_s + 1)}{Y_{r,k}^{\hat{\theta}}(\tau_s + 1)} \right] \leq \begin{cases} O(1) & \forall s, \\ O\left(\frac{1}{L^2 n}\right) & s \geq \log \left(l_{e,k}^* \right), \end{cases}$$

where $l_{e,k}^* := \left\lceil \log \left(\frac{288 \ln(L^2(n+e^{32}))}{\Delta_{e,k}^2} \right) \right\rceil$.

Fact 5. (Gaussian tail bound). Let X be a Gaussian distributed random variable with mean $\mathbb{E}[X]$ and variance σ^2 , then for any t we have

$$\mathbb{P} \{ X - \mathbb{E}[X] > t \} \leq e^{-\frac{t^2}{2\sigma^2}}. \quad (46)$$