

Optimal PAC Multiple Arm Identification with Applications to Crowdsourcing

Yuan Zhou
Carnegie Mellon University
yuanzhou@cs.cmu.edu

Xi Chen
UC Berkeley
xichen1987@berkeley.edu

Jian Li
Tsinghua University
lijian83@mail.tsinghua.edu.cn

November 13, 2013

Abstract

We study the problem of selecting K arms with the highest expected rewards in a stochastic N -armed bandit game. Instead of using existing evaluation metrics (e.g., misidentification probability (Bubeck et al., 2013) or the metric in EXPLORE- K (Kalyanakrishnan & Stone, 2010)), we propose to use *the aggregate regret*, which is defined as the gap between the average of the reward means of the optimal solution and that of our solution. Besides being a natural metric by itself, we argue that in many applications, such as our motivating example from the crowdsourcing, the aggregate regret bound is more suitable. In this paper, we propose a new PAC algorithm, which, with probability at least $1 - \delta$, identifies a set of K arms with regret at most ϵ . We provide a detailed analysis on the sample complexity of our algorithm. To complement, we establish a lower bound on the expected number of samples for the Bernoulli bandit and show that the sample complexity of our algorithm matches the lower bound. Finally, we report experimental results on both synthetic and real data sets, which demonstrates the superior performance of the proposed algorithm.

1 Introduction

We study the multiple arm identification problem in a stochastic multi-armed bandit game. More formally, assume that we are facing a bandit with n alternative arms, where each arm is associated with an unknown reward distribution supported on $[0, 1]$ with mean θ_i . Upon each sample (or “pull”) of a particular arm, the reward is an *i.i.d.* sample from the reward distribution. We sequentially decide which arm to pull next and then collect the reward by sampling that arm. The goal of our “top- K arm identification” problem is to identify a subset of K arms with the maximum total mean. The problem finds applications in a variety of areas, such as in industrial engineering (Koenig & Law, 1985), evolutionary computation (Schmidt et al., 2006) and medical domains (Thompson, 1933). Here, we highlight another application in *crowdsourcing*. In recent years, crowdsourcing services become increasingly popular for collecting labels of the data for many machine learning, data mining and analytical tasks. The readers may refer to (Raykar et al., 2010; Welinder et al., 2010; Karger et al., 2012; Zhou et al., 2012; Ho et al., 2013; Chen et al., 2013; Liu et al., 2013) and references therein for recent work on machine learning in crowdsourcing. In a typical crowdsourced labeling task, the requestor submits a batch of microtasks (e.g., unlabeled data) and workers from the crowd will complete the tasks and receive a small monetary reward for each task completion. Since some workers from the crowd are highly noisy and unreliable, it is important to first exclude those workers to obtain high-quality labels. An effective strategy is to test workers by a few gold samples, i.e., data with the known labels from the experts. Since testing workers with gold samples incurs nontrivial cost, it is desirable to select the best

K workers with the minimum number of queries. This problem can be cast into our top- K arm identification problem where each worker corresponds to an arm and the mean of the reward θ_i characterizes the worker's underlying reliability/quality.

More formally, assume that the arms are ordered by their means: $\theta_1 > \theta_2 > \dots > \theta_n$ and let T be the set of selected arms with size $|T| = K$. We define the *aggregate regret* (or *regret* for short) of T as:

$$\mathcal{L}_T = \frac{1}{K} \left(\sum_{i=1}^K \theta_i - \sum_{i \in T} \theta_i \right). \quad (1)$$

Our goal is to design an algorithm with low sample complexity and PAC (Probably Approximately Correct) style bounds. More specifically, given any fixed positive ϵ, δ , the algorithm should be able to identify a set T of K arms with $\mathcal{L}_T \leq \epsilon$ (we call such a solution an ϵ -optimal solution), with probability at least $1 - \delta$.

We first note that our problem strictly generalizes the previous work by (Even-Dar et al., 2006; Mannor & Tsitsiklis, 2004) for $K = 1$ to arbitrary positive integer K and hence is referred as multiple arm identification problem. Although the problem of choosing multiple arms has been studied in some existing work, e.g., (Bubeck et al., 2013; Audibert et al., 2013; Kalyanakrishnan & Stone, 2010; Kalyanakrishnan et al., 2012), our notion of aggregate regret is inherently different from previously studied evaluation metrics such as misidentification probability (MISID-PROB) (Bubeck et al., 2013) and EXPLORE- K (Kalyanakrishnan & Stone, 2010; Kalyanakrishnan et al., 2012). In particular, MISID-PROB controls the probability that the output set T is not exactly the same as the top K arms; and EXPLORE- K requires to return a set T where, with high confidence, the mean of each arm in T is ϵ -close to the K -th best arm. As we will explain in the related work section, our evaluation metric is a more suitable objective for many real applications, especially for the crowdsourcing.

We summarize our results in this paper as follows:

1. **Section 3 & 4:** We develop a new PAC algorithm with sample complexity $O\left(\frac{n}{\epsilon^2} \left(1 + \frac{\ln(1/\delta)}{K}\right)\right)$ for *arbitrarily small* δ in worst case scenarios. The analysis of the algorithm is in **Section 4**. It is interesting to compare this bound with the optimal $O\left(\frac{n}{\epsilon^2} \ln\left(\frac{1}{\delta}\right)\right)$ bound for $K = 1$ in (Even-Dar et al., 2006; Mannor & Tsitsiklis, 2004). For $K = 1$ (i.e., selecting the best arm), we match their result. When K is larger, our algorithm suggests that even less samples are needed. Intuitively, a larger K leads to a less stringent constraint for an ϵ -optimal solution and thus can tolerate more mistakes. Let us consider the following example. Assume all the arms have the same mean $1/2$, except for a random one with the mean $1/2 + 2\epsilon$. If $K = 1$, to obtain an ϵ -optimal solution, we essentially need to identify the special arm and thus need a lot of samples. However, if K is large, any subset of K arms would work fine since the regret is at most $2\epsilon/K$. Our algorithm bears some similarity with previous work, such as the halving technique in (Even-Dar et al., 2006; Kalyanakrishnan & Stone, 2010; Karnin et al., 2013) and idea of accept-reject in (Bubeck et al., 2013). However, the analysis is more involved than the case for $K = 1$ and needs to be done more carefully in order to achieve the above sample complexity.
2. **Section 5:** To complement the upper bound, we further establish a matching lower bound for Bernoulli bandits: any (deterministic or randomized) algorithm requires at least $\Omega\left(\frac{n}{\epsilon^2} \left(1 + \frac{\ln(1/\delta)}{K}\right)\right)$ samples to achieve an ϵ regret with the confidence $1 - \delta$. This shows that our algorithm achieves the optimal sample complexity for Bernoulli bandits. To this end, we show two different lower bounds, $\Omega\left(\frac{n}{\epsilon^2}\right)$ and $\Omega\left(\frac{n}{\epsilon^2} \frac{\ln(1/\delta)}{K}\right)$. The first bound is established via an reduction from our problem to the basic problem of distinguishing two similar Bernoulli arms (with mean $1/2$ and $1/2 + \epsilon$ respectively). The second one can be shown via a generalization of the argument in (Mannor & Tsitsiklis, 2004) for $K = 1$.
3. **Section 6:** Finally, we conduct extensive experiments on both simulated and real data sets. The experimental results demonstrate that our algorithm not only achieves lower regret than existing methods but also leads to a higher precision. Moreover, using our algorithm, the maximum number of samples taken from an individual arm is much smaller than that in the SAR algorithm (Bubeck et al., 2013).

. This property is desirable for crowdsourcing applications since it can be quite problematic, at least time-consuming, to test a single worker with too many gold samples.

2 Related Works

Multi-armed bandit problems have been extensively studied in the machine learning community over the past decade (see for example (Auer et al., 2002a,b; Beygelzimer et al., 2011; Bubeck & Cesa-Bianchi, 2012) and the references therein). In recent years, the multiple arm identification problem has received much attention and been investigated under different setups. For example, the work (Even-Dar et al., 2006; Mannor & Tsitsiklis, 2004; Audibert et al., 2010; Karnin et al., 2013) studied the special case when $K = 1$. When $K > 1$, Bubeck et al. (2013) proposed a SAR (Successive Accepts and Rejects) algorithm which tries to minimize the *misidentification probability* (MISID-PROB) given a fixed amount of budget. More precisely, SAR tries to minimize the value of $\Pr(T \neq \{1, \dots, K\})$. Another line of research (Kalyanakrishnan et al., 2012; Kalyanakrishnan & Stone, 2010) proposed to select a subset T of arms, such that with high probability, for all arms $i \in T$, $\theta_i > \theta_K - \epsilon$. Here, θ_K is the mean of the reward of the K -th best arm. We refer this metric to as the EXPLORE- K metric.

As compared to MISID-PROB and EXPLORE- K metrics, our notion of aggregate regret is inherently different from them and is a more suitable objective for many real applications. For example, MISID-PROB requires to identify the exact top- K arms, which is more stringent. When the gap of any consecutive pair θ_i and θ_{i+1} among the first $2K$ arms is extremely small (e.g., $o(\frac{1}{n})$), it requires a huge amount (e.g., $\omega(n^2)$) of samples to make the misidentification probability less than ϵ (Bubeck et al., 2013). While in our metric, any K arms among the first $2K$ arms forms an ϵ -optimal solution. In crowdsourcing applications, our main goal is not to select the exact top- K workers, but a pool of good enough workers with minimum amount of samples. We note that the *expected simple regret*, $\frac{1}{K}(\sum_{i=1}^K \theta_i - \mathbf{E}[\sum_{i \in T} \theta_i])$, was also considered in a number of prior works (Audibert et al., 2010; Bubeck et al., 2013; Audibert et al., 2013). In (Audibert et al., 2010; Bubeck et al., 2013), the expected simple regret was shown to be sandwiched by $\Delta \cdot \text{MISID-PROB}$ and MISID-PROB (for $K = 1$), where $\Delta = \theta_1 - \theta_2$. However, Δ can be arbitrarily small, hence MISID-PROB can be an arbitrarily bad bound for the regret. It is worthwhile noting that it is possible to obtain an expected simple regret of ϵ with at most $O(n^2/\epsilon)$ samples, using the semi-bandit regret bound in (Audibert et al., 2013)¹. In contrast, the goal of this paper is to develop an efficient algorithm to achieve an ϵ -regret with high probability, which is a stronger requirement than obtaining an ϵ -expected simple regret.

To compare our aggregate regret with the EXPLORE- K metric, let us consider another example where $\theta_1, \dots, \theta_{K-1}$ are much larger than θ_K and $\theta_{K+i} > \theta_K - \epsilon$ for $i = 1, \dots, K$. It is easy to see that the set $T = \{K+1, \dots, 2K\}$ also satisfies the requirement of EXPLORE- K ². However, the set T is far away from the optimal set with the regret much larger than ϵ . In crowdsourcing, the labeling performance can downgrade to a significant extent if the best set of workers (e.g., $\theta_1, \dots, \theta_{K-1}$ in the example) is left out of the solution.

3 Algorithm

Our algorithm OptMAI (Algorithm 1) takes three positive integers n, K, Q as the input, where n is the total number of arms, K is the number of arms we want to choose and Q is an upper bound on the total number of samples³. In Section 4, we show that $Q = O\left(\frac{n}{\epsilon^2} \left(1 + \frac{\ln(1/\delta)}{K}\right)\right)$ suffices to obtain an ϵ -optimal

¹ The result in (Audibert et al., 2013) was stated in terms of expected accumulative regret (i.e., the expected regret over Z time slots). By setting the number of time slots Z to be $\frac{n}{K\epsilon^2}$, and choosing a random action as the final solution among Z actions (see e.g., (Bubeck et al., 2009)), one can get an expected simple regret of ϵ .

² For this particular instance, it is unlikely that the algorithms proposed in (Kalyanakrishnan et al., 2012; Kalyanakrishnan & Stone, 2010) would choose $\{K+1, \dots, 2K\}$ as the solution, even though it is valid solution under their EXPLORE- K metric. However, it is not clear, from their theoretical analysis, how good their solution is, collectively, as compared with the best K arms.

³If Algorithm 1 stops at round $r = R$, the total number of samples is $(1 - \beta^R)Q < Q$.

Algorithm 1 Optimal Multiple Arm Identification (OptMAI)

Input: n, K, Q .**Initialization:** Active set of arms $S_0 = \{1, \dots, n\}$; $\beta = e^{0.2} \cdot \frac{3}{4}$; set of top arms $T_0 = \emptyset$. Let $r = 0$ **while** $|T_r| < K$ and $|S_r| > 0$ **do** **if** $|S_r| \geq 4K$ **then** $S_{r+1} = \text{QE}(S_r, K, \beta^r(1 - \beta)Q)$ $T_{r+1} = \emptyset$ **else** $(S_{r+1}, T_{r+1}) = \text{AR}(S_r, T_r, K, \beta^r(1 - \beta)Q)$ **end if** $r = r + 1$.**end while****Output:** The set of the selected K -arms T_r .

Algorithm 2 Quartile-Elimination(QE) (S, K, Q)

Input: S, K, Q .

1. Sample each arm $i \in S$ for $Q_0 = \frac{Q}{|S|}$ times and let $\hat{\theta}_i$ be the empirical mean of the i -th arm.
2. Find the first quartile (lower quartile) of the empirical mean $\hat{\theta}_a$, denoted by \hat{q} .

Output: The set $V = S \setminus \{i \in S : \hat{\theta}_i < \hat{q}\}$.

Algorithm 3 Accept-Reject(AR) (S, T, K, Q)

Input: S, T, K, Q and $s = |S|$.

1. Sample each arm $i \in S$ for $Q_0 = \frac{Q}{|S|}$ times and let $\hat{\theta}_i$ be the empirical mean of the i -th arm.
2. Let $K' = K - |T|$. Let $\hat{\theta}_{(K')}$ and $\hat{\theta}_{(K'+1)}$ be the K' -th and $(K' + 1)$ -th largest empirical mean. Define the empirical gap for each arm $i \in S$:

$$\hat{\Delta}_i = \max(\hat{\theta}_i - \hat{\theta}_{(K'+1)}, \hat{\theta}_{(K')} - \hat{\theta}_i) \quad (2)$$

while $|T| < K$ and $|S| > 3s/4$ **do** Let $a \in \arg \max_{i \in S} \hat{\Delta}_i$ and set $S = S \setminus \{a\}$. **if** $\hat{\theta}_a \geq \hat{\theta}_{(K'+1)}$ **then** Set $T = T \cup \{a\}$. **end if****end while****Output:** The set S and T .

solution with probability at least $1 - \delta$. OptMAI consists of two stages, the *Quartile-Elimination (QE) stage* and the *Accept-Reject (AR) stage*.

The QE stage proceeds in rounds. Each QE round calls the QE subroutine in [Algorithm 2](#), which requires three parameters S, K, Q . Here, S is the set of arms which we still want to pull and Q is the total number of samples in this round. We sample each arm in S for $Q/|S|$ times and then discard *a quarter* of arms with the minimum empirical mean. We note that in each call of the QE subroutine, we pass different Q values (exponentially decreasing), which keeps the total number of samples linear and is critical for achieving the optimal sample complexity. See [Algorithm 1](#) for the setting of the parameters. The QE stage repeatedly calls the QE subroutine until the number of remaining arms is at most $4K$.

Now, we enter the AR stage, which also runs in rounds. Each AR round ([Algorithm 3](#)) requires four

parameters, S, T, K, Q , where S, K, Q have the same meanings as in QE and T is the set of arms that we have decided to include in our final solution and thus will not be sampled any more. In each AR subroutine (Algorithm 3), we again sample each arm for $Q/|S|$ times. We define the *empirical gap* to be the absolute difference between the empirical mean of the i -th arm and the K' -th (or $K' + 1$ -th) largest empirical mean, where $K' = K - |T|$. We remove a quarter of arms with the largest empirical gaps. There are two types of those removed arms: those with the largest empirical means, which will be included in our final solution set T , and those with the smallest empirical means, which we discard from further consideration.

Remark 3.1. *To achieve the desired asymptotic bound on regret and sample complexity, the AR stage can be substituted by a simpler process which takes a uniform number of samples from each arm and chooses the K arms with the largest empirical means. However, we choose to present the AR subroutine because 1) it also meets the theoretical bound in Section 4; 2) the AR stage shows a much better empirical performance.*

4 Bounding the Regret and the Sample Complexity

We first provide an upper bound on the regret given the total budget Q ; then by a simple algebraic manipulation, we obtain the sample complexity to achieve ϵ -optimal solution with probability at least $1 - \delta$ for arbitrarily small δ . In Section 5, we will show that the obtained sample complexity is optimal.

Firstly, let us introduce some necessary notations. Let $\text{ind}_i(S)$ be the arm in the set S with the i -th largest mean and $\text{val}_C(S) = \mathbf{E}_{j \in [C]} \theta_{\text{ind}_i(S)} \triangleq \frac{1}{C} \sum_{i=1}^C \theta_{\text{ind}_i(S)}$ be the average of the reward means of the C best arms in S . Here $[C] \triangleq \{1, \dots, C\}$. Let $\text{tot}_C(S) = C \text{val}_C(S)$ be the sum of the means of the C best arms in S . In the next lemma, we first provide the regret bound for the QE algorithm.

Lemma 4.1. *Assume that $K \leq |S|/4$ and let V be the output of Algorithm 2, $\text{QE}(S, K, Q)$. For every $\delta > 0$, with probability $1 - \delta$, we have that $\text{val}_K(V) \geq \text{val}_K(S) - \epsilon$, where $\epsilon = \sqrt{\frac{|S|}{Q} \left(10 + \frac{4 \ln(2/\delta)}{K}\right)}$.*

The basic idea of the proof goes as follows. Let $p = \theta_{\text{ind}_{|S|/2}(S)}$ be the median of the θ in S and $\tau = \min_{i \in V} (\hat{\theta}_i)$ be the minimum empirical mean for the selected arms in V . For each arm i among the top K arms in S , we define the random variable $X_i = \mathbf{1}[\hat{\theta}_{\text{ind}_i(S)} < p + \frac{\epsilon}{2}]$ and $X = \mathbf{E}_{i \in [K]} (\theta_{\text{ind}_i(S)} - p) X_i$. We further define two events $\mathcal{E}_1 = \{X \leq \epsilon\}$ and $\mathcal{E}_2 = \{\tau < p + \frac{\epsilon}{2}\}$. Intuitively, \mathcal{E}_2 says that the threshold τ , as the third quartile, is not much larger than the expected median value (i.e., p). When \mathcal{E}_2 happens, \mathcal{E}_1 would give an upper bound on the regret $\text{val}_K(S) - \text{val}_K(V)$. We formalize this idea by first showing that the event \mathcal{E}_1 and \mathcal{E}_2 together imply that $\text{val}_K(V) \geq \text{val}_K(S) - \epsilon$. Then we then prove that with the definition of ϵ in the lemma statement, \mathcal{E}_1 and \mathcal{E}_2 both hold with probability at least $1 - \frac{\delta}{2}$, which concludes our proof by a union bound. The details are presented in the appendix.

Secondly, we provide the regret bound for the AR algorithm in the following lemma with the proof presented in the appendix.

Lemma 4.2. *Let (S', T') be the output of the algorithm $\text{AR}(S, T, K, Q)$. For every $\delta > 0$, with probability $1 - \delta$, we have that*

$$\text{tot}_{K-|T'|}(S') + \text{tot}_{|T'|}(T') \geq \text{tot}_{K-|T|}(S) + \text{tot}_{|T|}(T) - \epsilon K,$$

$$\text{where } \epsilon = \sqrt{\frac{|S|}{Q} \left(4 + \frac{\log(2/\delta)}{K}\right)}.$$

In each round of the AR-stage with $S = S_r$ and top arms $T = T_r$, the value $\frac{\text{tot}_{K-|T|}(S) + \text{tot}_{|T|}(T)}{K}$ is the best possible average of the mean of K arms. Lemma 4.2 provides the upper bound of the gap between this value on the output (T', S') by AR and the best possible one. This would further imply that this value of the output of Algorithm 1 is not far away from that of the real top- K arms.

With Lemma 4.1 and Lemma 4.2 in place, we prove the performance of Algorithm 1 in the next theorem.

Theorem 4.3. For every $\delta > 0$, with probability at least $1 - \delta$, the output of OptMAI algorithm T is an ϵ -optimal solution (i.e., $\text{val}_K(T) \geq \text{val}_K([n]) - \epsilon$) with

$$\begin{aligned}\epsilon &= \sqrt{\frac{n}{Q}} \left(44 \cdot \sqrt{10 + \frac{4 \ln(1/\delta) + 16}{K}} + \frac{72}{\sqrt{K}} \right) \\ &= O(1) \sqrt{\frac{n}{Q}} \left(1 + \sqrt{\frac{\ln(1/\delta)}{K}} \right).\end{aligned}\tag{3}$$

[Theorem 4.3](#) also provides us the sample complexity of [Algorithm 1](#) for pre-fixed (ϵ, δ) as stated in the next corollary.

Corollary 4.4. Fix positive constant $\epsilon, \delta > 0$. In order to get a set T of K elements such that $\text{val}_K(T) \geq \text{val}_K([n]) - \epsilon$ with probability $1 - \delta$, we need to run [Algorithm 1](#) with

$$Q = O\left(\frac{n}{\epsilon^2} \left(1 + \frac{\ln(1/\delta)}{K}\right)\right).\tag{4}$$

5 A Matching Lower Bound

In this section, we provide the lower bound for Bernoulli bandits where the reward of the i -th arm follows a Bernoulli distribution with mean θ_i . We prove that for any underlying $\{\theta_i\}_{i=1}^n$ and any randomized algorithm \mathcal{A} , the expected number of samples Q required to identify an ϵ -optimal solution with probability $1 - \delta$ is at least $\max\left\{\Omega\left(\frac{n \ln(1/\delta)}{\epsilon^2 K}\right), \Omega\left(\frac{n}{\epsilon^2}\right)\right\} = \Omega\left(\frac{n}{\epsilon^2} \left(\frac{\ln(1/\delta)}{K} + 1\right)\right)$. According to [Corollary 4.4](#), for Bernoulli bandits, our algorithm achieves this lower bound of the sample complexity. In particular, we separate the proof into two parts: in the first part, we show that $Q \geq \Omega\left(\frac{n}{\epsilon^2}\right)$; and $Q \geq \Omega\left(\frac{n \ln(1/\delta)}{\epsilon^2 K}\right)$ in the second.

5.1 First Lower Bound: $Q \geq \Omega\left(\frac{n}{\epsilon^2}\right)$

Theorem 5.1. Fix a real number ϵ , integers K, n , where $0 < \epsilon < 0.01$ and $10 \leq K \leq n/2$. Let \mathcal{A} be a possibly randomized algorithm, so that for any set of n Bernoulli arms with means $\theta_1, \theta_2, \dots, \theta_n$,

- \mathcal{A} makes at most Q samples in expectation;
- with probability at least 0.8, \mathcal{A} outputs a set T of size K with $\text{val}_K(T) \geq \text{val}_K([n]) - \epsilon$.

Then, we have that $Q \geq \Omega\left(\frac{n}{\epsilon^2}\right)$.

The proof of [Theorem 5.1](#) proceeds in two steps. In the first step, we show that we can construct an algorithm \mathcal{B} , using \mathcal{A} as a subroutine, to distinguish whether a single Bernoulli arm has mean $1/2$ or $1/2 + 4\epsilon$ with $O(Q/n)$ samples. In the second step, we utilize the well known lower bound result for distinguishing a Bernoulli arm to establish a lower bound for Q . Formally, we show the following lemma for the first step.

Lemma 5.2. Let \mathcal{A} be an algorithm in [Theorem 5.1](#), then there is an algorithm \mathcal{B} which correctly outputs whether a Bernoulli arm X has the mean $\frac{1}{2} + 4\epsilon$ or the mean $\frac{1}{2}$ with probability at least 0.51, and \mathcal{B} makes at most $\frac{200Q}{n}$ samples in expectation.

Given an algorithm \mathcal{A} stated in [Theorem 5.1](#), we construct the algorithm \mathcal{B} as follows:

- Choose a random subset $S \subseteq [n]$ such that $|S| = K$ and then choose a random element $j \in S$.
- For each $i \in [n], i \neq j$, let $\theta_i = \frac{1}{2} + 4\epsilon$ if $i \in S$, let $\theta_i = \frac{1}{2}$ otherwise.
- Simulate \mathcal{A} as follows: whenever \mathcal{A} samples the i -th arm, if $i = j$, we sample the Bernoulli arm X ; otherwise we sample the arm with mean θ_i .

- If the arm X is sampled by less than $\frac{200Q}{n}$ times and \mathcal{A} returns a set T such that $j \notin T$, we decide that X has the mean of $\frac{1}{2}$; otherwise we decide that X has the mean of $\frac{1}{2} + 4\epsilon$.

We note that the number of samples of \mathcal{B} increases by one whenever X is sampled. Since \mathcal{B} stops and outputs the mean $\frac{1}{2} + 4\epsilon$ if the number of samples on X reaches $\frac{200Q}{n}$, \mathcal{B} makes at most $\frac{200Q}{n}$ samples. The intuition why the above algorithm can separate X is as follows. If X has mean $1/2 + \epsilon$, X is no different from any other arm in S . Similarly, if X has mean $1/2$, X is the same as any other arm in $[n] \setminus S$. If \mathcal{A} satisfies the requirement in [Theorem 5.1](#), \mathcal{A} can identify a significant proportion of arms with mean $1/2 + 4\epsilon$. So if X has mean $1/2 + 4\epsilon$, there is a good chance (noticeably larger than 0.5) that X will be chosen by \mathcal{A} . In the appendix, we formally prove the correctness of \mathcal{B} , i.e., it can correctly output the mean of X with probability at least 0.51; and thus conclude the proof of [Lemma 5.2](#).

The second step of the proof of [Theorem 5.1](#) is a well-known lower bound on the expected sample complexity for separating a single Bernoulli arm ([Chernoff, 1972](#); [Anthony & Bartlett, 1999](#)).

Lemma 5.3. *Fix ϵ such that $0 < \epsilon < 0.01$ and let X be a Bernoulli random variable with the mean either $\frac{1}{2} + 4\epsilon$ or $\frac{1}{2}$. If an algorithm \mathcal{B} can output the correct mean of X with probability at least 0.51, then expected number of samples performed by \mathcal{B} is at least $\Omega(\frac{1}{\epsilon^2})$.*

By combining [Lemma 5.2](#) and [Lemma 5.3](#), we have $\frac{200Q}{n} \geq \Omega(\frac{1}{\epsilon^2})$; and therefore prove the claim that $Q \geq \Omega(\frac{n}{\epsilon^2})$ in [Theorem 5.1](#).

5.2 Second Lower Bound: $Q \geq \Omega\left(\frac{n \ln(1/\delta)}{\epsilon^2 K}\right)$

Lemma 5.4. *Fix real numbers δ, ϵ such that $0 < \delta, \epsilon \leq 0.01$, and integers K, n such that $K \leq n/2$. Let \mathcal{A} be a deterministic algorithm (i.e., the only randomness comes from the sampling the arms), so that for any set of n Bernoulli arms with means $\theta_1, \theta_2, \dots, \theta_n$,*

- \mathcal{A} makes at most Q samples in expectation;
- with probability at least $1 - \delta$, \mathcal{A} outputs a set T of size K with $\text{val}_K(T) \geq \text{val}_K([n]) - \epsilon$.

Then, we have that $Q \geq \frac{n \ln(1/\delta)}{20000\epsilon^2 K}$.

Now, we provide a sketch of our proof, which generalizes the previous proof for the lower bound when $K = 1$ ([Mannor & Tsitsiklis, 2004](#)). Let $t = \lfloor \frac{n}{K} \rfloor \geq 2$ and we divide the first tK arms into t groups, where the j -th group consists of the arms with the indices in $[(j-1)K + 1, jK]$. We first construct t hypotheses H_1, H_2, \dots, H_t as follows. In H_1 , we let $\theta_i = 1/2 + 4\epsilon$ for arms in the first group and let $\theta_i = 1/2$ for the remaining arms. In H_j , where $2 \leq j \leq t$, we let $\theta_i = 1/2 + 4\epsilon$ when i is in the first group, $\theta_i = 1/2 + 8\epsilon$ when i is in the j -th group, and $\theta_i = 1/2$ otherwise.

For each $j \in [t]$, let $\Pr_{H_j}[\cdot]$ ($\mathbf{E}_{H_j}[\cdot]$ resp.) denote the probability of the event in $[\cdot]$ (the expected value of the random variable in $[\cdot]$ resp.) under the hypothesis H_j . Let \tilde{q}_j be the total number of samples taken from the arms in the j -th group. By an averaging argument, there must exist a group $j_0 \geq 2$ such that $\mathbf{E}_{H_1}[\tilde{q}_{j_0}] \leq \frac{Q}{t-1} \leq \frac{2Q}{t}$.

After fixing j_0 , we focus on the hypothesis H_1 and H_{j_0} . Let $\text{val}_K^{H_j}(T)$ ($\text{val}_K^{H_j}([n])$ resp.) be the $\text{val}_K(T)$ value ($\text{val}_K([n])$ resp.) computed using θ values defined in hypothesis H_j . Note that $\text{val}_K^{H_j}(T)$ (for any j) is always well defined no matter which hypothesis is the true underlying probability measure. Our proof strategy is to assume for contradiction that $Q < \frac{n \ln(1/\delta)}{20000\epsilon^2 K}$ and use the assumption $\Pr_{H_1}[\text{val}_K^{H_1}(T) \geq \text{val}_K^{H_1}([n]) - \epsilon] \geq 1 - \delta$ to first prove that:

$$\Pr_{H_{j_0}}[\text{val}_K^{H_1}(T) \geq \text{val}_K^{H_1}([n]) - \epsilon] \geq \frac{\sqrt{\delta}}{4} \quad (5)$$

To this end, we construct the likelihood ratio between events under the hypothesis H_1 and H_{j_0} . The intuition is that H_1 and H_{j_0} are similar, thus for any sampling outcomes y , the probability that H_1 generates y is close

to H_{j_0} . Since \mathcal{A} is deterministic, the sampling outcomes determine the next action and the final decision. Using this argument, we can show that if the event $\{\text{val}_K^{H_1}(T) \geq \text{val}_K^{H_1}([n]) - \epsilon\}$ happens under H_1 with probability at least $1 - \delta$, it would happen under H_{j_0} with a significant probability (i.e., at least $\frac{\sqrt{\delta}}{4}$).

We further observe that when $\text{val}_K^{H_1}(T) \geq \text{val}_K^{H_1}([n]) - \epsilon$, T must consist of more than $\frac{3}{4}K$ arms from the first group; while when $\text{val}_K^{H_{j_0}}(T) \geq \text{val}_K^{H_{j_0}}([n]) - \epsilon$, T must consist of more than $\frac{3}{4}K$ arms from the j_0 -th group. Therefore the two events are mutually exclusive and we have: $\Pr_{H_{j_0}} \left[\text{val}_K^{H_{j_0}}(T) \geq \text{val}_K^{H_{j_0}}([n]) - \epsilon \right] \leq 1 - \Pr_{H_{j_0}} \left[\text{val}_K^{H_1}(T) \geq \text{val}_K^{H_1}([n]) - \epsilon \right] \leq 1 - \frac{\sqrt{\delta}}{4} \leq 1 - 2\delta$, where the last inequality is because of $\delta < 0.01$. This contradicts the performance guarantees of Algorithm \mathcal{A} and thus we conclude our proof. The details are provided in the appendix.

The proof of Lemma 5.4 can be easily generalized to the case where \mathcal{A} is randomized, which allows us to prove the following stronger lower bound statement. The proof of Theorem 5.5 is relegated to the appendix.

Theorem 5.5. *Fix real numbers δ, ϵ such that $0 < \delta, \epsilon \leq 0.01$, and integers K, n , where $K \leq n/2$. Let \mathcal{A} be a (possibly randomized) algorithm so that for any set of n arms with the mean $\theta_1, \theta_2, \dots, \theta_n$,*

- *\mathcal{A} makes at most Q samples in expectation;*
- *With probability at least $1 - \delta$, \mathcal{A} outputs a set T of size K with $\text{val}_K(T) \geq \text{val}_K([n]) - \epsilon$.*

We have that $Q = \Omega\left(\frac{n \ln(1/\delta)}{\epsilon^2 K}\right)$.

6 Experiments

In this experiment, we assume that arms follow independent Bernoulli distributions with different means. To make a fair comparison, we fix the total budget Q and compare our algorithm (OptMAI) with the uniform sampling strategy and two other state-of-the-art algorithms: SAR (Bubeck et al., 2013) and LUCB (Kalyanakrishnan et al., 2012), in terms of the aggregate regret in (1).

The implementation of our algorithm is slightly different from its description in Section 3. While we choose to present a variant of our implementation only because of its simplicity of exposition, we describe the full details of the differences as follows. It is easy to check that our implementation still meets the theoretical bound proved in Section 4.

First, observe that in OptMAI, Q is an upper bound of the number of samples, while $(1 - \beta^R)Q < Q$ is the actual number of samples used, where R is the total number of rounds run by the algorithm. To fully utilize the budget, we run OptMAI with a parameter slightly greater than Q to ensure that the actual number of samples roughly equals to (but no greater than) Q .

Second, in each round of QE or AR, when computing the empirical mean $\hat{\theta}_i$, our implementation uses all the samples obtained for the i -th arm (i.e. including the samples from previous rounds). This will lead to be better empirical performance especially when the budget is very limited.

Third, in each round of OptMAI, the ratio of the number of samples between two consecutive rounds is set to be $\beta = e^{0.2} \cdot 0.75 \approx 0.91$. In the real implementation, one could treat this quantity as a tuning parameter to make the algorithm more flexible (as long as $\beta \in (0.75, 1)$). In this experiment, we report the results for both $\beta = 0.8$ and $\beta = 0.9$. Based on our experimental results, one could simply set $\beta = 0.8$, which will lead to reasonably good performance under different scenarios. We propose the following strategy to tune β as a future work. In the first stage, we sample each arm for a few times and then use the empirical estimate of θ_i to generate as much simulated data as we want. Then, we choose the best β based on the simulated data. Finally, we apply the carefully tuned β to the real data using the remaining budget.

6.1 Simulated Experiments

In our simulated experiment, the number of total arms is set to $n = 1000$. We vary the total budget $Q = 20n, 50n, 100n$ and $K = 10, 20, \dots, 500$. We use different ways to generate $\{\theta_i\}_{i=1}^n$ and report the comparison results among different algorithms:

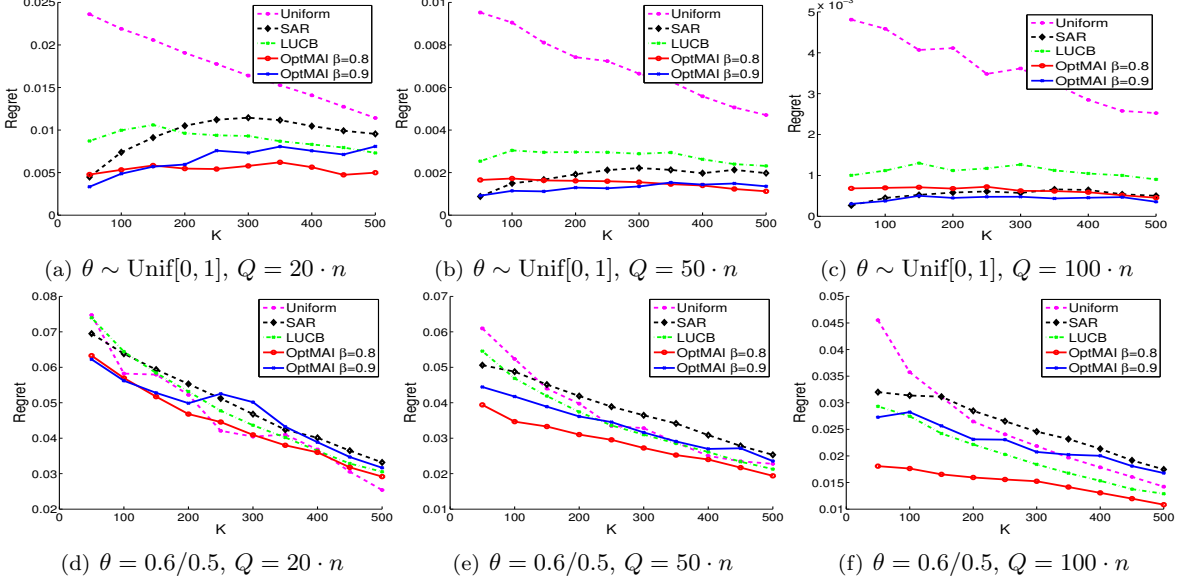


Figure 1: Performance comparison on simulated data.

1. $\theta_i \sim \text{Unif}[0, 1]$: each θ_i is uniformly distributed on $[0, 1]$ (see Figure 1(a) to Figure 1(b)).
2. $\theta_i = 0.5/0.6$: $\theta_i = 0.6$ for $i = 1, \dots, K$ and $\theta_i = 0.5$ for $i = K + 1, \dots, n$. We note that such a two level setting of θ_i is more challenging for the selection of top- K arms (see Figure 1(d) to Figure 1(f)).

It can be seen from Figure 1 that the uniform sampling performs the worst and our method outperforms SAR and LUCB in most of the scenarios. We also observe that when K is large, the setting of $\beta = 0.8$ (red line) outperforms that of $\beta = 0.9$; while for small K , $\beta = 0.9$ (blue line) is a better choice.

We also generate θ_i from the truncated normal distribution and the Beta distribution and have similar observations. The comparison results are presented in the appendix due to space constraints.

6.2 Real RTE Data

We generate θ from a real recognizing textual entailment (RTE) dataset (Section 4.3 in (Snow et al., 2008)). There are 800 task and each task is a sentence pair. Each sentence pair is presented to 10 different workers to acquire binary choices of whether the second hypothesis sentence can be inferred from the first one. There are in total 164 different workers. Since there are true labels of tasks, we set each θ_i for the i -th worker to be his/her labeling accuracy. The histogram of θ_i is presented in Figure 2(a). We vary the total budget $Q = 10n, 20n, 50n$ and K from 10 to 100 and report the comparison of the regret for different approaches in Figure 2(b) to Figure 2(d). As we can see, our method with $\beta = 0.8$ (red line) outperforms other competitors for most of K 's and Q 's. SAR performs the best for $K = 10, Q = 10n$; while our method with $\beta = 0.9$ performs the best for $K = 10$ and $Q = 20n$.

In addition, we would like to highlight an interesting property of our method from the empirical study. As shown in Figure 2(e) and Figure 2(f) with $Q = 10n$ and $K = 20$, the empirical distribution of the number of samples (i.e., tasks) assigned to a worker using SAR is much more skewed than that using our method. This property makes our method particularly suitable for crowdsourcing applications since it will be extremely time-consuming if a single worker is assigned with too many tasks (e.g., gold samples). For example, for SAR, a worker could receive up to 143 tasks (Figure 2(e)) while for our method, a worker receives at most 48 tasks (Figure 2(f)). In crowdsourcing, a single worker will take a long time and soon lose patience when performing nearly 150 testing tasks.

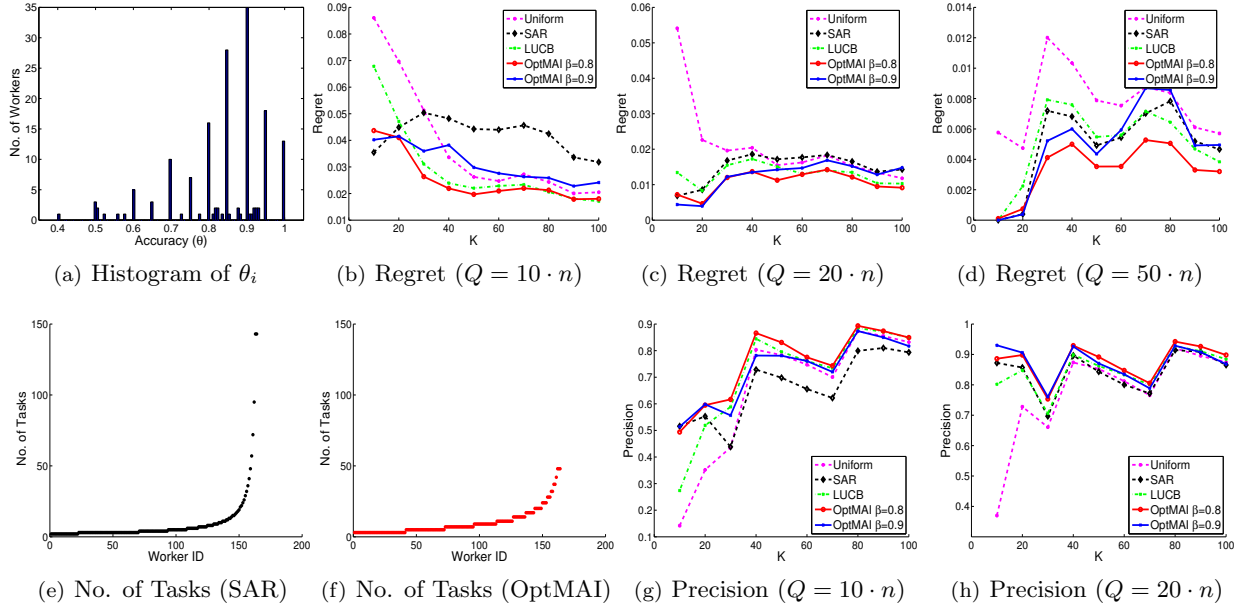


Figure 2: Performance comparison on the RTE data.

In Figure 2(g) and Figure 2(h), we compare different algorithms in terms of the precision, which is defined as the number of arms in T which belong to the set of the top K arms over K , i.e., $\frac{|T \cap [K]|}{K}$. As we can see, our method with $\beta = 0.8$ achieves the highest precision followed by LUCB.

7 Conclusions and Future Work

We study the problem of identifying the (approximate) top K -arms in a stochastic multi-armed bandit game. We propose to use the aggregate regret as the evaluation metric, which fits to the PAC framework. We argue that in many real applications, our metric is more suitable. Our algorithm can identify an ϵ -optimal solution with the sample complexity $O\left(\frac{n}{\epsilon^2} \left(1 + \frac{\ln(1/\delta)}{K}\right)\right)$ with probability at least $1 - \delta$, which matches the lower bound provided in this paper.

There are several directions that we would like to explore in the future. Firstly, our algorithm provides the guarantees for the worst case scenarios and does not depend on the actual reward distributions (i.e., the value of θ_i). In many real data sets, the means of the arms are well separated, and might be easier than the worst case instances (such as the instances constructed in our lower bound proof). Inspired by the work (Audibert et al., 2010; Bubeck et al., 2013; Karnin et al., 2013), our next step is to design new adaptive algorithms and provide more refined distribution dependent upper and lower bounds. Secondly, our lower bound instances are based on Bernoulli bandits. It would be interesting to establish the lower bound for other distributions supported on $[0, 1]$ or even more general distributions, such as sub-Gaussian. Finally, it would be of great interest to test our algorithm on real crowdsourcing platforms and apply it to many other real-world applications.

References

Anthony, Martin and Bartlett, Peter L. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, 1999. 7

- Audibert, Jean-Yves, Bubeck, Sébastien, and Lugosi, Gábor. Regret in online combinatorial optimization. *Mathematics of Operations Research*, 2013. [2](#), [3](#)
- Audibert, J.Y., Bubeck, S., and Munos, R. Best arm identification in multi-armed bandits. In *Proceedings of the 23rd Annual Conference on Learning Theory (COLT)*, 2010. [3](#), [10](#)
- Auer, P., Cesa-Bianchi, N., Freund, Y., and Schapire, R.E. The nonstochastic multiarmed bandit problem. *SIAM J. on Comput.*, 32(1):48–77, 2002a. [3](#)
- Auer, Peter, Cesa-Bianchi, Nicolo, and Fischer, Paul. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47:235–256, 2002b. [3](#)
- Beygelzimer, Alina, Langford, John, Li, Lihong, Reyzin, Lev, and Schapire, Robert E. Contextual bandit algorithms with supervised learning guarantees. In *AISTATS*, 2011. [3](#)
- Bubeck, S. and Cesa-Bianchi, N. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012. [3](#)
- Bubeck, Sebastian, Wang, Tengyao, and Viswanathan, Nitin. Multiple indentifications in multi-armed bandits. In *Proceedings of the 30th International Conference on Machine Learning*, 2013. [1](#), [2](#), [3](#), [8](#), [10](#)
- Bubeck, Sébastien, Munos, Rémi, and Stoltz, Gilles. Pure exploration in multi-armed bandits problems. In *Algorithmic Learning Theory*, pp. 23–37. Springer, 2009. [3](#)
- Chen, Xi, Lin, Qihang, and Zhou, Dengyong. Optimistic knowledge gradient for optimal budget allocation in crowdsourcing. In *ICML*, 2013. [1](#)
- Chernoff, Herman. *Sequential Analysis and Optimal Design*. Society for Industrial and Applied Mathematics (SIAM), 1972. [7](#)
- Even-Dar, Eyal, Mannor, Shie, and Mansour, Yishay. Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *Journal of machine learning research*, 7:1079–1105, 2006. [2](#), [3](#)
- Ho, Chien-Ju, Jabbari, Shahin, , and Vaughan, Jennifer Wortman. Adaptive task assignment for crowdsourced classification. In *ICML*, 2013. [1](#)
- Kalyanakrishnan, Shivaram and Stone, Peter. Efficient selection of multiple bandit arms: Theory and practice. In *Proceedings of International Conference of Machine Learning*, 2010. [1](#), [2](#), [3](#)
- Kalyanakrishnan, Shivaram, Tewari, Ambuj, Auer, Peter, and Stone, Peter. Pac subset selection in stochastic multi-armed bandits. In *Proceedings of the 29th International Conference on Machine Learning (ICML)*, 2012. [2](#), [3](#), [8](#)
- Karger, D. R., Oh, S., and Shah, D. Budget-optimal task allocation for reliable crowdsourcing systems. arXiv:1110.3564v3, 11 2012. [1](#)
- Karnin, Zohar, Koren, Tomer, and Somekh, Oren. Almost optimal exploration in multi-armed bandits. In *International Conference on Machine Learning*, 2013. [2](#), [3](#), [10](#)
- Koenig, Lloyd W. and Law, Averill M. A procedure for selecting a subset of size m containing the l best of k independent normal populations, with applications to simulation. *Communications in statistics. Simulation and computation*, 14:719–734, 1985. [1](#)
- Liu, Qiang, Steyvers, Mark, and Ihler, Alexander. Scoring workers in crowdsourcing: How many control questions are enough? In *NIPS*, 2013. [1](#)

- Mannor, Shie and Tsitsiklis, John N. The sample complexity of exploration in the multi-armed bandit problem. *Journal of Machine Learning Research*, 5:623–648, 2004. [2](#), [3](#), [7](#)
- Raykar, V. C., Yu, S., Zhao, L. H., Valadez, G. H., Florin, C., Bogoni, L., and Moy, L. Learning from crowds. *JMLR*, 11:1297–1322, 2010. [1](#)
- Schmidt, Christian, Branke, Jürgen, and Chick, Stephen E. Integrating techniques from statistical ranking into evolutionary algorithms. 2006. [1](#)
- Snow, R., Connor, B. O., Jurafsky, D., and Ng., A. Y. Cheap and fast - but is it good? evaluating non-expert annotations for natural language tasks. In *EMNLP*, 2008. [10](#)
- Thompson, W.R. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25:285–294, 1933. [1](#)
- Welinder, P., Branson, S., Belongie, S., and Perona, P. The multidimensional wisdom of crowds. In *NIPS*, 2010. [1](#)
- Zhou, D., Basu, S., Mao, Y., and Platt, J. Learning from the wisdom of crowds by minimax conditional entropy. In *NIPS*, 2012. [1](#)

A Proof of the Correctness of the QE Algorithm ([Lemma 4.1](#) in the Main Text)

Let $p = \theta_{\text{ind}_{|S|/2}(S)}$ be the median of the θ in S and $\tau = \min_{i \in V}(\hat{\theta}_i)$ be the minimum empirical mean for the selected arms in V . For each arm i among the top K arms in S , we define the random variable $X_i = \mathbf{1}\{\hat{\theta}_{\text{ind}_i(S)} < p + \frac{\epsilon}{2}\}$ and $X = \mathbf{E}_{i \in [K]}(\theta_{\text{ind}_i(S)} - p)X_i$. Note that we use $\mathbf{E}_{i \in [K]}$ to denote the average operator $\frac{1}{K} \sum_{i=1}^K$ and \mathbf{E} (without subscript) to denote the expectation operator with respect to the randomness of the arms. We further define two events $\mathcal{E}_1 = \{X \leq \epsilon\}$ and $\mathcal{E}_2 = \{\tau < p + \frac{\epsilon}{2}\}$. Our first claim is that \mathcal{E}_1 and \mathcal{E}_2 together imply our conclusion $\text{val}_K(V) \geq \text{val}_K(S) - \epsilon$.

Claim A.1. \mathcal{E}_1 and \mathcal{E}_2 imply that $\text{val}_K(V) \geq \text{val}_K(S) - \epsilon$.

Proof. For each arm i that is among the best K arms of S , if $X_i = 0$, by \mathcal{E}_2 , we have

$$\hat{\theta}_{\text{ind}_i(S)} \geq p + \frac{\epsilon}{2} > \tau = \min_{i \in V}(\hat{\theta}_i),$$

and hence $i \in V$. Therefore, i is also one of the best K arms of V . On the other hand, since there are $|S|/2 \geq K + |S|/4$ arms with means greater or equal to p , after removing $|S|/4$ of them, there are still at least K such arms. Therefore we know that the best K arms of V all have means greater than or equal to p . Now, we can see that

$$\begin{aligned} \text{val}_K(V) &= \mathbf{E}_{i \in [K]} \theta_{\text{ind}_i(V)} \geq \mathbf{E}_{i \in [K]} ((1 - X_i)\theta_{\text{ind}_S(i)} + X_i p) \\ &= \mathbf{E}_{i \in [K]} \theta_{\text{ind}_i(S)} - \mathbf{E}_{i \in [K]} (\theta_{\text{ind}_i(S)} - p)X_i \geq \text{val}_K(S) - \epsilon \end{aligned}$$

where the last inequality is due to \mathcal{E}_1 . □

To prove the correctness of the QE algorithm, it suffices to lower bound the probabilities of the events \mathcal{E}_1 and \mathcal{E}_2 .

Claim A.2. $\Pr[\mathcal{E}_1] \geq 1 - \frac{\delta}{2}$.

Before proceeding to the proof of claim, we first prove a concentration inequality.

Proposition A.3. Let $X_i (1 \leq i \leq K)$ be independent random variables. Each X_i takes value a_i ($a_i \geq 0$) with probability at most $\exp(-a_i^2 t)$ and 0 otherwise, where $t \geq 0$. Let $X = \mathbf{E}_{i \in [K]} X_i$. For every $\epsilon > 0$, when $t \geq \frac{2}{\epsilon^2}$, we have

$$\Pr[X > \epsilon] < \exp\left(-\frac{\epsilon^2 t K}{2}\right).$$

To prove [Proposition A.3](#), we first provide a technical lemma.

Lemma A.4. Let Z be a random variable that takes value a ($a \geq 0$) with probability at most $\exp(-a^2 t)$ and 0 otherwise. For every $\epsilon > 0$, when $t \geq \frac{2}{\epsilon^2}$, we have

$$\mathbf{E} e^{\epsilon t Z} \leq \exp\left(\frac{\epsilon^2 t}{2}\right).$$

Proof. By the definition of Z , we know that

$$\mathbf{E} e^{\epsilon t Z} \leq \exp(\epsilon a t - a^2 t) + 1 = \exp(a(\epsilon - a)t) + 1 \leq \exp\left(\frac{\epsilon^2 t}{4}\right) + 1.$$

When $\epsilon^2 t \geq 2$, we have $\exp\left(\frac{\epsilon^2 t}{2}\right) - \exp\left(\frac{\epsilon^2 t}{4}\right) > 1.06 > 1$ and hence

$$\mathbf{E} e^{\epsilon t Z} \leq \exp\left(\frac{\epsilon^2 t}{2}\right).$$

□

Proof of [Proposition A.3](#). Observe that

$$\begin{aligned} \Pr[X > \epsilon] &= \Pr\left[\sum_{i=1}^K X_i > \epsilon K\right] = \Pr\left[\sum_{i=1}^K \epsilon t X_i > \epsilon^2 t K\right] = \Pr\left[\exp\left(\sum_{i=1}^K \epsilon t X_i\right) > \exp(\epsilon^2 t K)\right] \\ &\leq \mathbf{E} \frac{\exp\left(\sum_{i=1}^K \epsilon t X_i\right)}{\exp(\epsilon^2 t K)} = \frac{\prod_{i=1}^K \mathbf{E} \exp(\epsilon t X_i)}{\exp(\epsilon^2 t K)} \end{aligned}$$

Now apply [Lemma A.4](#), and we get

$$\Pr[X > \epsilon] \leq \frac{\prod_{i=1}^K \exp\left(\frac{\epsilon^2 t}{2}\right)}{\exp(\epsilon^2 t K)} = \exp\left(-\frac{\epsilon^2 t K}{2}\right).$$

□

With the concentration inequality in place, we prove the [Claim A.2](#).

Proof of [Claim A.2](#). Recall that $Q_0 = \frac{Q}{|S|}$ is the number of samples taken from each arm in S . Since

$$\epsilon = \sqrt{\frac{1}{Q_0} \left(10 + \frac{4 \ln(2/\delta)}{K}\right)},$$

we trivially have that $\epsilon \geq \max\left\{\sqrt{\frac{10}{Q_0}}, \sqrt{\frac{4 \ln(2/\delta)}{Q_0 K}}\right\}$.

For each $i \in [K]$, let $\eta_i = \max\{\theta_{\text{ind}_i(S)} - p - \frac{\epsilon}{2}, 0\}$; and let $Y_i = \eta_i X_i$. By Hoeffding's inequality, we have

$$\Pr[Y_i = \eta_i] = \Pr[X_i = 1] = \Pr\left[\hat{\theta}_{\text{ind}_i(S)} < p + \frac{\epsilon}{2}\right] \leq \exp\left(-2 \left(\theta_{\text{ind}_i(S)} - p - \frac{\epsilon}{2}\right)^2 \cdot Q_0\right) \leq \exp(-\eta_i^2 \cdot 2Q_0).$$

Now apply [Proposition A.3](#) on Y_i 's, we get

$$\Pr \left[\mathbf{E}_{i \in [K]} Y_i > \frac{\epsilon}{2} \right] \leq \exp \left(-\frac{\epsilon^2 Q_0 K}{4} \right) \leq \frac{\delta}{2},$$

where the last inequality is because $\epsilon \geq \sqrt{\frac{4 \ln(2/\delta)}{Q_0 K}}$. Observe that $Y_i \geq (\theta_{\text{ind}_S(i)} - p)X_i - \frac{\epsilon}{2}$ for all $i \in [K]$. Therefore, with probability at least $(1 - \frac{\delta}{2})$, we have

$$X = \mathbf{E}_{i \in [K]} (\theta_{\text{ind}_i(S)} - p)X_i \leq \mathbf{E}_{i \in [K]} Y_i + \frac{\epsilon}{2} \leq \epsilon.$$

□

Claim A.5. $\Pr[\mathcal{E}_2] \geq 1 - \frac{\delta}{2}$.

Proof. Since the set V contains at least $|S|/4$ arms with the indices between $[|S|/2, |S|]$, to prove $\tau < p + \frac{\epsilon}{2}$, we only need to show that there are at most $|S|/4$ indices i with $i \in [|S|/2, |S|]$ such that $\hat{\theta}_{\text{ind}_i(S)} \geq p + \frac{\epsilon}{2}$ (with high probability).

For each $i \in [|S|/2, |S|]$, we have $\theta_{\text{ind}_i(S)} \leq p$. By Hoeffding's inequality, we have

$$\Pr \left[\hat{\theta}_{\text{ind}_i(S)} \geq p + \frac{\epsilon}{2} \right] \leq \Pr \left[\hat{\theta}_{\text{ind}_i(S)} \geq \theta_{\text{ind}_i(S)} + \frac{\epsilon}{2} \right] \leq \exp \left(-\frac{\epsilon^2}{2} \cdot Q_0 \right).$$

Define the indicator random variable $Z_i = \mathbf{1}\{\hat{\theta}_{\text{ind}_i(S)} \geq p + \frac{\epsilon}{2}\}$ with its mean $\mathbf{E}(Z_i) < \exp \left(-\frac{\epsilon^2}{2} \cdot Q_0 \right)$. Let $\mu = \max_{i \in [|S|/2, |S|]} \mathbf{E}(Z_i)$. We have $\mu < \exp \left(-\frac{\epsilon^2}{2} \cdot Q_0 \right)$. Then, by Chernoff bound, we have

$$\begin{aligned} \Pr \left[\sum_{i=|S|/2}^{|S|} Z_i > \frac{|S|}{4} \right] &\leq \left(\left(\frac{\mu}{1/2} \right)^{1/2} \left(\frac{1-\mu}{1/2} \right)^{1/2} \right)^{|S|/2} \leq \left(\sqrt{2\mu} \cdot \sqrt{2} \right)^{|S|/2} \\ &\leq \exp \left(\frac{|S|}{2} \left(\ln(2) - \frac{\epsilon^2}{2} \cdot Q_0 \right) \right) \leq \exp \left(-\frac{|S|}{2} \cdot \frac{\epsilon^2}{4} \cdot Q_0 \right) \leq \exp \left(-\frac{\epsilon^2 K}{2} \cdot Q_0 \right) \leq \frac{\delta}{2}. \end{aligned}$$

where the third last inequality is because of $\epsilon > \sqrt{\frac{10}{Q_0}}$, the second inequality uses the assumption that $|S|/4 \geq K$ and the last inequality is because

$$\epsilon \geq \sqrt{\frac{4 \ln(2/\delta)}{Q_0 K}} \geq \sqrt{\frac{2 \ln(2/\delta)}{Q_0 K}}.$$

□

Proof of [Lemma 4.1](#). By [Claim A.2](#), [Claim A.5](#), and a union bound, we have $\Pr[\mathcal{E}_1 \text{ and } \mathcal{E}_2] \geq 1 - \delta$. By [Claim A.1](#), we have $\Pr[\text{val}_K(T) \geq \text{val}_K(S) - \epsilon] \geq \Pr[\mathcal{E}_1 \text{ and } \mathcal{E}_2]$; and the lemma follows. □

B Proof of the Correctness of the AR Algorithm ([Lemma 4.2](#) in the Main Text)

Proof. Recall that $Q_0 = \frac{Q}{|S|}$ is the number of samples for each arm in S . Also recall that $K' = K - |T|$.

Let $U_1 = T' \setminus T$ denote the set of arms we added to T' in this round. Let $U_2 = S \setminus (S' \cup U_1)$ be the set of arms we discarded in this round. Let U_1^* be the set of $|U_1|$ arms in S with largest θ_i 's; let U_2^* be the set of $|U_2|$ arms in S with smallest θ_i 's. One can check that

$$\begin{aligned}
& (\text{tot}_{K-|T'|}(S') + \text{tot}_{|T'|}(T')) - (\text{tot}_{K-|T|}(S) + \text{tot}_{|T|}(T)) \\
&= \text{tot}_{K-|T'|}(S') - \text{tot}_{K-|T|}(S) + \sum_{i \in U_1} \theta_i \\
&= \text{tot}_{K-|T'|}(S') - \text{tot}_{K-|T'|}(S \setminus U_1^*) + \left(\sum_{i \in U_1} \theta_i - \sum_{i \in U_1^*} \theta_i \right) \\
&\geq \text{tot}_{K-|T'|}(S') - \text{tot}_{K-|T'|}(S \setminus U_1) + \left(\sum_{i \in U_1} \theta_i - \sum_{i \in U_1^*} \theta_i \right). \tag{6}
\end{aligned}$$

Let \tilde{U}_2 be the $|U_2|$ arms with the smallest means in $S \setminus U_1$. By definition we have 1) $|\tilde{U}_2| = |U_2^*|$; 2) $\tilde{U}_2 \cap U_1 = U_2 \cap U_1 = \emptyset$; 3) $\sum_{i \in \tilde{U}_2} \theta_i \geq \sum_{i \in U_2^*} \theta_i$.

Since $|U_1| + |U_2| + (K - |T'|) \leq |S|$, the $(K - |T'|)$ arms with largest means in $S \setminus U_1$ do not intersect with the $|U_2|$ arms with smallest means in $S \setminus U_1$ (namely \tilde{U}_2). Therefore we have

$$\text{tot}_{K-|T'|}(S \setminus U_1) = \text{tot}_{K-|T'|}((S \setminus U_1) \setminus \tilde{U}_2). \tag{7}$$

On the other hand, for every set W of arms, define $\text{tot}_t^{\min}(W)$ to be the sum of the t smallest means among the arms in W . Let $t = |S| - |U_1| - |U_2| - (K - |T'|)$. Since \tilde{U}_2 consists of the arms with the smallest means in $S \setminus U_1$, we have

$$\text{tot}_t^{\min}((S \setminus U_1) \setminus \tilde{U}_2) \geq \text{tot}_t^{\min}((S \setminus U_1) \setminus U_2).$$

Together with

$$\begin{aligned}
\text{tot}_t^{\min}((S \setminus U_1) \setminus \tilde{U}_2) &= \left(\sum_{i \in (S \setminus U_1) \setminus \tilde{U}_2} \theta_i \right) - \text{tot}_{K-|T'|}((S \setminus U_1) \setminus \tilde{U}_2), \\
\text{tot}_t^{\min}((S \setminus U_1) \setminus U_2) &= \left(\sum_{i \in (S \setminus U_1) \setminus U_2} \theta_i \right) - \text{tot}_{K-|T'|}((S \setminus U_1) \setminus U_2),
\end{aligned}$$

we have

$$\left(\sum_{i \in (S \setminus U_1) \setminus \tilde{U}_2} \theta_i \right) - \text{tot}_{K-|T'|}((S \setminus U_1) \setminus \tilde{U}_2) \geq \left(\sum_{i \in (S \setminus U_1) \setminus U_2} \theta_i \right) - \text{tot}_{K-|T'|}((S \setminus U_1) \setminus U_2),$$

i.e.

$$\text{tot}_{K-|T'|}((S \setminus U_1) \setminus U_2) - \text{tot}_{K-|T'|}((S \setminus U_1) \setminus \tilde{U}_2) \geq \sum_{i \in \tilde{U}_2} \theta_i - \sum_{i \in U_2} \theta_i. \tag{8}$$

By combining (6), (7) and (8), and the observations that $S' = (S \setminus U_1) \setminus U_2$ and $\sum_{i \in \tilde{U}_2} \theta_i \geq \sum_{i \in U_2^*} \theta_i$, we have

$$(6) \geq \left(\sum_{i \in U_2^*} \theta_i - \sum_{i \in U_2} \theta_i \right) - \left(\sum_{i \in U_1} \theta_i - \sum_{i \in U_1^*} \theta_i \right). \tag{9}$$

For every $t \leq K$, for every set $U \subseteq S$ of t arms (i.e. $|U| = t$), by Hoeffding's inequality, we have

$$\Pr \left[\left| \sum_{i \in U} \hat{\theta}_i - \sum_{i \in U} \theta_i \right| > \frac{\epsilon K}{4} \right] \leq 2 \exp \left(-\frac{\epsilon^2}{8} \cdot Q_0 \frac{K^2}{t} \right) \leq 2 \exp \left(-\frac{\epsilon^2}{8} \cdot Q_0 K \right).$$

By a union bound, we have

$$\begin{aligned} \Pr \left[\forall U \subseteq S, |U| \leq K : \left| \sum_{i \in U} \hat{\theta}_i - \sum_{i \in U} \theta_i \right| \leq \frac{\epsilon K}{4} \right] &\geq 1 - 2 \cdot 2^{|S|} \exp \left(-\frac{\epsilon^2}{8} \cdot Q_0 K \right) \\ &\geq 1 - 2 \exp \left(|S| - \frac{\epsilon^2}{8} \cdot Q_0 K \right) \geq 1 - \delta, \end{aligned}$$

where we used $|S| < 4K$ and $\epsilon \geq \sqrt{\frac{1}{Q_0} \left(4 + \frac{\ln(2/\delta)}{K} \right)}$.

Thus, with probability at least $1 - \delta$, we have

$$\begin{aligned} \left| \sum_{i \in U_1} \hat{\theta}_i - \sum_{i \in U_1} \theta_i \right| &\leq \frac{\epsilon K}{4}, & \left| \sum_{i \in U_1^*} \hat{\theta}_i - \sum_{i \in U_1^*} \theta_i \right| &\leq \frac{\epsilon K}{4}, \\ \left| \sum_{i \in U_2} \hat{\theta}_i - \sum_{i \in U_2} \theta_i \right| &\leq \frac{\epsilon K}{4}, & \left| \sum_{i \in U_2^*} \hat{\theta}_i - \sum_{i \in U_2^*} \theta_i \right| &\leq \frac{\epsilon K}{4}. \end{aligned}$$

Therefore we have

$$\sum_{i \in U_1} \theta_i \geq \sum_{i \in U_1} \hat{\theta}_i - \frac{\epsilon K}{4} \geq \sum_{i \in U_1^*} \hat{\theta}_i - \frac{\epsilon K}{4} \geq \sum_{i \in U_1^*} \theta_i - \frac{\epsilon K}{2}, \quad (10)$$

and

$$\sum_{i \in U_2} \theta_i \leq \sum_{i \in U_2} \hat{\theta}_i - \frac{\epsilon K}{4} \leq \sum_{i \in U_2^*} \hat{\theta}_i - \frac{\epsilon K}{4} \leq \sum_{i \in U_2^*} \theta_i - \frac{\epsilon K}{2}. \quad (11)$$

Combining (9), (10) and (11), we get (6) $\geq -\epsilon K$, which concludes the proof. \square

C Proof of the Main Theorem (Theorem 4.3 in the Main Text)

Proof. Recall r is the counter of the number of iterations in Algorithm 1. Let r_0 be the first r such that we have $|S_r| < 4K$. Let r_1 be the final value of r .

For $r < r_0$, by Lemma 4.1, with probability $(1 - e^{-.1r}(1 - e^{-.1})\delta)$ we have

$$\text{val}_K(S_{r+1}) \geq \text{val}_K(S_r) - \sqrt{\frac{(\frac{3}{4})rn}{(1-\beta)\beta^r Q} \left(10 + \frac{4 \ln(2/(e^{-.1r}(1 - e^{-.1})\delta))}{K} \right)}.$$

By a union bound, with probability $1 - \sum_{r=0}^{r_0-1} e^{-.1r}(1 - e^{-.1})\delta$ we have

$$\text{val}_K(S_{r_0}) \geq \text{val}_K([n]) - \sum_{r=0}^{r_0-1} \sqrt{\frac{(\frac{3}{4})rn}{(1-\beta)\beta^r Q} \left(10 + \frac{4 \ln(2/(e^{-.1r}(1 - e^{-.1})\delta))}{K} \right)}. \quad (12)$$

For $r : r_0 \leq r < r_1$, by [Lemma 4.2](#), with probability $(1 - e^{-.1r}(1 - e^{-.1}))$ we have

$$\begin{aligned} & \left(\text{tot}_{K-|T_{r+1}|}(S_{r+1}) + \text{tot}_{|T_{r+1}|}(T_{r+1}) \right) - \left(\text{tot}_{K-|T_r|}(S_r) + \text{tot}_{|T_r|}(T_r) \right) \\ & \geq K \cdot \sqrt{\frac{(\frac{3}{4})^r n}{(1-\beta)\beta^r Q} \left(4 + \frac{\ln(2/(e^{-.1r}(1 - e^{-.1})\delta))}{K} \right)} \\ & \geq K \cdot \sqrt{\frac{(\frac{3}{4})^r n}{(1-\beta)\beta^r Q} \left(10 + \frac{4 \ln(2/(e^{-.1r}(1 - e^{-.1})\delta))}{K} \right)}. \end{aligned}$$

Since T_{r_1} has exactly K elements and $T_{r_0} = \emptyset$, by a union bound, with probability $(1 - \sum_{r=r_0}^{r_1-1} e^{-.1r}(1 - e^{-.1})\delta)$ we have

$$\text{val}_K(T_{r_1}) \geq \text{val}_K(S_{r_0}) - \sum_{r=r_0}^{r_1-1} \sqrt{\frac{(\frac{3}{4})^r n}{(1-\beta)\beta^r Q} \left(10 + \frac{4 \ln(2/(e^{-.1r}(1 - e^{-.1})\delta))}{K} \right)}. \quad (13)$$

Now by a union bound over (12) and (13), we have that with probability $1 - \sum_{r=0}^{r_1-1} e^{-.1r}(1 - e^{-.1})\delta \geq 1 - \delta$,

$$\begin{aligned} \text{val}_K([n]) - \text{val}_K(T_{r_1}) & \leq \sum_{r=0}^{r_1-1} \sqrt{\frac{(\frac{3}{4})^r n}{(1-\beta)\beta^r Q} \left(10 + \frac{4 \ln(2/(e^{-.1r}(1 - e^{-.1})\delta))}{K} \right)} \\ & \leq \sum_{r=0}^{r_1-1} \sqrt{\frac{e^{-.2r} n}{.08Q} \left(10 + \frac{4 \ln(1/\delta) + 4 \cdot 3.1 + 4 \cdot .1r}{K} \right)} \\ & \leq \sqrt{\frac{n}{.08Q}} \sum_{r=0}^{r_1-1} \sqrt{e^{-.2r} \left(10 + \frac{4 \ln(1/\delta) + 16}{K} \right)} + \sqrt{\frac{n}{.08QK}} \sum_{r=0}^{r_1-1} \sqrt{.4r e^{-.2r}} \\ & \leq 11 \cdot \sqrt{\frac{n}{.08Q}} \sqrt{10 + \frac{4 \ln(1/\delta) + 16}{K}} + 18 \cdot \sqrt{\frac{n}{.08QK}} \\ & \leq \sqrt{\frac{n}{Q}} \left(44 \cdot \sqrt{10 + \frac{4 \ln(1/\delta) + 16}{K}} + \frac{72}{\sqrt{K}} \right). \end{aligned}$$

□

D Lower bounds

D.1 Proof of the First Lower Bound ([Lemma 5.2](#) in the Main Text)

Proof. Given an algorithm \mathcal{A} stated in [Theorem 5.1](#), we construct the algorithm \mathcal{B} as follows:

- Choose a random subset $S \subseteq [n]$ such that $|S| = K$ and then choose a random element $j \in S$.
- For each $i \in [n], i \neq j$, let $\theta_i = \frac{1}{2} + 4\epsilon$ if $i \in S$, let $\theta_i = \frac{1}{2}$ otherwise.
- Simulate \mathcal{A} as follows: whenever \mathcal{A} samples the i -th arm, if $i = j$, we sample the Bernoulli arm X ; otherwise we sample an arm with the mean θ_i .
- If the arm X is sampled by less than $\frac{200Q}{n}$ times and \mathcal{A} returns a set T such that $j \notin T$, we decide that X has the mean $\frac{1}{2}$; otherwise we decide that X has the mean $\frac{1}{2} + 4\epsilon$.

We note that the number of trials of \mathcal{B} adds one whenever X is sampled. Since \mathcal{B} stops and output the mean $\frac{1}{2} + 4\epsilon$ if the number of trials on X reaches $\frac{200Q}{n}$, \mathcal{B} makes at most $\frac{200Q}{n}$ trials. Now, we further prove

the correctness of \mathcal{B} , i.e., it can correctly output the mean of X with the probability at least 0.51; and thus conclude the proof of [Lemma 5.2](#).

We first show that when the Bernoulli arm X has the mean (i.e., bias) $\frac{1}{2}$, \mathcal{B} decides correctly with probability at least .51. Assuming that X has the mean $\frac{1}{2}$, among the n arms in the algorithm \mathcal{A} , the arms in $S \setminus \{j\}$ have the mean $\frac{1}{2} + 4\epsilon$, while others have the mean $\frac{1}{2}$. For each $i \in [n]$, let the random variable q_i be the number of trials of the i -th arm performed by \mathcal{A} . We have

$$\sum_{i \in [n]} \mathbf{E}[q_i] \leq Q.$$

Let the random variable q_X be the number of trials that the arm X is performed and $S' = S \setminus \{j\}$. Observe that when conditioned on S' , the means (θ_i 's) for \mathcal{A} are constants and j is uniformly distributed among $[n] \setminus S'$. We have

$$\mathbf{E}[q_X] = \mathbf{E}_{S'} \mathbf{E}[q_X | S'] = \mathbf{E}_{S'} \left[\frac{1}{n - K + 1} \sum_{i \in [n] \setminus S'} \mathbf{E}[q_i] \right] \leq \frac{2}{n} \sum_{i \in [n]} \mathbf{E}[q_i] \leq \frac{2Q}{n},$$

where we used assumption that $K \leq n/2$. Therefore, by Markov's inequality,

$$\Pr \left[q_X \geq \frac{200Q}{n} \right] < .01.$$

Let T be the output of the algorithm \mathcal{A} . Since $\text{val}_K([n]) = \frac{1}{2} + 4\epsilon \cdot (1 - \frac{1}{K})$, and $\text{val}_K(T) = \frac{1}{2} + 4\epsilon \cdot \frac{|(S \setminus \{j\}) \cap T|}{K}$, when $\text{val}_K(T) \geq \text{val}_K([n]) - \epsilon$, we have

$$\frac{1}{2} + 4\epsilon \cdot \frac{|(S \setminus \{j\}) \cap T|}{K} \geq \frac{1}{2} + 4\epsilon \cdot \left(1 - \frac{1}{K}\right) - \epsilon \Rightarrow |(S \setminus \{j\}) \cap T| \geq \frac{3}{4}K - 1$$

By the correctness property of \mathcal{A} , for every S' , we have

$$\Pr_T \left[|S' \cap T| \geq \frac{3}{4}K - 1 \right] \geq \Pr_T[\text{val}_K(T) \geq \text{val}_K([n]) - \epsilon] \geq .8.$$

Conditioned on S' , j is uniformly distributed among $[n] \setminus S'$ and is independent from T . Therefore, we have

$$\Pr[j \in T] = \mathbf{E}_{S'} \Pr[j \in T | S'] = \mathbf{E}_{S', T} \left[\frac{|([n] \setminus S') \cap T|}{|[n] \setminus S'|} \right] \leq .2 \cdot 1 + .8 \cdot \frac{\frac{1}{4}K + 1}{n - K + 1} \leq .2 + .8 \cdot (.25 + .1) = .48,$$

where in the last inequality, we used the assumption that $10 \leq K \leq n/2$.

Therefore, when X has the mean $\frac{1}{2}$, we have

$$\Pr \left[\mathcal{A} \text{ decides that } X \text{ has the mean } \frac{1}{2} \right] = \Pr \left[j \notin T \text{ and } q_X \leq \frac{200Q}{n} \right] \geq .52 - .01 = .51.$$

Now we assume that X has the mean $\frac{1}{2} + 4\epsilon$. Among the n arms, the arms in S have the mean $\frac{1}{2} + 4\epsilon$, while others have the mean $\frac{1}{2}$. Again, let T be the output of the algorithm \mathcal{A} . Since with probability at least .8, we have that $\text{val}_K(T) \geq \text{val}_K([n]) - \epsilon$, we have

$$\Pr \left[|S \cap T| \geq \frac{3}{4}K \right] \geq .8.$$

Since j is a random variable uniformly sampled from S , and conditioned on S , j is independent from T , we have

$$\Pr[j \in T] = \mathbf{E}_S \Pr[j \in T | S] = \mathbf{E}_{S, T} \left[\frac{|S \cap T|}{|S|} \right] \geq .8 \cdot \frac{3}{4} \geq .6.$$

In sum, assuming that X has the mean $\frac{1}{2} + 4\epsilon$, we have that

$$\Pr \left[\mathcal{A} \text{ decides that } X \text{ has the mean } \frac{1}{2} + 4\epsilon \right] \geq \Pr[j \in T] \geq .6 > .51.$$

□

D.2 Proof of the Second Lower Bound (Lemma 5.4 and Theorem 5.5 in the Main Text)

Proof of Lemma 5.4. Let $t = \lfloor \frac{n}{K} \rfloor \geq 2$ and we divide the first tK arms into t groups. The j -th group consists of the arms with the index in $[(j-1)K+1, jK]$ for $j \in [t]$. We first construct t hypotheses H_1, H_2, \dots, H_t as follows. In H_1 , we let $\theta_i = 1/2 + 4\epsilon$ for arms in the first group and let $\theta_i = 1/2$ for the remaining arms. In H_j , where $2 \leq j \leq t$, we let $\theta_i = 1/2 + 4\epsilon$ when i is in the first group, $\theta_i = 1/2 + 8\epsilon$ when i is in the j -th group, and $\theta_i = 1/2$ otherwise. For each $j \in [t]$, let $\Pr_{H_j}[\cdot]$ denote the probability of the event in $[\cdot]$ under the hypothesis H_j and $\mathbf{E}_{H_j}[\cdot]$ the expected value of the random variable in $[\cdot]$ under the hypothesis H_j .

For each arm $i \in [n]$, let the random variable q_i be the number of times that \mathcal{A} samples the i -th arm before termination. For each $j \in [t]$, let the random variable $\tilde{q}_j = \sum_{i=(j-1)K+1}^{jK} q_i$ be the total number of trials of the arms in the j -th group. We know that

$$\mathbf{E}_{H_1} \left[\sum_{j=2}^t \tilde{q}_j \right] \leq \mathbf{E}_{H_1} \left[\sum_{i \in [n]} q_i \right] \leq Q.$$

By an averaging argument, there exists $j_0 : 2 \leq j_0 \leq t$ such that

$$\mathbf{E}_{H_1}[\tilde{q}_{j_0}] \leq \frac{Q}{t-1} \leq \frac{2Q}{t}.$$

Using Markov's inequality and letting $Q_0 = \frac{8Q}{t}$, we have

$$\Pr_{H_1}[\tilde{q}_{j_0} \geq Q_0] \leq \frac{\mathbf{E}_{H_1}[\tilde{q}_{j_0}]}{Q_0} \leq \frac{1}{4},$$

and hence,

$$\Pr_{H_1}[\tilde{q}_{j_0} \leq Q_0] \geq \frac{3}{4}. \quad (14)$$

Now we only focus on the hypotheses H_1 and H_{j_0} . Let $\text{val}_K^{H_j}(T)$ ($\text{val}_K^{H_j}([n]$) resp.) be $\text{val}_K(T)$ ($\text{val}_K([n]$) resp.) under the hypothesis H_j . Our proof strategy is to assume for contradiction that $Q < \frac{n \ln(1/\delta)}{20000\epsilon^2 K}$ (i.e., $Q_0 < \frac{\ln(1/\delta)}{1250\epsilon^2}$) and use the assumption

$$\Pr_{H_1}[\text{val}_K^{H_1}(T) \geq \text{val}_K^{H_1}([n]) - \epsilon] \geq 1 - \delta \quad (15)$$

to first prove that:

$$\Pr_{H_{j_0}}[\text{val}_K^{H_1}(T) \geq \text{val}_K^{H_1}([n]) - \epsilon] \geq \frac{\sqrt{\delta}}{4}. \quad (16)$$

We further observe that when $\text{val}_K^{H_1}(T) \geq \text{val}_K^{H_1}([n]) - \epsilon$, T must consist of more than $\frac{3}{4}K$ arms from the first group; while when $\text{val}_K^{H_{j_0}}(T) \geq \text{val}_K^{H_{j_0}}([n]) - \epsilon$, T must consist of more than $\frac{3}{4}K$ arms from the j_0 -th group. Therefore the two events are mutually exclusive and we have:

$$\begin{aligned} & \Pr_{H_{j_0}} \left[\text{val}_K^{H_{j_0}}(T) \geq \text{val}_K^{H_{j_0}}([n]) - \epsilon \right] \\ & \leq 1 - \Pr_{H_{j_0}} \left[\text{val}_K^{H_1}(T) \geq \text{val}_K^{H_1}([n]) - \epsilon \right] \\ & \leq 1 - \frac{\sqrt{\delta}}{4} \leq 1 - 2\delta \end{aligned} \quad (17)$$

where the last inequality is because of $\delta < .01$. This contradicts the performance guarantees of the Algorithm \mathcal{A} and thus we conclude our proof.

Therefore, the remaining task is to prove the correctness of Eq. (16). We first define a sequence of random variables $Z_0, Z_1, Z_2, \dots, Z_{Q_0}$ where $Z_0 = 0$. For each $i \in [Q_0]$, if the i -th trial of the j_0 -th group by \mathcal{A} results in 1, let $Z_i = Z_{i-1} + 1$; if the result is 0, let $Z_i = Z_{i-1} - 1$; if \mathcal{A} terminates before the i -th trial of the j_0 -th group, let $Z_i = Z_{i-1}$. Under hypothesis H_0 , the sequence $\{Z_0, Z_1, Z_2, \dots, Z_{Q_0}\}$ forms a martingale since arms in the j_0 -th group are independent zero-mean random variables. Therefore, by Azuma-Hoeffding's inequality, we have

$$\Pr_{H_1} \left[|Z_{Q_0}| \leq \sqrt{5Q_0} \right] > 1 - 2 \exp \left(-\frac{(\sqrt{5Q_0})^2}{2Q_0} \right) > \frac{3}{4}. \quad (18)$$

By a union bound over (14), (15) and (18), we have

$$\Pr_{H_1} \left[\text{val}_K^{H_1}(T) \geq \text{val}_K^{H_1}([n]) - \epsilon \text{ and } \tilde{q}_{j_0} \leq Q_0 \text{ and } |Z_{Q_0}| \leq \sqrt{5Q_0} \right] \geq 1 - \delta - \frac{1}{4} - \frac{1}{4} \geq \frac{1}{4}. \quad (19)$$

We call a string $y = ((i_1, b_1), (i_2, b_2), \dots, (i_{Q'}, b_{Q'}))$ a transcript for \mathcal{A} if the r -th trial ($1 \leq r \leq Q'$) performed by \mathcal{A} is the i_r -th arm and the result is $b_r \in \{0, 1\}$; and \mathcal{A} uses exactly Q' trials. Let \mathcal{Y} be the set of transcripts for \mathcal{A} . For each $y \in \mathcal{Y}$,

- Let $u_0^i(y)$ be the number of $(i, 0)$ pairs in y and $u_1^i(y)$ be the number of $(i, 1)$ pairs in y ;
- Let $q_i(y) = u_0^i(y) + u_1^i(y)$ be the number of times to sample the i -th arm in y ;
- For all $j \in [t]$, let $\tilde{q}_j(y) = \sum_{i=(j-1)K+1}^{jK} q_i(y)$, be the number of times to sample the j -th group in y . Let $\tilde{u}_0^j(y) = \sum_{i=(j-1)K+1}^{jK} u_0^i(y)$ be the number of times that sampling the j -th group with result 1; $\tilde{u}_1^j(y) = \sum_{i=(j-1)K+1}^{jK} u_1^i(y)$ be the number of times that sampling the j -th group with result 0;
- let $T(y)$ be the output of \mathcal{A} when the transcript generated by \mathcal{A} is y (note that the output of \mathcal{A} is completed determined by y since \mathcal{A} is deterministic).

Let the random variable Y be the transcript generated by \mathcal{A} . By (19), we have

$$\sum_{y \in \mathcal{Y}} \mathbf{1} \left[\text{val}_{T(y)}^{H_1}(K) \geq \text{val}^{H_1}(K) - \epsilon \text{ and } \tilde{q}_{j_0} \leq Q_0 \text{ and } |u_0^{j_0}(y) - u_1^{j_0}(y)| \leq \sqrt{5Q_0} \right] \cdot \Pr_{H_1}[Y = y] \geq \frac{1}{4}. \quad (20)$$

Fix a $y \in \mathcal{Y}$, since

$$\Pr_{H_1}[Y = y] = \left(\frac{1}{2} + 4\epsilon \right)^{\tilde{u}_1^1(y)} \left(\frac{1}{2} - 4\epsilon \right)^{\tilde{u}_0^1(y)} \prod_{j=2}^t \left(\frac{1}{2} \right)^{\tilde{q}_j(y)},$$

and

$$\Pr_{H_{j_0}}[Y = y] = \left(\frac{1}{2} + 4\epsilon \right)^{\tilde{u}_1^1(y)} \left(\frac{1}{2} - 4\epsilon \right)^{\tilde{u}_0^1(y)} (1 - 16\epsilon)^{\tilde{u}_0^{j_0}(y)} (1 + 16\epsilon)^{\tilde{u}_1^{j_0}(y)} \prod_{j=2}^t \left(\frac{1}{2} \right)^{\tilde{q}_j(y)},$$

we have

$$\begin{aligned} \frac{\Pr_{H_{j_0}}[Y = y]}{\Pr_{H_1}[Y = y]} &= (1 - 16\epsilon)^{\tilde{u}_0^{j_0}(y)} (1 + 16\epsilon)^{\tilde{u}_1^{j_0}(y)} = (1 - 16\epsilon)^{\frac{\tilde{q}_{j_0}(y)}{2} + \frac{\tilde{u}_0^{j_0}(y) - \tilde{u}_1^{j_0}(y)}{2}} (1 + 16\epsilon)^{\frac{\tilde{q}_{j_0}(y)}{2} - \frac{\tilde{u}_0^{j_0}(y) - \tilde{u}_1^{j_0}(y)}{2}} \\ &= (1 - 256\epsilon^2)^{\frac{\tilde{q}_{j_0}(y)}{2}} \left(\frac{1 - 16\epsilon}{1 + 16\epsilon} \right)^{\frac{\tilde{u}_0^{j_0}(y) - \tilde{u}_1^{j_0}(y)}{2}} \geq (1 - 256\epsilon^2)^{\frac{\tilde{q}_{j_0}(y)}{2}} (1 - 32\epsilon)^{\left| \frac{\tilde{u}_0^{j_0}(y) - \tilde{u}_1^{j_0}(y)}{2} \right|}. \end{aligned} \quad (21)$$

When $\tilde{q}_{j_0} \leq Q_0$ and $|u_0^{j_0}(y) - u_1^{j_0}(y)| \leq \sqrt{5Q_0}$ holds, we have (recall that $Q_0 \leq \frac{\ln(1/\delta)}{1250\epsilon^2}$)

$$\begin{aligned} (1 - 256\epsilon^2)^{\frac{\tilde{q}_{j_0}(y)}{2}} (1 - 32\epsilon)^{\left| \frac{\tilde{u}_0^{j_0}(y) - \tilde{u}_1^{j_0}(y)}{2} \right|} &\geq (1 - 256\epsilon^2)^{Q_0/2} (1 - 32\epsilon)^{\sqrt{5Q_0}/2} \\ &\geq (1 - 256\epsilon^2)^{\frac{\ln(1/\delta)}{2500\epsilon^2}} (1 - 32\epsilon)^{\frac{\sqrt{\ln(1/\delta)}}{30\epsilon}} \geq \delta^{1/4} \cdot \delta^{1/4} = \sqrt{\delta}, \end{aligned} \quad (22)$$

where in the penultimate inequality we used the assumption that $0 < \epsilon, \delta \leq .01$. Therefore we have

$$\begin{aligned} &\Pr_{H_{j_0}} \left[\text{val}_K^{H_1}(T) \geq \text{val}_K^{H_1}([n]) - \epsilon \right] \\ &\geq \Pr_{H_{j_0}} \left[\text{val}_K^{H_1}(T) \geq \text{val}_K^{H_1}([n]) - \epsilon \text{ and } \tilde{q}_{j_0} \leq Q_0 \text{ and } |u_0^{j_0}(y) - u_1^{j_0}(y)| \leq \sqrt{5Q_0} \right] \\ &= \sum_{y \in \mathcal{Y}} \mathbf{1} \left[\text{val}_K^{H_1}(T(y)) \geq \text{val}_K^{H_1}([n]) - \epsilon \text{ and } \tilde{q}_{j_0} \leq Q_0 \text{ and } |u_0^{j_0}(y) - u_1^{j_0}(y)| \leq \sqrt{5Q_0} \right] \cdot \Pr_{H_{j_0}}[Y = y] \\ &\geq \sum_{y \in \mathcal{Y}} \mathbf{1} \left[\text{val}_K^{H_1}(T(y)) \geq \text{val}_K^{H_1}([n]) - \epsilon \text{ and } \tilde{q}_{j_0} \leq Q_0 \text{ and } |u_0^{j_0}(y) - u_1^{j_0}(y)| \leq \sqrt{5Q_0} \right] \cdot \Pr_{H_1}[Y = y] \cdot \sqrt{\delta} \\ &\geq \frac{\sqrt{\delta}}{4}, \end{aligned}$$

where the penultimate inequality is because of (21) and (22); and the last inequality is because of (20). Therefore, we finish the proof of the Eq. (16); and by Eq. (18), we conclude the proof of Lemma 5.4. \square

Proof of Theorem 5.5. We show that essentially the same lower bound also holds for any randomized algorithm. The following argument is standard and we include it for completeness. Fix $0 < \epsilon, \delta < 1/2$. We assume, for contradiction, that there is a randomized algorithm \mathcal{A} which can achieve the same performance guarantee stated as in the theorem, but the expected number Q of samples is no more than $\frac{n \log(1/\delta)}{100000\epsilon^2 K}$. We can view the randomized algorithm \mathcal{A} as a deterministic algorithm with a sequence S of random bits. We use R to denote the randomness from the arms. Note that if we fix S and R , the execution and the output of the algorithm are fixed. We use $\mathcal{A}(S, R) = 1$ to denote the event that the output of \mathcal{A} is an ϵ -optimal solution. Let us use $Q(S, R)$ to denote the number of samples taken by \mathcal{A} . The performance guarantee of \mathcal{A} is that

$$\Pr_{S,R}[\mathcal{A}(S, R) = 1] = \mathbf{E}_{S,R}[\mathcal{A}(S, R)] = \mathbf{E}_S \mathbf{E}_R[\mathcal{A}(S, R) \mid S] \geq 1 - \delta.$$

This is equivalent to say that $\mathbf{E}_S \mathbf{E}_R[1 - \mathcal{A}(S, R) \mid S] \leq \delta$. By Markov inequality, we have that

$$\Pr_S \left[\mathbf{E}_R[1 - \mathcal{A}(S, R) \mid S] \geq 2\delta \right] \leq 1/2.$$

Equivalently, we have that

$$\Pr_S \left[\mathbf{E}_R[\mathcal{A}(S, R) \mid S] \geq 1 - 2\delta \right] \geq 1/2. \quad (23)$$

By our assumption, we have $\mathbf{E}_{S,R} Q(S, R) \leq \frac{n \log(1/\delta)}{100000\epsilon^2 K}$. So, by Markov inequality,

$$\Pr_S \left[\mathbf{E}_R[Q(S, R) \mid S] \leq \frac{n \log(1/\delta)}{40000\epsilon^2 K} \right] \geq \frac{3}{5}. \quad (24)$$

Combining (23) and (24), we know there is a particular random sequence S such that both $\mathbf{E}_R[\mathcal{A}(S, R) \mid S] \geq 1 - 2\delta$ and $\mathbf{E}_R[Q(S, R) \mid S] \leq \frac{n \log(1/\delta)}{40000\epsilon^2 K}$ hold. Since \mathcal{A} with a particular sequence S is simply a deterministic algorithm, this contradicts the lower bound we proved for any deterministic algorithm in Lemma 5.4. \square

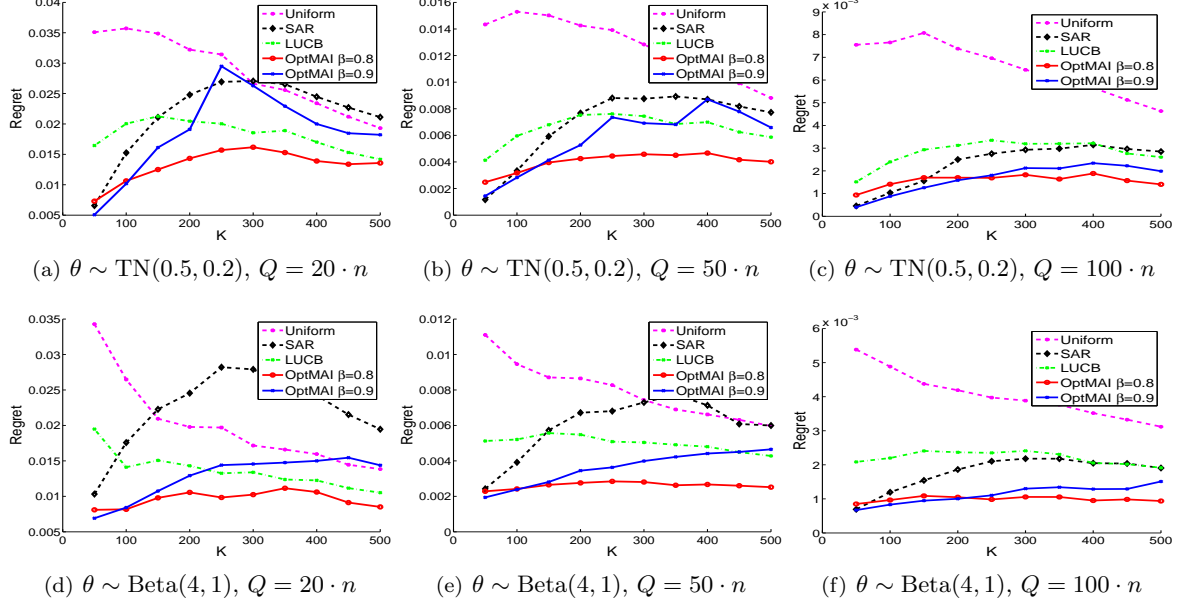


Figure 3: Performance comparison on simulated data.

E Additional Experiments

In this section, we provide additional simulated experimental results when using the following two different ways to generate $\{\theta_i\}_{i=1}^n$:

1. $\theta \sim \text{TN}(0.5, 0.2)$: each θ_i is generated from a truncated normal distribution with mean 0.5, the standard deviation 0.2 and the support $[0, 1]$ (Figure 3(a) to Figure 3(c)).
2. $\theta \sim \text{Beta}(4, 1)$: each θ_i is generated from a Beta distribution with the parameters (4, 1). The $\{\theta_i\}$ from Beta(4, 1) are close to the workers' accuracy in real crowdsourcing applications, where most workers perform reasonably well and the averaged accuracy is around 80% (Figure 3(d) to Figure 3(f)).

We note that the number of total arms is set to $n = 1000$. We vary the total budget $Q = 20n, 50n, 100n$ and $K = 10, 20, \dots, 500$. We use different ways to generate $\{\theta_i\}_{i=1}^n$ and report the comparison among different algorithms. It can be seen from Figure 1 that our method outperforms the SAR and LUCB in most of the scenarios. In addition, we also observe that when K is large, the setting of $\beta = 0.8$ outperforms that of $\beta = 0.9$; while for small K , $\beta = 0.9$ is a better choice.