
Globally Optimal Algorithms for Fixed-Budget Best Arm Identification

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We consider the fixed-budget best arm identification problem where the goal
2 is to find the arm of the largest mean with a fixed number of samples. It is
3 known that the probability of misidentifying the best arm is exponentially
4 small to the number of rounds. However, limited characterizations have
5 been discussed on the rate (exponent) of this value. In this paper, we
6 characterize the optimal rate as a result of global optimization over all
7 possible parameters. We introduce two rates, R^{go} and R_{∞}^{go} , corresponding to
8 lower bounds on the misidentification probability, each of which is associated
9 with a proposed algorithm. The rate R^{go} is associated with R^{go} -tracking,
10 which can be efficiently implemented by a neural network and is shown to
11 outperform existing algorithms. However, this rate requires a nontrivial
12 condition to be achievable. To deal with this issue, we introduce the second
13 rate R_{∞}^{go} . We show that this rate is indeed achievable by introducing a
14 conceptual algorithm called delayed optimal tracking (DOT).

15 1 Introduction

16 We consider K -armed best arm identification problem with T samples. In this problem,
17 each arm $i \in [K] = \{1, 2, \dots, K\}$ is associated with (unknown) distribution $P_i \in \mathcal{P}$ for some
18 class of distributions \mathcal{P} . Upon choosing arm i , the forecaster observes reward $X(t)$, which is
19 independently drawn from P_i . The forecaster then tries to identify (one of) the best arm¹
20 $\mathcal{I}^* = \arg \max_i \mu_i$ with the largest mean $\mu^* = \max_i \mu_i$ for $\mu_i = \mathbb{E}_{X \sim P_i}[X]$. The problem² is
21 called the best arm identification (BAI, Audibert et al. (2010)), or the ranking and selection
22 (R&S, Hong et al. (2021)).

23 To this aim, the forecaster uses some algorithm that would adaptively choose an arm based
24 on its history of rewards. At each round t , the algorithm chooses one of the arms $I(t) \in [K]$
25 and receives the corresponding reward $X(t)$. After the T -th round, the algorithm outputs a
26 recommendation arm $J(T) \in [K]$, which corresponds to an estimator of the best arm. The
27 misidentification probability is expressed by $\mathbb{P}[J(T) \notin \mathcal{I}^*]$, which will be referred to as the
28 probability of the error (PoE) throughout the paper. Best arm identification has two settings.
29 In the fixed confidence setting, the forecaster minimizes the number of draws T until the
30 confidence level on the PoE reaches a given value $\delta \in (0, 1)$. In this case, T is a stopping time
31 that can be chosen adaptively. In the fixed-budget setting, the forecaster tries to minimize
32 the PoE given a constant T . In this paper, we shall focus on the fixed-budget setting. In

¹We use $\mathcal{I}^* = \mathcal{I}^*(\mathbf{P}) \subset [K]$ as the set of best arms and $i^* = i^*(\mathbf{P}) \in \mathcal{I}^*(\mathbf{P})$ as one of them (ties are broken in an arbitrary way). These differences does not matter much in this paper.

²See Section 1.3 regarding the related work on BAI and R&S.

33 general, a good algorithm for the fixed-confidence setting is very different from that for the
 34 fixed-budget setting. To be more specific, an algorithm for the fixed-confidence setting can
 35 be instance-wise optimal³. Namely, several algorithms exist (Garivier and Kaufmann, 2016)
 36 that can be optimized for each instance of distributions $\mathbf{P} = (P_1, P_2, \dots, P_K)$ as far as we
 37 consider algorithms called δ -PAC. By contrast, an algorithm for the fixed-budget setting
 38 requires consideration of the possibility that improving the PoE for an instance \mathbf{P} worsens
 39 the PoE for another instance \mathbf{P}' . Thus, we must consider a kind of a global optimization of
 40 the performance over all possible $\mathbf{P} \in \mathcal{P}^K$.

41 1.1 Global optimality in the fixed-budget setting

42 In the fixed-budget setting, the PoE decays exponentially to T as $\exp(-RT)$ for some *rate*
 43 $R > 0$. The instance-wise optimality given above is no longer available here. To demonstrate
 44 this, assume that we make an estimate of \mathbf{P} based on the initial $o(T)$ rounds, say, \sqrt{T}
 45 rounds. In this case, we can obtain the estimate of \mathbf{P} that is ϵ -correct with probability
 46 $\exp(-\epsilon^2 O(\sqrt{T})) = \exp(-o(T))$. However, this estimation does not help to improve the rate
 47 of exponential convergence. In other words, estimating \mathbf{P} requires non-negligible (i.e. $O(T)$)
 48 cost for exploration. As a result, we cannot optimize the PoE for each instance \mathbf{P} unlike the
 49 fixed-confidence setting. Instead, to discuss optimality in the fixed-budget setting, we must
 50 choose a complexity function $H(\mathbf{P})$, and the performance of an algorithm must be evaluated
 51 on the rate normalized by the complexity.

52 In literature, little is known about the optimal rate of the exponent. Audibert et al.
 53 (2010) proposed the successive rejects (SR) algorithm, which has the rate of $1/(\log K)$
 54 with the complexity function $H_2(\mathbf{P}) := \max_{i \in [K]} \frac{i}{\Delta_i^2}$ for $\Delta_i = \max_j \mu_j - \mu_i$ satisfying
 55 $\Delta_1 \leq \Delta_2 \leq \dots \leq \Delta_K$. Carpentier and Locatelli (2016) showed a particular set of instances
 56 such that this rate matches the lower bound up to a constant factor. However, the constant
 57 used there is by far loose⁴, and there is limited discussion on the actual rate of such algorithms.

58 1.2 Contributions

59 This paper tightly characterizes the optimal minimax rate of the PoE as a result of a
 60 global optimization given \mathbf{P} . Let $H = H(\mathbf{P}) : \mathcal{P}^K \rightarrow \mathbb{R}^+$ be any continuous complexity
 61 measure. We then discuss the best possible rate $R > 0$ such that the PoE is bounded by
 62 $\exp(-RT/H(\mathbf{P}) + o(T))$ for all $\mathbf{P} \in \mathcal{P}^K$ and make the following contributions.

- 63 • We derive an upper bound on R (corresponding to a lower bound of the PoE), denoted
 64 by R^{go} , which we obtain by considering a class of oracle algorithms that can determine
 65 the allocation of trials to each arm knowing the final empirical distribution after T
 66 rounds (Theorem 1).
- 67 • We propose an algorithm (R^{go} -tracking) that greedily tracks this oracle allocation
 68 based on the current empirical distribution (Section 2.1). Though this oracle allocation
 69 is expressed by a complicated minimax optimization, we propose a technique to learn
 70 this by a neural network and empirically confirm that the PoE of the learned algorithm
 71 is close to the lower bound (Sections 3 and 4). We also discuss that the algorithm is
 72 unlikely to provably achieve the bound even when the minimax problem is perfectly
 73 solved because of the impossibility of the tracking.
- 74 • We tighten the PoE lower bound by weakening the oracle algorithms. Based on
 75 this refined bound, we propose the delayed optimal tracking (DOT) algorithm that
 76 asymptotically achieves the tightened lower bound for Bernoulli and Gaussian arms,
 77 though the algorithm is computationally almost infeasible (Sections 2.2 and 2.3).

78 In summary, we propose a nearly-tight PoE lower bound with a computationally feasible
 79 algorithm that is empirically close to this bound. We also propose a provably tight lower
 80 bound and matching algorithm in a computationally infeasible form. Notation is listed in
 81 the supplementary material.

³A more complete discussion on this topic can be found in Section D of supplementary material.

⁴Theorem 1 therein includes a large constant 400.

82 **1.3 Related work**

83 Compared with the works of the fixed-confidence BAI, less is known about the fixed-budget
 84 BAI. For example, a book on this subject (Lattimore and Szepesvári, 2020) spends only two
 85 pages on the fixed-budget BAI.⁵ Many algorithms designed for the fixed-confidence BAI,
 86 such as D-tracking (Kaufmann et al., 2016), do not have a finite-time PoE guarantee when we
 87 apply them to the fixed-budget setting. Nevertheless, there are two well-known fixed-budget
 88 BAI algorithms: Successive rejection (SR, Audibert et al. (2010)) and successive halving (SH,
 89 Shahrampour et al. (2017)). Both SR and SH progressively narrow the candidate of the best
 90 arm at the end of each segment. While SR discards one arm after each segment, SH discards
 91 half of the remaining arms after each segment. SR and SH have the guarantee on PoE of the
 92 rate $\exp(-RT/H_2(\mathbf{P}))$ for some constant $R > 0$. Other fixed-budget BAI algorithms, such
 93 as UCB-E (Audibert et al., 2010) and UGapE (Gabillon et al., 2012), require the knowledge
 94 of minimum gap $\min_i \Delta_i$, and thus are not universal to all best arm identification instances.

95 Another literature on this topic is the ranking and selection (R&S) problems (Powell and
 96 Ryzhov, 2018; Hong et al., 2021). Although the goal of R&S problems is to identify the best
 97 arm, many R&S papers do not consider the estimation error of \mathbf{P} in a finite time. As a result,
 98 algorithms therein do not have the guarantee on the PoE in the best arm identification setting.
 99 The optimal computing budget allocation (OCBA Chen et al. (2000); Glynn and Juneja
 100 (2004)) algorithm tries to minimize the PoE assuming the plug-in estimator matches the true
 101 parameter. Bayesian R&S algorithms try to solve the dynamic programming of minimizing
 102 the PoE given a prior, which is computationally prohibitive, and thus approximated solutions
 103 have been sought (Frazier et al., 2008; Powell and Ryzhov, 2018).

104 **2 Globally optimal algorithm**

105 In this section, we derive several lower bounds on the PoE and propose algorithms to
 106 empirically or theoretically achieve these bounds.

107 First, we formalize the problem. Let \mathcal{P} be a known class of reward distributions. We consider
 108 the case where \mathcal{P} is the set of Bernoulli distributions with mean $\Theta \subset [0, 1]$ (including the
 109 case $\Theta = [0, 1]$), or Gaussian distributions with mean in $\Theta \subset \mathbb{R}$ (including the case $\Theta = \mathbb{R}$)
 110 and known variance $\sigma^2 > 0$. It should be noted that many parts of results in this paper can
 111 be generalized to much wider classes of distributions, but it makes the notation much longer
 112 and is discussed in Appendix E.

113 When we derive lower bounds and construct algorithms, we introduce \mathcal{Q} as a class of
 114 distributions corresponding to the estimated distributions of the arms. Namely, we set \mathcal{Q} as
 115 the set of all Bernoulli (resp. Gaussian) distributions with mean in $[0, 1]$ (resp. \mathbb{R}) when \mathcal{P} is
 116 the set of Bernoulli (resp. Gaussian) distributions with mean in Θ . As such, we take $\mathcal{Q} \supset \mathcal{P}$
 117 so that the estimator of P_i is always in \mathcal{Q} . In these models, we identify the distribution
 118 $P_i \in \mathcal{P}$ with its mean parameter in $\Theta \subset \mathbb{R}$.

119 Our interest lies in the rate $\lim_{T \rightarrow \infty} \frac{1}{T} \log(1/\mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})])$ of convergence of the PoE.
 120 Since we are interested in lower and upper bounds of the rate of algorithms including those
 121 requiring the knowledge of T , we define the rate for a sequence of algorithms $\{\pi_T\}$ by

$$R(\{\pi_T\}) = \inf_{\mathbf{P} \in \mathcal{P}^\kappa} H(\mathbf{P}) \liminf_{T \rightarrow \infty} \frac{1}{T} \log(1/\mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})]). \quad (1)$$

122 Here, a larger $R(\{\pi_T\})$ corresponds to a faster convergence of the PoE.

123 **2.1 PoE for oracle algorithms**

124 First, we derive a lower bound on the PoE that is unlikely to be achievable but strongly
 125 related to an optimal algorithm. Let $D(P||Q) = \mathbb{E}_{X \sim P}[\frac{dP}{dQ}(X)]$ be the Kullback–Leibler
 126 (KL) divergence between P and Q . Then we have the following bound.

⁵Section 33.3 therein.

Algorithm 1: R^{go} -Tracking

input : (ϵ) -optimal solution $(\mathbf{r}^*(\cdot), J^*(\cdot))$ of (2).

1 Draw each arm once.

2 **for** $t = K + 1, 2, \dots, T$ **do**3 \lfloor Draw arm $\arg \max_{i \in [K]} \{r_i^*(\mathbf{Q}(t-1)) - N_i(t-1)/(t-1)\}$.4 **return** $J(T) = J^*(\mathbf{Q})$.

127 **Theorem 1.** Under any sequence of algorithm $\{\pi_T\}$ it holds that

$$R(\{\pi_T\}) \leq \sup_{\mathbf{r}(\cdot) \in \Delta^K, J(\cdot) \in [K]} \inf_{\mathbf{Q} \in \mathcal{Q}^K} \inf_{\mathbf{P} \in \mathcal{P}^K: J(\mathbf{Q}) \notin \mathcal{I}^*(\mathbf{P})} H(\mathbf{P}) \sum_{i \in [K]} r_i(\mathbf{Q}) D(Q_i \| P_i) =: R^{\text{go}}, \quad (2)$$

128 where the outer supremum is taken over all functions $\mathbf{r}(\cdot) : \mathcal{Q}^K \rightarrow \Delta^K$, $J(\cdot) : \mathcal{Q}^K \rightarrow [K]$.

129 All proofs are provided in the supplementary material owing to page limitation. This theorem
130 states that under any algorithm there exists an instance \mathbf{P} such that the PoE is at least
131 $\exp(-TR^{\text{go}}/H(\mathbf{P}) + o(T))$. Intuitively speaking, the bound in Theorem 1 corresponds to the
132 best possible rate of oracle algorithms that can determine the allocation as $\mathbf{r} = \mathbf{r}^*(\mathbf{Q}) \in \Delta^K$
133 knowing the final empirical mean $\mathbf{Q} = \mathbf{Q}(T)$, where $\mathbf{r}^*(\cdot)$ is the (ϵ) -optimal⁶ solution of (2).

134 From the technical viewpoint, the main difference from the lower bound on the fixed-
135 confidence setting is that we also have to consider candidates of empirical distributions \mathbf{Q}
136 as well as the true distributions \mathbf{P} . This makes the analysis much more difficult, because
137 a slight difference of the empirical distribution might (possibly discontinuously) affect the
138 allocation unlike the difference of the true distribution \mathbf{P} unknown to the algorithm. A
139 naive analysis just depending on the empirical distribution fails because of this discontinuity
140 of the allocation. To overcome this difficulty, we adopt a technique inspired by the *typical*
141 *set analysis* often used in the information theory (Cover and Thomas, 2006). We define the
142 *typical allocation* for each candidate of empirical distribution \mathbf{Q} and prove the theorem by
143 evaluating the error probability based on the typical allocation.

144 **Remark 1.** We can take arbitrary $H(\mathbf{P}) > 0$ as a complexity measure, but R^{go} might
145 become zero if $H(\mathbf{P})$ is not taken reasonably. When $R^{\text{go}} = 0$ any algorithm trivially satisfies
146 $\text{PoE} \leq \exp(-TR^{\text{go}}/H(\mathbf{P}) + o(T))$. This means that any algorithm is minimax-optimal in
147 terms of $H(\mathbf{P})$, that is, such choice of $H(\mathbf{P})$ gives meaningless results.

148 In the actual trial, the algorithm can only know the empirical mean $\mathbf{Q}(t-1)$ at the beginning
149 of the current round t and we cannot ensure the achievability of the bound for oracle
150 algorithms. Despite this, one reasonable choice of the algorithm would be to keep tracking
151 this optimal allocation $\mathbf{r}^*(\mathbf{Q}(t-1))$, expecting that the current empirical mean $\mathbf{Q}(t-1)$
152 is close to $\mathbf{Q}(T)$. R^{go} -tracking in Algorithm 1 is the algorithm based on this idea. Here,
153 $N_i(t-1)$ is the number of times that the arm i is drawn at the beginning of the t -th round,
154 and it draws the arm such that the current fraction of the allocation $N_i(t-1)/(t-1)$ is the
155 most insufficient compared with the estimated optimal allocation $\mathbf{r}^*(\mathbf{Q}(t-1))$.

156 As we will see in Section 4, the empirical performance of Algorithm 1 is very close to the
157 PoE lower bound stated above. However, it is difficult to expect that this algorithm provably
158 achieves this bound in general because of the following: We could prove that R^{go} -tracking is
159 optimal if the fraction of allocation always satisfies $N_i(t)/t = \mathbf{r}(\mathbf{Q}(t)) + o(1)$, that is, the
160 algorithm can track the ideal allocation $\mathbf{r}(\mathbf{Q}(t))$. However, this does not generally hold. For
161 example, the empirical mean $\mathbf{Q}(t)$ sometimes changes rapidly in the Gaussian case. Whilst
162 this event occurs with exponentially small probability, the PoE itself is also an exponentially
163 small probability and it is highly nontrivial to specify in which case the tracking failure
164 probability becomes negligible.

165 **Remark 2.** Eq. (2) also involves the optimization of the recommendation arm $J(\mathbf{Q})$ as
166 well as $\mathbf{r}(\mathbf{Q})$. We can easily see that it is optimal to set $J(\mathbf{Q}) = i^*(\mathbf{Q})$, that is, taking

⁶This paper uses $\epsilon > 0$ as an arbitrarily small gap to the optimal solution. An asterisk is used to denote optimality.

167 the empirical best arm as the recommendation arm when $\mathcal{P} = \mathcal{Q}$ since $R(\pi)$ becomes zero
 168 otherwise. However, $J(\mathbf{Q}) = i^*(\mathbf{Q})$ might not hold for $\mathbf{Q} \notin \mathcal{P}$ when $\mathcal{P} \subsetneq \mathcal{Q}$.

169 2.2 PoE considering trackability

170 To construct an algorithm that is provably optimal, we begin with refining the PoE lower
 171 bound by weakening the “strength” of the oracle algorithm.

172 We consider splitting T rounds into B batches of size $\lfloor T/B \rfloor$ or $\lfloor T/B \rfloor + 1$. Let

$$\mathbf{r}^B = (\mathbf{r}_1(\mathbf{Q}_1), \mathbf{r}_2(\mathbf{Q}_1, \mathbf{Q}_2), \mathbf{r}_3(\mathbf{Q}_1, \mathbf{Q}_2, \mathbf{Q}_3), \dots, \mathbf{r}_B(\mathbf{Q}_1, \dots, \mathbf{Q}_B))$$

173 be a sequence of B functions, where $\mathbf{r}_b : \mathcal{Q}^{Kb} \rightarrow \Delta^K$ corresponds the allocation in the b -th
 174 batch when the empirical means of the first b batches are $\mathbf{Q}^b = (\mathbf{Q}_1, \mathbf{Q}_2, \dots, \mathbf{Q}_b)$. Based on
 175 this class of allocation rule, we have the following PoE lower bound.

176 **Theorem 2.** (PoE Bound for batch-oracle algorithms) Under any sequence of algorithms
 177 π_T and $B \in \mathbb{N}$,

$$R(\{\pi_T\}) \leq \sup_{\mathbf{r}^B(\cdot), J(\cdot)} \inf_{\mathbf{Q}^B \in \mathcal{Q}^{KB}} \inf_{\mathbf{P}: J(\mathbf{Q}^B) \notin \mathcal{I}^*(\mathbf{P})} \frac{H(\mathbf{P})}{B} \sum_{i \in [K], b \in [B]} r_{b,i} D(Q_{b,i} \| P_i) =: R_B^{\text{go}}. \quad (3)$$

178 Here, the outer supremum is taken over all functions $\mathbf{r}^B(\cdot) = (\mathbf{r}_1(\cdot), \mathbf{r}_2(\cdot), \dots, \mathbf{r}_B(\cdot))$ for
 179 $\mathbf{r}_b(\cdot) : \mathcal{Q}^{Kb} \rightarrow \Delta^K$ and $J(\cdot) : \mathcal{Q}^{KB} \rightarrow [K]$.

180 Theorem 1 is the special case of this theorem with $B = 1$. This bound corresponds to
 181 the best bound of oracle algorithms that can determine the allocation of the b -th batch
 182 knowing the empirical distribution of this batch. It is tighter than that given in Theorem 1,
 183 as the oracle considered here cannot know the empirical distribution of the later batches
 184 $b + 1, b + 2, \dots, B$. It follows that we can obtain the following result.

185 **Corollary 3.** We have $R_B^{\text{go}} \leq R^{\text{go}}$ for any $B \in \mathbb{N}$.

186 We will show that $R_\infty^{\text{go}} := \lim_{B \rightarrow \infty} R_B^{\text{go}}$ exists and is the best possible rate.

187 2.3 Matching algorithm

188 In this section, we derive an algorithm that has a rate that almost matches R_B^{go} . For any $\epsilon > 0$,
 189 let an ϵ -optimal solution of Eq. (3) be $(\mathbf{r}^{B,*}(\cdot), J^*(\cdot)) = (\mathbf{r}_1^*(\cdot), \mathbf{r}_2^*(\cdot), \mathbf{r}_3^*(\cdot), \dots, \mathbf{r}_B^*(\cdot), J^*(\cdot))$
 190 with its objective at least

$$\inf_{\mathbf{Q}^B \in \mathcal{Q}^{KB}} \inf_{\mathbf{P}: J^*(\mathbf{Q}^B) \notin \mathcal{I}^*(\mathbf{P})} \frac{H(\mathbf{P})}{B} \sum_{i \in [K], b \in [B]} r_{b,i}^*(\mathbf{Q}^b) D(Q_{b,i} \| P_i) \geq R_B^{\text{go}} - \epsilon.$$

191 We cannot naively follow the allocation $r_{b,i}^*(\mathbf{Q}^b)$ because it requires the empirical mean
 192 of the current batch \mathbf{Q}_b , which is not fully available until the end of the current batch.
 193 The delayed optimal tracking algorithm (DOT, Algorithm 2) addresses this issue. This
 194 algorithm divides T rounds into $B + K - 1$ batches, where the b -th batch corresponds to
 195 $(bT_B + 1, bT_B + 2, \dots, (b + 1)T_B)$ -th rounds for $T_B = T/(B + K - 1)$. Here, for simplicity,
 196 we assume that T is a multiple of $B + K - 1$. In the other case, we can reach almost the
 197 same result by just ignoring the last $T - (B + K - 1)\lfloor T/(B + K - 1) \rfloor$ rounds.

198 The crux of Algorithm 2 is to determine allocation \mathbf{r}_b by using the *stored* empirical mean
 199 $\mathbf{Q}'_1, \mathbf{Q}'_2, \dots, \mathbf{Q}'_B$ rather than the true empirical mean $\mathbf{Q}_1, \mathbf{Q}_2, \dots, \mathbf{Q}_{B+K-1}$: The first K batches
 200 are devoted to uniform exploration and the samples are stored in a queue (though this
 201 explanation is not strict, in that the actual procedure is done after taking the mean of the
 202 stored samples). At the b -th batch for $b \geq K + 1$, we draw each arm i for $n_{b,i} \approx T_B \mathbf{r}_b$ times⁷,
 203 where \mathbf{r}_b is determined based on the stored samples in the queue. When drawing arm i
 204 for $n_{b,i}$ times, we dequeue and open $n_{b,i}$ stored samples instead of opening the actual $n_{b,i}$
 205 samples, the latter of which are enqueued and kept unopened.

206 By the nature of this algorithm we can ensure the following property.

⁷The $-K$ in Line 6 of Algorithm 2 is for the ceiling fractional values. This is reflected in the term T' in Theorem 5. If T is large compared to B, K , the difference between T and T' does not matter.

Algorithm 2: Delayed optimal tracking (DOT)

input : ϵ -optimal solution $\mathbf{r}^{B,*}(\cdot) = (\mathbf{r}_1^*(\cdot), \mathbf{r}_2^*(\cdot), \dots, \mathbf{r}_B^*(\cdot), J^*(\cdot))$ of (3).

- 1 **for** $b = 1, 2, \dots, K$ **do**
- 2 | Set $r_{b,i} = \mathbf{1}[i = b]$ for $i \in [K]$ and draw arm b for T_B times.
- 3 Set $\mathbf{Q}'_1 := \mathbf{Q}_K$ for the empirical mean \mathbf{Q}_K .
- 4 **for** $b = K + 1, K + 2, \dots, B + K - 1$ **do**
- 5 | Compute $\mathbf{r}_b = (r_{b,1}, r_{b,2}, \dots, r_{b,K}) = \mathbf{r}_{b-K}^*(\mathbf{Q}'_1, \mathbf{Q}'_2, \dots, \mathbf{Q}'_{b-K})$.
- 6 | Draw each arm i for $n_{b,i}$ times, where $n_{b,i} \geq r_{b,i}(T_B - K)$ is taken so that $\sum_{i \in [K]} n_{b,i} = T_B$.
- 7 | Observe empirical mean \mathbf{Q}_b of the batch.
- 8 | Update the *stored* empirical average as

$$\mathbf{Q}'_{b-K+1} = \mathbf{Q}'_{b-K} + \mathbf{r}_b(\mathbf{Q}_b - \mathbf{Q}'_{b-K}),$$
 where $\mathbf{r}_b \mathbf{Q}$ denotes the element-wise product.
- 9 Recommend $J(T) = J^*(\mathbf{Q}'_1, \mathbf{Q}'_2, \dots, \mathbf{Q}'_B)$.

207 **Lemma 4.** Assume that we run Algorithm 2. Then, the following inequality always holds:

$$\frac{1}{B + K - 1} \sum_{i \in [K], b \in [B + K - 1]} r_{b,i} D(Q_{b,i} \| P_i) \geq \frac{B}{B + K - 1} \frac{R_B^{\text{go}} - \epsilon}{H(\mathbf{P})}. \quad (4)$$

208 Lemma 4 states that the empirical divergence of DOT given in the LHS of (4) almost matches
 209 the upper bound $R_B^{\text{go}}/H(\mathbf{P})$ for sufficiently large B despite the delayed allocation. Using
 210 this property we obtain the following achievability bound.

211 **Theorem 5.** (Performance bound of Algorithm 2) The PoE of the DOT algorithm satisfies

$$\mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \leq \exp\left(-\frac{BT'}{B + K - 1} \frac{R_B^{\text{go}} - \epsilon}{H(\mathbf{P})} + f(K, B, T)\right),$$

212 where $T' = T - (B + K - 1)K$ and $f(K, B, T) = 2BK \log(2T)$.

213 The following corollary is immediate since $f(K, B, T) = o(T)$ holds for fixed K, B .

214 **Corollary 6.** The worst-case rate of the DOT algorithm $\pi_{\text{DOT}, T}$ satisfies

$$R(\{\pi_{\text{DOT}, T}\}) \geq \frac{B}{B + K - 1} (R_B^{\text{go}} - \epsilon).$$

215 2.4 Optimality

216 In this section, we show the rate $R(\pi_{\text{DOT}})$ of DOT becomes arbitrarily close to optimal when
 217 we take a sufficiently large number of batches B .

218 **Theorem 7.** (Optimality of DOT) Assume $H(\mathbf{P})$ be such that $R^{\text{go}} < \infty$. Then, the limit

$$R_\infty^{\text{go}} := \lim_{B \rightarrow \infty} R_B^{\text{go}} \quad (5)$$

219 exists. Moreover, for any $\eta > 0$, there exist parameters B, ϵ such that the following holds on
 220 the performance of the DOT algorithm:

$$R(\{\pi_{\text{DOT}, T}\}) = \inf_{\mathbf{P} \in \mathcal{P}} H(\mathbf{P}) \liminf_{T \rightarrow \infty} \frac{\log(1/\mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})])}{T} \geq R_\infty^{\text{go}} - \eta. \quad (6)$$

221 **Remark 3.** (η -optimality) Since $R(\{\pi_T\}) \leq R_B^{\text{go}}$ holds for any sequence of algorithms $\{\pi_T\}$
 222 and $b \in \mathbb{N}$, we have

$$R(\{\pi_T\}) \leq \inf_{B \in \mathbb{N}} R_B^{\text{go}} \leq \liminf_{B \rightarrow \infty} R_B^{\text{go}} = R_\infty^{\text{go}}$$

223 from (5). Therefore the rate of the DOT algorithm given in (6) is optimal up to η for
 224 arbitrary small $\eta > 0$. This essentially states that, no algorithm can be η -better than DOT
 225 in terms of the rate against the worst-case instance \mathbf{P} .

Algorithm 3: Gradient descent method for θ

input : learning rate η

```
1 while not converged do
  /* Compute  $\mathbf{P}^{\min}$  and  $\mathbf{Q}^{\min}$  which minimizes the negative exponent */
2   Set  $E^{\min} \leftarrow \infty$ .
3   for  $n_1 = 1, 2, \dots, N^{\text{true}}$  do
4     Sample true parameters  $\mathbf{P}$  from  $\mathcal{P}$  uniformly at random.
5     for  $n_2 = 1, 2, \dots, N^{\text{emp}}$  do
6       Sample  $\mathbf{Q}$  from  $\{\mathbf{Q} \in \mathcal{Q}^K : \mathcal{I}^*(\mathbf{Q}) \cap \mathcal{I}^*(\mathbf{P}) = \emptyset\}$ .
7       if  $E(\mathbf{P}, \mathbf{Q}; \theta) < E^{\min}$  then
8          $\mathbf{P}^{\min} \leftarrow \mathbf{P}, \mathbf{Q}^{\min} \leftarrow \mathbf{Q}, E^{\min} \leftarrow E(\mathbf{P}^{\min}, \mathbf{Q}^{\min}; \theta)$ 
9   Update parameters  $\theta \leftarrow \theta - \eta \nabla_{\theta} E(\mathbf{P}^{\min}, \mathbf{Q}^{\min}; \theta)$ .
```

Algorithm 4: R^{go} -Tracking by Neural Network (TNN)

```
1 Draw each arm once.
2 for  $t = K + 1, 2, \dots, T$  do
3   Draw arm  $\arg \max (r_{\theta, i}(\mathbf{Q}(t-1)) - N_i(t-1)/(t-1))$ .
4 return  $J(T) = \arg \max Q_i(T)$  (empirical best arm).
```

226 **Remark 4.** (Utility of the DOT algorithm) Although DOT (Algorithm 2) has an asymptoti-
227 cally optimal rate R_{∞}^{go} , it is difficult to calculate, or to even approximate, the optimal solution
228 of (3) since it is not an optimization of a finite-dimensional vector but an optimization of
229 function \mathbf{r}^B , which has high input dimension proportional to B . In this sense, DOT algorithm
230 as well as Theorem 7 is purely theoretical thus far, and the existence of a computationally
231 tractable and provably optimal algorithm is an important open question.

232 3 Learning

233 In this section, we propose a method to learn $\mathbf{r}(\mathbf{Q})$ of Eq. (2) by utilizing a neural network to
234 practically realize R^{go} -tracking in Algorithm 1. Throughout this section, we assume a class
235 of algorithms satisfying $J(\mathbf{Q}) \in \mathcal{I}^*(\mathbf{Q})$, that is, algorithms that recommend the empirical
236 best arm, which is guaranteed to be optimal when $\mathcal{P} = \mathcal{Q}$ (see Remark 2).

237 3.1 Learning allocation

238 Let $\mathbf{r}_{\theta}(\mathbf{Q}) : \mathcal{Q}^K \rightarrow \Delta^K$ be a neural network with a set of parameters θ . We consider
239 alternately optimizing $\mathbf{r}_{\theta}(\cdot)$ and (\mathbf{P}, \mathbf{Q}) , and we update θ via mini-batch gradient descent.
240 Given a complexity function $H(\mathbf{P})$, Eq. (2) is defined as the minimum over all (\mathbf{P}, \mathbf{Q}) such
241 that the best arm is different. Our learning method (Algorithm 3) uses L mini-batches⁸. Let

$$E(\mathbf{P}, \mathbf{Q}; \theta) := H(\mathbf{P}) \sum_{i=1}^K r_{\theta, i}(\mathbf{Q}) D(Q_i \| P_i). \quad (7)$$

242 Given allocation \mathbf{r}_{θ} , Eq. (7) is the negative log-likelihood (rate) of the bandit instance \mathbf{P}
243 given the empirical means \mathbf{Q} . At each batch, it obtains the pair $\mathbf{P}^{\min}, \mathbf{Q}^{\min}$ such that Eq. (7)
244 is minimized. Specifically, for each iteration, we sample N^{true} candidates of true means \mathbf{P}
245 uniformly from \mathcal{P}^K , then for each \mathbf{P} , we sample N^{emp} values of empirical means $\mathbf{Q} \in \mathcal{Q}^K$
246 such that $\mathcal{I}^*(\mathbf{Q}) \cap \mathcal{I}^*(\mathbf{P}) = \emptyset$ uniformly at random.

⁸Algorithm 3 is a conceptual explanation and the actual implementation used a momentum method. See Section 4.1 for implementation details.

247 3.2 Tracking by neural network

248 Having trained \mathbf{r}_θ , we propose the R^{go} -Tracking by Neural Network (TNN) algorithm
249 (Algorithm 4), which is an implementation of R^{go} -Tracking by the trained neural network.
250 This algorithm draws the arm such that the current fraction of samples $N_i(t-1)/(t-1)$ is
251 the most insufficient compared with the learned allocation $\mathbf{r}_\theta(\mathbf{Q}(t-1))$.

252 4 Simulation

253 This section tests numerically the performance of TNN algorithm. We compared the
254 performance of TNN (Algorithm 4) with two algorithms: Uniform algorithm, which samples
255 each arm in a round-robin fashion, and Successive Rejects (SR, Audibert et al., 2010), where
256 the entire trial is divided into segments before the game starts, and one arm with the smallest
257 estimated mean reward is removed for each segment.

258 We consider Bernoulli bandits with $K = 3$ arms, where each mean parameter is in $[0, 1]$. In
259 particular, we consider the three sets of true parameters: (instance 1) $\mathbf{P} = (0.5, 0.45, 0.3)$,
260 (instance 2) $\mathbf{P} = (0.5, 0.45, 0.05)$, and (instance 3) $\mathbf{P} = (0.5, 0.45, 0.45)$. The number of the
261 rounds T is fixed to 2000, and we repeated the experiments for 10^5 times.

262 4.1 Training neural networks

263 Here, we show experimental details for training neural networks for the TNN algorithm
264 discussed in Section 3.2.

265 We used the complexity measure $H_1(\mathbf{P}) = \sum_{i \neq i^*(\mathbf{P})} (P^* - P_i)^{-2}$ as a standard choice of $H(\mathbf{P})$.
266 We used the neural network with four layers (including the input layer and output layer),
267 where we used the ReLU for the activation functions and introduced the skip-connection (He
268 et al., 2016) between each hidden layer to make training the network easier. To obtain the
269 map to Δ^K , we adopted the softmax function. The number of nodes in the hidden layers
270 was fixed to $K \times 3$. We used AdamW (Loshchilov and Hutter, 2019) with a learning rate
271 10^{-3} and weight decay 10^{-7} to update the parameters.

272 For training the neural network, we ran Algorithm 3 with $N^{\text{true}} = 32$ and $N^{\text{emp}} = 90$. Addi-
273 tionally, to allow the neural network to easily learn \mathbf{r} , the elements of $\mathbf{P} = (P_1, P_2, \dots, P_K)$
274 were sorted beforehand. Other details of the implementation is given in Appendix C.

275 4.2 Experimental results

276 Figure 1 illustrates the results of our simulations. Each column corresponds to the result for
277 each instance.

278 The first row ((a)–(c)) shows the PoE of the compared methods when the arm with the largest
279 empirical mean is regarded as the estimated best arm $J(t)$ at each round t . Here, the black
280 line represents $\exp(-t \inf_{\mathbf{Q}} \sum_i r_{\theta,i}(\mathbf{Q}) D(Q_i \| P_i))$, which corresponds to the exponent of the
281 oracle algorithm that can perfectly track the allocation $r_{i,\theta}(\mathbf{Q})$. Therefore, the asymptotic
282 slope of TNN cannot be better than that of the black line. We can see from the figures that
283 the slope of the TNN is close to the oracle algorithm and performs better than or comparable
284 to the other algorithms. Note that this is the result for fixed time horizon T . Though the
285 final slope of SR may look outperforming TNN, it just comes from the fact that SR is not
286 anytime and is an algorithm that divides T rounds into several segments.

287 The second row ((d)–(f)) shows the tracking error of the TNN algorithm, which is defined as
288 $\text{disc}(t) = \max_{i \in [K]} |r_i(\mathbf{Q}(t)) - N_i(t)/t|$, which measures the discrepancy between the ideal
289 allocation $r_i(\mathbf{Q}(t))$ and the actual allocation $N_i(t)/t$. If this quantity is $o(T)$ in almost all
290 trials (including the ones where the algorithm failed to recommend the best arm) and all
291 instances, then we can guarantee $R^{\text{go}} = R_{\infty}^{\text{go}}$. The labels TNN (average), TNN (worst), and
292 TNN (average in fail) corresponds to the average tracking error of all trials, the worst-case
293 tracking error and the average tracking error of all failed trials, respectively. The fact that
294 ‘TNN (worst)’ is small at $T = 2,000$ implies that the gap between R^{go} and R_{∞}^{go} is small,
295 which supports the reasonableness of algorithms based on R^{go} .

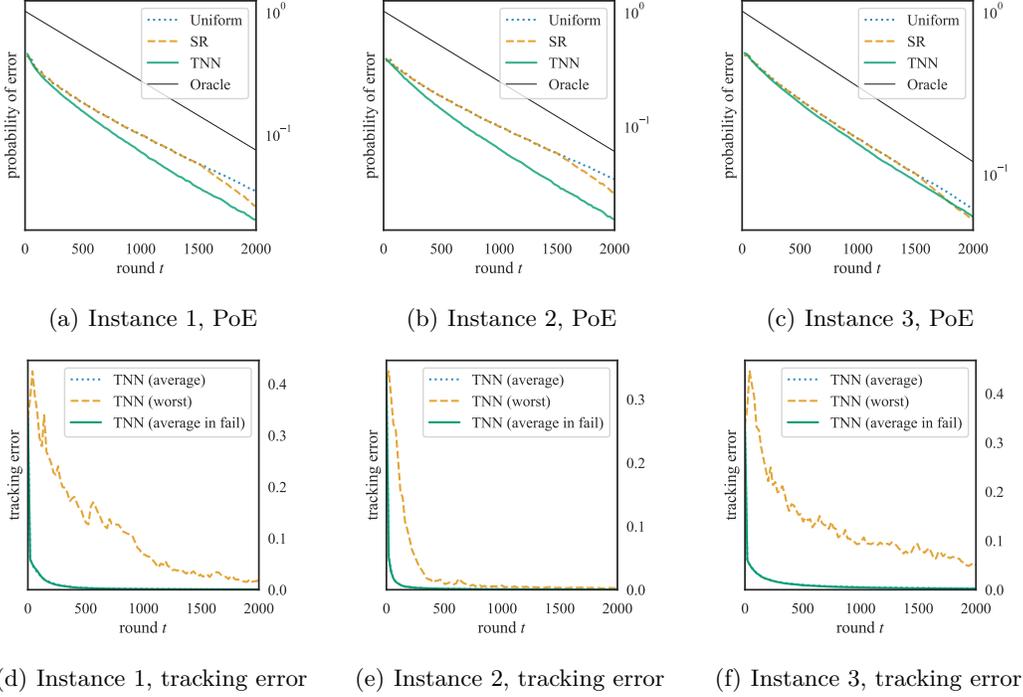


Figure 1: Bernoulli bandits, $K = 3, T = 2000$, average over 10^5 trials.

296 5 Conclusion

297 This paper considered the fixed-budget best arm identification problem. We identified the
 298 minimax rate R_∞^{go} on the exponent of the probability of error by introducing a matching
 299 algorithm (DOT algorithm). Optimization on the rate R_∞^{go} is very challenging to implement,
 300 and we considered the learning of a simpler optimization problem of rate R^{go} by using a
 301 neural network (TNN algorithm). The TNN algorithm outperformed existing algorithms. A
 302 number of possible lines of future work include the following points.

- 303 • A more scalable learning of $\mathbf{r}(\mathbf{Q})$: TNN adopted a neural network to obtain the
 304 oracle allocation $\mathbf{r}(\mathbf{Q})$ associated with the rate bound. While its empirical results are
 305 promising and support our theoretical findings, the current experiment is limited to
 306 the case of $K = 3$ arms because the learning is very costly even for small K . A more
 307 sophisticated learning algorithm is desired to realize R^{go} -tracking for larger K .
- 308 • Identifying the existence (or non-existence) of the gap: though the empirical results
 309 suggest that R^{go} is very close (or maybe equal) to the optimal rate R_∞^{go} for the
 310 Bernoulli case, a formal analysis of this gap for general cases is demanded since the
 311 DOT algorithm to achieve R_∞^{go} is computationally almost infeasible.
- 312 • A bound for another rate measure: we defined the worst-case rate of convergence by
 313 (1), which first takes the limit of T and then takes the worst-case instance \mathbf{P} . Another
 314 natural choice of the rate would be to exchange them, that is, to consider

$$R'(\{\pi_T\}) = \liminf_{T \rightarrow \infty} \inf_{\mathbf{P} \in \mathcal{P}^K} \frac{H(\mathbf{P})}{T} \log(1/\mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})]) \leq R(\{\pi_T\}).$$

315 Whereas Theorems 1 and 2 on the upper bounds of $R(\pi)$ are still valid for $R'(\{\pi_T\}) \leq$
 316 $R(\{\pi_T\})$, the current achievability analysis does not apply and analyzing the tightness
 317 of R_∞^{go} for $R'(\{\pi_T\})$ is an open problem.

318 **References**

- 319 Jean-Yves Audibert, Sébastien Bubeck, and Rémi Munos. Best arm identification in multi-
320 armed bandits. In *COLT 2010 - The 23rd Conference on Learning Theory*, pages 41–53.
321 Omnipress, 2010.
- 322 L. Jeff Hong, Weiwei Fan, and Jun Luo. Review on ranking and selection: A new perspective.
323 *Frontiers of Engineering Management*, 8(3):321–343, Sep 2021.
- 324 Aurélien Garivier and Emilie Kaufmann. Optimal best arm identification with fixed confidence.
325 In *Conference on Learning Theory*, Proceedings of Machine Learning Research, 2016.
- 326 Alexandra Carpentier and Andrea Locatelli. Tight (lower) bounds for the fixed budget best
327 arm identification bandit problem. In *Conference on Learning Theory*, 2016.
- 328 Tor Lattimore and Csaba Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
329 doi: 10.1017/9781108571401.
- 330 Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best-arm
331 identification in multi-armed bandit models. *Journal of Machine Learning Research*, 17
332 (1):1–42, 2016.
- 333 Shahin Shahrampour, Mohammad Noshad, and Vahid Tarokh. On sequential elimination
334 algorithms for best-arm identification in multi-armed bandits. *IEEE Transactions on*
335 *Signal Processing*, 65(16):4281–4292, 2017.
- 336 Victor Gabillon, Mohammad Ghavamzadeh, and Alessandro Lazaric. Best arm identification:
337 A unified approach to fixed budget and fixed confidence. In *Advances in Neural Information*
338 *Processing Systems*, 2012.
- 339 Warren Buckler Powell and Ilya O. Ryzhov. *Optimal Learning*. Second edition. Unpub-
340 lished Manuscript, March 2018. URL [https://castlelab.princeton.edu/wp-content/
341 uploads/2019/02/Powell-OptimalLearningWileyMarch112018.pdf](https://castlelab.princeton.edu/wp-content/uploads/2019/02/Powell-OptimalLearningWileyMarch112018.pdf).
- 342 Chun-Hung Chen, Jianwu Lin, Enver Yücesan, and Stephen E. Chick. Simulation budget
343 allocation for further enhancing the efficiency of ordinal optimization. *Discrete Event*
344 *Dynamic Systems*, 10(3):251–270, July 2000.
- 345 Peter Glynn and Sandeep Juneja. A large deviations perspective on ordinal optimization. In
346 *Winter Simulation Conference*, volume 1. IEEE, 2004.
- 347 Peter I. Frazier, Warren B. Powell, and Savas Dayanik. A knowledge-gradient policy for
348 sequential information collection. *SIAM J. Control Optim.*, 47(5):2410–2439, sep 2008.
- 349 Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley-Interscience,
350 second edition, July 2006. ISBN 0471241954.
- 351 K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016*
352 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778,
353 Los Alamitos, CA, USA, jun 2016. IEEE Computer Society.
- 354 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International*
355 *Conference on Learning Representations*, 2019.

356 **Checklist**

- 357 1. For all authors...
- 358 (a) Do the main claims made in the abstract and introduction accurately reflect
359 the paper’s contributions and scope? [Yes] See Section 1.2
- 360 (b) Did you describe the limitations of your work? [Yes] See Section 5
- 361 (c) Did you discuss any potential negative societal impacts of your work? [N/A]

- 362 (d) Have you read the ethics review guidelines and ensured that your paper conforms
363 to them? [Yes] This is a methodology paper and ethical concerns do not directly
364 apply here.
- 365 2. If you are including theoretical results...
- 366 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See
367 Section 2
- 368 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendices
- 369 3. If you ran experiments...
- 370 (a) Did you include the code, data, and instructions needed to reproduce the main
371 experimental results (either in the supplemental material or as a URL)? [No]
372 We will publish the source code upon acceptance.
- 373 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how
374 they were chosen)? [Yes] See Section 4
- 375 (c) Did you report error bars (e.g., with respect to the random seed after running
376 experiments multiple times)? [Yes] Yes
- 377 (d) Did you include the total amount of compute and the type of resources used
378 (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix B.
- 379 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new
380 assets...
- 381 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 382 (b) Did you mention the license of the assets? [N/A]
- 383 (c) Did you include any new assets either in the supplemental material or as a
384 URL? [N/A]
- 385 (d) Did you discuss whether and how consent was obtained from people whose data
386 you're using/curating? [N/A]
- 387 (e) Did you discuss whether the data you are using/curating contains personally
388 identifiable information or offensive content? [N/A]
- 389 5. If you used crowdsourcing or conducted research with human subjects...
- 390 (a) Did you include the full text of instructions given to participants and screenshots,
391 if applicable? [N/A]
- 392 (b) Did you describe any potential participant risks, with links to Institutional
393 Review Board (IRB) approvals, if applicable? [N/A]
- 394 (c) Did you include the estimated hourly wage paid to participants and the total
395 amount spent on participant compensation? [N/A]

Table 1: Major notation

symbol	definition
K	number of the arms
T	number of the rounds
B	number of the batches
T_B	$= T/(B + K - 1)$
T'	$= T - (B + K - 1)K$
$I(t)$	arm selected at round t
$X(t)$	reward at round t
$J(T)$	recommendation arm at the end of round T
$\mathbf{P} \in \mathcal{P}^K$	true parameters
$P_i \in \mathcal{P}$	i -th component of \mathbf{P}
$\mathcal{I}^* = \mathcal{I}^*(\mathbf{P})$	Set of best arms under parameter \mathbf{P}
$i^*(\mathbf{P})$	one arm in $\mathcal{I}^*(\mathbf{P})$ (taken arbitrary in a deterministic way)
$\mathbf{Q} \in \mathcal{Q}^K$	estimated parameters of \mathbf{P}
$Q_i \in \mathcal{Q}$	i -th component of \mathbf{Q}
$\mathbf{Q}_b \in \mathcal{Q}^K$	estimated parameters of b -th batch
$Q_{b,i} \in \mathcal{Q}$	i -th component of \mathbf{Q}_b
\mathbf{Q}^b	$= (Q_1, Q_2, \dots, Q_b)$
\mathbf{Q}'_b	stored parameters (in Algorithm 2)
$Q'_{b,i} \in \mathcal{Q}$	i -th component of \mathbf{Q}'_b
\mathcal{P}	hypothesis class of \mathbf{P}
\mathcal{Q}	distribution of estimated parameter of \mathbf{Q}
$D(Q\ P)$	KL divergence between Q and P
Δ^K	probability simplex in K dimensions
$\mathbf{r} \in \Delta^K$	allocation (proportion of arm draws)
r_i	i -th component of \mathbf{r}
$\mathbf{r}_b \in \Delta^K$	allocation at b -th batch
$r_{b,i}$	i -th component of \mathbf{r}_b
\mathbf{r}^b	$= (\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_b)$
n_b	Number of draws of Algorithm 2 at b -th batch
$n_{b,i}$	i -th component of n_b . Note that $n_{b,i} \geq r_{b,i}(T_B - K)$ holds.
$J(\mathbf{Q}^B)$	recommendation arm given \mathbf{Q}^B
$(\mathbf{r}^{B,*}, J^*)$	ϵ -optimal allocation
$H(\cdot)$	complexity measure of instances
$R(\{\pi_T\})$	worst-case rate of PoE of sequence of algorithms $\{\pi_T\}$ in (1)
R^{go}	best possible $R(\{\pi_T\})$ for oracle algorithms in (2)
R_B^{go}	best possible $R(\{\pi_T\})$ for B -batch oracle algorithms in (3)
R_∞^{go}	$\lim_{B \rightarrow \infty} R_B^{\text{go}}$. Limit exists (Theorem 7)
$\boldsymbol{\theta}$	model parameter of the neural network
$\mathbf{r}_\boldsymbol{\theta}$	allocation by a neural network with model parameters $\boldsymbol{\theta}$
$r_{\boldsymbol{\theta},i}$	i -th component of $\mathbf{r}_\boldsymbol{\theta}$

396 A Notation table

397 Table 1 summarizes our notation.

398 B Computational resources

399 We used a modern laptop (Macbook Pro) for learning $\boldsymbol{\theta}$. It took less than one hour to learn
400 $\boldsymbol{\theta}$. For conducting a large number of simulations (i.e., Run TNN and existing algorithms for

401 10^5 times), we used a 2-CPU Xeon server of sixteen cores. It took less than twelve hours to
 402 complete simulations. We did not use a GPU for computation.

403 C Implementation details

404 To speed up computation, the same \mathbf{Q} was used for each \mathbf{P} with the same optimal arm $i^*(\mathbf{P})$
 405 in the mini-batches.

406 The final model θ of the neural network is chosen as follows. We stored sequence of models
 407 $\theta^{(1)}, \theta^{(2)}, \dots$ during training (Algorithm 3). Among these models, we chose the one with
 408 the maximum objective function $\arg \max_l \min_{(\mathbf{P}, \mathbf{Q}) \in (\mathcal{P}^{\text{emp}}, \mathcal{Q}^{\text{emp}})} E(\mathbf{P}, \mathbf{Q}; \theta^{(l)})$. Here, the
 409 minimum is taken over a finite dataset of size $|\mathcal{P}^{\text{emp}}| = 32$ and $|\mathcal{Q}^{\text{emp}}| = 10^5$.

410 The black lines in Figure 1 (a)–(c) representing $\exp(-t \inf_{\mathbf{Q}} \sum_i r_{\theta, i}(\mathbf{Q}) D(Q_i \| P_i))$ are com-
 411 puted by the grid search of \mathbf{Q} with each Q_i separated by intervals of 5.0×10^{-3} .

412 D Instance optimality in the fixed-confidence setting

413 For sufficiently small $\delta > 0$, the asymptotic sample complexity for fixed-confidence setting is
 414 known. Namely, any fixed-confidence algorithm is required to draw at least

$$\liminf_{\delta \rightarrow +0} \frac{T}{\log(\delta^{-1})} \geq C^{\text{conf}}(\mathbf{P}) \quad (8)$$

415 times, where

$$C^{\text{conf}}(\mathbf{P}) = \left(\sup_{\mathbf{r}(\mathbf{P}) \in \Delta^K} \inf_{\mathbf{P}': i^*(\mathbf{P}') \notin \mathcal{I}^*(\mathbf{P})} \sum_{i=1}^K r_i D(P_i \| P'_i) \right)^{-1}.$$

416 Garivier and Kaufmann (2016) proposed C -Tracking and D -Tracking algorithms that have
 417 a sample complexity bound that matches Eq. (8). This bound implies that an algorithm
 418 adapts the true parameter \mathbf{P} without paying essential cost of exploration. In fact, building
 419 an optimal algorithm such that Eq. (8) holds is not very difficult.

420 Roughly speaking, a $o(\log(1/\delta))$ cost, say, uniform exploration of $\sqrt{\log(1/\delta)}$ rounds, enables
 421 us to obtain enough accuracy the bound of

$$|\hat{\mathbf{P}} - \mathbf{P}| \sim (\log(1/\delta))^{-1/4} = o(1) \quad (9)$$

422 with probability $1 - o(1)$. The expected value of the stopping time is bounded as:

$$\underbrace{\sqrt{\log(1/\delta)}}_{\text{uniform exploration}} + \underbrace{(C^{\text{conf}}(\mathbf{P}) + o(1)) \log(\delta^{-1})}_{\text{stopping time bound under Eq. (9)}} + \underbrace{o(1)}_{\text{probability of Eq. (9) does not hold}} \times O(\log(\delta^{-1})).$$

423 The first and the third terms does not hurt the optimal rate, and thus the bound of Eq. (8)
 424 is derived.

425 E Extension to wider models

426 In the main body of the paper, we assumed that $P \in \mathcal{P}$ and $Q \in \mathcal{Q}$ are Bernoulli or Gaussian
 427 distributions. Many parts of the results of the paper can be extended to exponential families
 428 or distributions over a support set $\mathcal{S} \subset \mathbb{R}$.

429 Let us consider an exponential family of form

$$dP(x|\theta) = \exp(\theta^\top Y(x) - A(\theta)) dF(x),$$

430 where F is a base measure and $\theta \in \Theta \subset \mathbb{R}^d$ is a natural parameter. We assume that
 431 $A'(\theta) = \mathbb{E}_{X \sim F(\cdot|\theta)}[Y(X)]$ has the inverse $(A')^{-1} : \text{im}(Y) \rightarrow \Theta$, where $\text{im}(Y)$ is the image of
 432 Y .

433 Let \mathcal{P} be a class of reward distributions. \mathcal{P} can be the family of distributions over a known
 434 support $\mathcal{S} \subset \mathbb{R}$. We can also consider the case where \mathcal{P} is the above exponential family
 435 with possibly restricted parameter set $\Theta' \subset \Theta$. For example, \mathcal{P} can be the set of Gaussian
 436 distributions with mean parameters in $[0, 1]$ and variances in $(0, \infty)$.

437 When we derive lower bounds and construct algorithms, we introduce \mathcal{Q} as a class of
 438 distributions corresponding to the estimated reward distributions of the arms. We set $\mathcal{Q} = \mathcal{P}$
 439 when \mathcal{P} is a family of distributions over a known support $\mathcal{S} \subset \mathbb{R}$. When we consider a natural
 440 exponential family with parameter set $\Theta' \subset \Theta$, we set \mathcal{Q} as this exponential family with
 441 parameter set Θ , so that the estimator of P_i is always within \mathcal{Q} . For example, if we consider
 442 \mathcal{P} as a class of Gaussians with means in $[0, 1]$ and variances in $(0, \infty)$, \mathcal{Q} is the class of all
 443 Gaussians with means in $(-\infty, \infty)$ and variances in $(0, \infty)$.

444 In Algorithm 2, we use a convex combination of distributions Q and Q' . The key property
 445 used in the analysis is the convexity of KL divergence between distributions. When we
 446 consider the family \mathcal{P} of distributions over support set \mathcal{S} , the convexity

$$D(\alpha Q + (1 - \alpha)Q' \| P) \leq \alpha D(Q \| P) + (1 - \alpha)D(Q' \| P)$$

447 holds for any $P, Q, Q' \in \mathcal{Q}$ when we define $\alpha Q + (1 - \alpha)Q'$ as the mixture of Q and Q' with
 448 weight $(\alpha, 1 - \alpha)$. When \mathcal{P} is the exponential family, the convexity of the KL divergence holds
 449 when $\alpha Q + (1 - \alpha)Q'$ is defined as the distribution in this family such that the expectation
 450 of the sufficient statistics $Y(X)$ is equal to $\alpha \mathbb{E}_{X \sim Q}[Y(X)] + (1 - \alpha) \mathbb{E}_{X \sim Q'}[Y(X)]$. Note that
 451 this corresponds to taking the convex combination of the empirical means when we consider
 452 Bernoulli distributions or Gaussian distributions with a known variance.

453 By the convexity of the KL divergence, most parts of the analysis apply to \mathcal{P} in this section
 454 and we straightforwardly obtain the following result.

455 **Proposition 8.** Theorems 1 and 2, Corollary 3, and Lemma 4 hold under the models \mathcal{P}
 456 with the definition of the convex combination in this section.

457 The only part where the analysis is limited to Bernoulli or Gaussian is Theorem 5 on the PoE
 458 upper bound of the DOT algorithm. The subsequent results immediately follow if Theorem 5
 459 is extended to the models in this section. Since the key property of the DOT algorithm in
 460 Lemma 4 on the trackability of the empirical divergence is still valid for these models, we
 461 expect that Theorem 5 can also be extended though it remains as an open question.

462 F Proofs

463 F.1 Proofs of Theorems 1 and 2

464 We only give the proof of Theorem 2 since Theorem 1 is a special case of this theorem with
 465 $B = 1$.

466 In this proof, we consider many candidates of the true distributions $\mathbf{P} = (P_1, P_2, \dots, P_K)$
 467 and we write $\mathbf{P}[A]$ to denote the probability of the event A when the reward of each
 468 arm i follows P_i . We divide T rounds into B batches, and the b -th batch corresponds to
 469 $(t_b, t_b + 1, \dots, t_{b+1} - 1)$ -th rounds for $b \in [B]$ and $t_b = \lfloor (b-1)T/B \rfloor + 1$. We define the history
 470 of the b -th batch by $\mathcal{H}_b = ((I(t_b), X(t_b)), (I(t_b + 1), X(t_b + 2)), \dots, (I(t_{b+1} - 1), X(t_{b+1} - 1)))$.
 471 The entire history is denoted by $\mathcal{H}^B = (\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_B)$.

472 By slight abuse of notation, we interchangeably write

$$\mathcal{H}_b = ((X_{b,1,1}, X_{b,1,2}, \dots, X_{b,1,N_{b,1}}), (X_{b,2,1}, X_{b,2,2}, \dots, X_{b,2,N_{b,2}}), \dots, (X_{b,K,1}, X_{b,K,2}, \dots, X_{b,K,N_{b,K}})),$$

473 where $X_{b,k,n}$ is the reward of the n -th draw of arm k in the b -th batch and $N_{b,k}$ is the number
 474 of draws of arm k in the b -th batch.

475 We adopt the formulation of the random rewards such that every $X_{b,k,m}$, the m -th reward
 476 of arm k in the b -th batch, is randomly generated before the game begins, and if an arm is
 477 drawn then this reward is revealed to the player. Then $Y_{b,k,m}$ is well-defined even if arm k is
 478 not drawn m times in the b -th batch.

479 Fix an arbitrary $\epsilon > 0$. We define sets of “typical” rewards under \mathbf{Q}^B : we write $\mathcal{T}_\epsilon(\mathbf{Q}^B)$ to
 480 denote the event such that rewards (a part of which might be unrevealed as noted above)
 481 satisfy

$$\sum_{k=1}^K \left| \left(n_{b,k} D(Q_{b,k} \| P_k) - \sum_{m=1}^{n_{b,k}} \log \frac{dQ_{b,k}}{dP_k}(X_{b,k,m}) \right) \right| \leq \epsilon T/B \quad (10)$$

482 for any $b \in [B]$ and $\mathbf{n}_b = (n_{b,1}, n_{b,2}, \dots, n_{b,K})$ such that $\sum_{k \in [K]} n_{b,k} = t_{b+1} - t_b$. By the
 483 strong law of large numbers, $\lim_{T \rightarrow \infty} \mathbf{Q}^B[\mathcal{T}_\epsilon^B(\mathbf{Q}^B)] = 1$, where $\mathbf{Q}^B[\cdot]$ denotes the probability
 484 under which $X_k(t)$ follows distribution $Q_{b,k}$ for $t \in \{t_b, t_b + 1, \dots, t_{b+1} - 1\}$.

485 We define $\mathbf{r}^B = \mathbf{r}^B(\mathcal{H}^B) = (\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_B)$ for $\mathbf{r}_b = \mathbf{n}_b / (t_{b+1} - t_b)$, where $\mathbf{n}_b =$
 486 $(n_{b,1}, n_{b,2}, \dots, n_{b,K})$. In other words, \mathbf{r}_b is the fractions of arm-draws in the b -th batch
 487 under history \mathcal{H}_b .

488 Let $\mathcal{R}_{T,B} \subset (\Delta^K)^B$ be the set of all possible $\mathbf{r}^B(\mathcal{H}^B)$. Since $n_{b,k} \in \{0, 1, \dots, t_{b+1} - t_b\}$ and
 489 $t_{b+1} - t_b \leq T/B + 1$, we see that

$$|\mathcal{R}_{T,B}| \leq (T/B + 2)^{KB},$$

490 which is polynomial in T .

491 Consider an arbitrary algorithm π and define the “typical” allocation $\mathbf{r}^b(\mathbf{Q}^b; \pi, \epsilon)$ and decision
 492 $J(\mathbf{Q}^b; \pi, \epsilon)$ of the algorithm for distributions $\mathbf{Q}^b = (Q_1, Q_2, \dots, Q_b)$ as

$$\begin{aligned} \mathbf{r}_1(\mathbf{Q}^1; \pi, \epsilon) &= \arg \max_{\mathbf{r} \in \mathcal{R}_{T,1}} \mathbf{Q}^1 [\mathbf{r}_1(\mathcal{H}_1) = \mathbf{r} | \mathcal{T}_\epsilon(\mathbf{Q}^B)], \\ \mathbf{r}_b(\mathbf{Q}^b; \pi, \epsilon) &= \arg \max_{\mathbf{r} \in \mathcal{R}_{T,b}} \mathbf{Q}^b [\mathbf{r}_b(\mathcal{H}_b) = \mathbf{r} | \mathbf{r}^{b-1}(\mathcal{H}^{b-1}) = \mathbf{r}^{b-1}(\mathbf{Q}^{b-1}; \pi, \epsilon), \mathcal{T}_\epsilon(\mathbf{Q}^B)], \\ & \hspace{15em} b = 2, 3, \dots, B, \\ J(\mathbf{Q}^B; \pi, \epsilon) &= \arg \max_{i \in [K]} \mathbf{Q}^B [J(T) = i | \mathbf{r}^B(\mathcal{H}^B) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon), \mathcal{T}_\epsilon(\mathbf{Q}^B)]. \end{aligned}$$

493 Then we have

$$\mathbf{Q}^B [\mathbf{r}^B(\mathcal{H}^B) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon) | \mathcal{T}_\epsilon(\mathbf{Q}^B)] \geq \frac{1}{|\mathcal{R}_{T,B}|}, \quad (11)$$

$$\mathbf{Q}^B [J(T) = J(\mathbf{Q}^B; \pi, \epsilon) | \mathbf{r}^B(\mathcal{H}^B) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon), \mathcal{T}_\epsilon(\mathbf{Q}^B)] \geq \frac{1}{K}. \quad (12)$$

494 **Lemma 9.** Let $\epsilon > 0$ and algorithm π be arbitrary. Then, for any \mathbf{P}, \mathbf{Q}^B be such that
 495 $J(\mathbf{Q}^B; \pi, \epsilon) \neq \mathcal{I}^*(\mathbf{P})$ it holds that

$$\frac{1}{T} \log \mathbf{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \geq -\frac{1}{B} \sum_{b=1}^B \sum_{k=1}^K \mathbf{r}_{b,k}(\mathbf{Q}^b; \pi, \epsilon) D(Q_{b,k} \| P_k) - \epsilon - \delta_{\mathbf{P}, \mathbf{Q}^B, \epsilon}(T)$$

496 for a function $\delta_{\mathbf{P}, \mathbf{Q}^B, \epsilon}(T)$ satisfying $\lim_{T \rightarrow \infty} \delta_{\mathbf{P}, \mathbf{Q}^B, \epsilon}(T) = 0$.

497 *Proof.* For arbitrary \mathbf{Q}^B we obtain by a standard argument of a change of measures that

$$\begin{aligned} & \mathbf{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \\ & \geq \mathbf{P}[\mathcal{T}_\epsilon(\mathbf{Q}^B), \mathbf{r}^B(\mathcal{H}^B) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon), J(T) = J(\mathbf{Q}^B; \pi, \epsilon)] \\ & = \mathbf{P}[\mathcal{T}_\epsilon(\mathbf{Q}^B), \mathbf{r}^B(\mathcal{H}^B) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon)] \\ & \quad \times \mathbf{P}[J(T) = J(\mathbf{Q}^B; \pi, \epsilon) | \mathcal{H}^B \in \mathcal{T}_\epsilon(\mathbf{Q}^B), \mathbf{r}^B(\mathcal{H}^B) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon)] \\ & = \mathbf{P}[\mathcal{T}_\epsilon(\mathbf{Q}^B), \mathbf{r}^B(\mathcal{H}^B) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon)] \\ & \quad \times \mathbf{Q}^B[J(T) = J(\mathbf{Q}^B; \pi, \epsilon) | \mathcal{H}^B \in \mathcal{T}_\epsilon(\mathbf{Q}^B), \mathbf{r}^B(\mathcal{H}^B) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon)] \quad (13) \\ & \geq \frac{1}{K} \mathbf{P}[\mathcal{T}_\epsilon(\mathbf{Q}^B), \mathbf{r}^B(\mathcal{H}^B) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon)] \quad (\text{by (12)}) \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{K} \mathbb{E}_{\mathbf{P}} [\mathbf{1}[\mathcal{H}^B \in \mathcal{T}_\epsilon(\mathbf{Q}^B), \mathbf{r}^B(\mathcal{H}^B) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon)]] \\
&= \frac{1}{K} \mathbb{E}_{\mathbf{Q}^B} \left[\mathbf{1}[\mathcal{T}_\epsilon(\mathbf{Q}^B), \mathbf{r}^B(\mathcal{H}^B) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon)] \prod_{b=1}^B \prod_{t=t_b}^{t_{b+1}-1} \frac{dP_{I(t)}}{dQ_{b,I(t)}}(X(t)) \right] \\
&\geq \frac{1}{K} \mathbb{E}_{\mathbf{Q}^B} [\mathbf{1}[\mathcal{H}^B \in \mathcal{T}_\epsilon(\mathbf{Q}^B), \mathbf{r}^B(\mathcal{H}^B) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon)]] \\
&\quad \times \exp \left(-\frac{T}{B} \sum_{b=1}^B \sum_{k=1}^K r_{b,k}(\mathbf{Q}^b; \pi, \epsilon) D(Q_{b,k} \| P_k) - \epsilon T \right) \quad (\text{by (10)}) \\
&= \frac{1}{K} \mathbf{Q}^B [\mathcal{T}_\epsilon(\mathbf{Q}^B), \mathbf{r}^B(\mathcal{H}^B) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon)] \\
&\quad \times \exp \left(-\frac{T}{B} \sum_{b=1}^B \sum_{k=1}^K r_{b,k}(\mathbf{Q}^b; \pi, \epsilon) D(Q_{b,k} \| P_k) - \epsilon T \right) \\
&\geq \frac{\mathbf{Q}^B[\mathcal{H}^B \in \mathcal{T}_\epsilon(\mathbf{Q}^B)]}{K|\mathcal{R}_{T,B}|} \exp \left(-\frac{T}{B} \sum_{b=1}^B \sum_{k=1}^K r_{b,k}(\mathbf{Q}^b; \pi, \epsilon) D(Q_{b,k} \| P_k) - \epsilon T \right), \quad (\text{by (11)})
\end{aligned}$$

498 where (13) holds since $J(T)$ does not depend on the true distribution \mathbf{P} given the history
499 \mathcal{H}^B . The proof is completed by letting $\delta_{\mathbf{P}, \mathbf{Q}^B, \epsilon} = \log \frac{\mathbf{Q}^B[\mathcal{H}^B \in \mathcal{T}_\epsilon(\mathbf{Q}^B)]}{K|\mathcal{R}_{T,B}|}$. \square

500 *Proof of Theorem 2.* For each \mathbf{Q}^B , let $\mathbf{r}^B(\mathbf{Q}^B; \{\pi_T\}, \epsilon)$, $J(\mathbf{Q}^B; \{\pi_T\}, \epsilon)$ be such that there
501 exists a subsequence $\{T_n\}_n \subset \mathbb{N}$ satisfying

$$\begin{aligned}
\lim_{n \rightarrow \infty} \mathbf{r}^B(\mathbf{Q}^B; \pi^{T_n}, \epsilon) &= \mathbf{r}^B(\mathbf{Q}^B; \{\pi_T\}, \epsilon), \\
J(\mathbf{Q}^B; \pi^{T_n}, \epsilon) &= J(\mathbf{Q}^B; \{\pi_T\}, \epsilon), \quad \forall n.
\end{aligned}$$

502 Such $\mathbf{r}^B(\mathbf{Q}^B; \{\pi_T\}, \epsilon) \in (\Delta^K)^B$ and $J(\mathbf{Q}^B; \{\pi_T\}, \epsilon) \in [K]$ exist since $(\Delta^K)^B$ and $[K]$ are
503 compact. By Lemma 9, for any $J(\mathbf{Q}^B; \{\pi_T\}, \epsilon) \notin \mathcal{I}^*(\mathbf{P})$ we have

$$\begin{aligned}
\liminf_{T \rightarrow \infty} \frac{1}{T} \log 1/\mathbf{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] &\leq \liminf_{n \rightarrow \infty} \frac{1}{T_n} \log 1/\mathbf{P}[J(T_n) \notin \mathcal{I}^*(\mathbf{P})] \\
&\leq \frac{1}{B} \sum_{b=1}^B \sum_{k=1}^K r_{b,k}(\mathbf{Q}^b; \{\pi_T\}, \epsilon) D(Q_{b,k} \| P_k) + \epsilon. \quad (14)
\end{aligned}$$

504 By taking the worst case we have

$$\begin{aligned}
R(\{\pi_T\}) &= \inf_{\mathbf{P}} H(\mathbf{P}) \liminf_{T \rightarrow \infty} \frac{1}{T} \log 1/\mathbf{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \\
&\leq \inf_{\mathbf{P} \in \mathcal{P}^K, \mathbf{Q}^B \in \mathcal{Q}^{KB}: J(\mathbf{Q}^B; \{\pi_T\}, \epsilon) \notin \mathcal{I}^*(\mathbf{P})} \frac{H(\mathbf{P})}{B} \sum_{b=1}^B \sum_{k=1}^K r_{b,k}(\mathbf{Q}^b; \{\pi_T\}, \epsilon) D(Q_{b,k} \| P_k) + \epsilon.
\end{aligned}$$

505 By optimizing $\{\pi^T\}$ we have

$$\begin{aligned}
R(\{\pi_T\}) &\leq \sup_{\{\pi_T\}} \inf_{\mathbf{P} \in \mathcal{P}^K} H(\mathbf{P}) \liminf_{T \rightarrow \infty} \frac{1}{T} \log 1/\mathbf{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \\
&= \sup_{\mathbf{r}^B(\cdot), J(\cdot)} \sup_{\{\pi_T\}: \mathbf{r}^B(\cdot; \{\pi_T\}, \epsilon) = \mathbf{r}^B(\cdot)} \inf_{\mathbf{P} \in \mathcal{P}^K} \frac{H(\mathbf{P})}{B} \liminf_{T \rightarrow \infty} \frac{1}{T} \log 1/\mathbf{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \\
&\leq \sup_{\mathbf{r}^B(\cdot), J(\cdot)} \sup_{\{\pi_T\}: \mathbf{r}^B(\cdot; \{\pi_T\}, \epsilon) = \mathbf{r}^B(\cdot)} \inf_{\mathbf{P} \in \mathcal{P}^K, \mathbf{Q}^B \in \mathcal{Q}^{KB}: J(\mathbf{Q}^B) \notin \mathcal{I}^*(\mathbf{P})} \frac{H(\mathbf{P})}{B} \sum_{b=1}^B \sum_{k=1}^K r_{b,k}(\mathbf{Q}^b) D(Q_{b,k} \| P_k) + \epsilon \\
&\hspace{15em} (\text{by (14)})
\end{aligned}$$

$$\leq \sup_{\mathbf{r}^B(\cdot), J(\cdot)} \inf_{\mathbf{P} \in \mathcal{P}^K, \mathbf{Q}^B \in \mathcal{Q}^{KB}: J(\mathbf{Q}^B) \notin \mathcal{I}^*(\mathbf{P})} \frac{H(\mathbf{P})}{B} \sum_{b=1}^B \sum_{k=1}^K r_{b,k}(\mathbf{Q}^b) D(Q_{b,k} \| P_k) + \epsilon.$$

506 We obtain the desired result since $\epsilon > 0$ is arbitrary. \square

507 **F.2 Proof of Corollary 3**

508 *Proof of Corollary 3.* We have

$$\begin{aligned}
& R_B^{\text{go}} \\
& := \sup_{\mathbf{r}^B(\mathbf{Q}^B), J(\mathbf{Q}^B)} \inf_{\mathbf{Q}^B} \inf_{\mathbf{P}: J(\mathbf{Q}^B) \notin \mathcal{I}^*(\mathbf{P})} \frac{H(\mathbf{P})}{B} \sum_{i \in [K], b \in [B]} r_{b,i} D(Q_{b,i} \| P_i) \\
& \leq \sup_{\mathbf{r}^B(\mathbf{Q}^B), J(\mathbf{Q}^B)} \inf_{\mathbf{Q}^B: \mathbf{Q}_1 = \mathbf{Q}_2 = \dots = \mathbf{Q}_B} \inf_{\mathbf{P}: J(\mathbf{Q}^B) \notin \mathcal{I}^*(\mathbf{P})} \frac{H(\mathbf{P})}{B} \sum_{i \in [K], b \in [B]} r_{b,i} D(Q_{b,i} \| P_i) \quad (\text{inf over a subset}). \\
& = \sup_{\mathbf{r}^B(\mathbf{Q}), J(\mathbf{Q})} \inf_{\mathbf{Q}} \inf_{\mathbf{P}: J(\mathbf{Q}) \notin \mathcal{I}^*(\mathbf{P})} H(\mathbf{P}) \sum_{i \in [K]} \left(\frac{1}{B} \sum_{b \in [B]} r_{b,i} \right) D(Q_i \| P_i) \\
& \quad (\text{by denoting } \mathbf{Q} = \mathbf{Q}_1 = \mathbf{Q}_2 = \dots = \mathbf{Q}_B) \\
& = \sup_{\mathbf{r}(\mathbf{Q}), J(\mathbf{Q})} \inf_{\mathbf{Q}} \inf_{\mathbf{P}: J(\mathbf{Q}) \notin \mathcal{I}^*(\mathbf{P})} H(\mathbf{P}) \sum_{i \in [K]} r_i D(Q_i \| P_i) \\
& \quad (\text{by letting } r_i = (1/B) \sum_b r_{b,i}) \\
& = R^{\text{go}} \quad (\text{by definition}).
\end{aligned}$$

509

□

510 **F.3 Additional Lemma**

511 The following lemma is used to derive the regret bound.

512 **Lemma 10.** Assume that we run Algorithm 2. Then, for any $B_C \in K, K+1, \dots, B$, it
513 follows that

$$\sum_{i, b \in [B_C]} r_{b,i} D(Q_{b,i} \| P_i) \geq \sum_{i, a \in [B_C - K]} r_{a,i}^* D(Q'_{a,i} \| P_i) + \sum_{i \in [K]} D(Q'_{B_C - K + 1, i} \| P_i). \quad (15)$$

514 *Proof of Lemma 10.* We use induction over $B_C = K, K+1, \dots, B$. (i) It is trivial to derive
515 Eq. (15) for $B_C = K$. (ii) Assume that Eq. (15) holds for B_C . In batch $B_C + 1$, the algorithm
516 draws arms in accordance with allocation $\mathbf{r}_{B_C + 1} = \mathbf{r}_{B_C - K + 1}^*$. We have,

$$\begin{aligned}
& \sum_{i \in [K], b \in [B_C + 1]} r_{b,i} D(Q_{b,i} \| P_i) \\
& \geq \sum_{i \in [K], a \in [B_C - K]} r_{a,i}^* D(Q'_{a,i} \| P_i) + \sum_{i \in [K]} D(Q'_{B_C - K + 1, i} \| P_i) + \underbrace{\sum_i r_{B_C + 1, i} D(Q_{B_C + 1, i} \| P_i)}_{\text{Batch } B_C + 1} \\
& \quad (\text{by the assumption of the induction}) \\
& = \sum_i \left(\sum_{a \in [B_C - K]} r_{a,i}^* D(Q'_{a,i} \| P_i) + r_{B_C - K + 1, i}^* D(Q'_{B_C - K + 1, i} \| P_i) \right) + \sum_i (1 - r_{B_C - K + 1, i}^*) D(Q'_{B_C - K + 1, i} \| P_i) \\
& \quad + \sum_i r_{B_C + 1, i} D(Q_{B_C + 1, i} \| P_i) \\
& = \sum_i \left(\sum_{a \in [B_C - K]} r_{a,i}^* D(Q'_{a,i} \| P_i) + r_{B_C - K + 1, i}^* D(Q'_{B_C - K + 1, i} \| P_i) \right) + \sum_i (1 - r_{B_C + 1, i}) D(Q'_{B_C - K + 1, i} \| P_i)
\end{aligned}$$

$$\begin{aligned}
& + \sum_i r_{B_C+1,i} D(Q_{B_C+1,i} \| P_i) \\
& \text{(by definition)} \\
& = \sum_i \left(\sum_{a \in [B_C-K]} r_{a,i}^* D(Q'_{a,i} \| P_i) + r_{B_C-K+1,i}^* D(Q'_{B_C-K+1,i} \| P_i) \right) + \sum_i D(Q'_{B_C-K+2,i} \| P_i) \\
& \text{(by Jensen's inequality and } Q'_{B_C-K+2,i} = r_{B_C+1,i} Q_{B_C+1,i} + (1 - r_{B_C+1,i}) Q'_{B_C-K+1,i} \text{)} \\
& = \sum_i \sum_{a \in [B_C-K+1]} r_{a,i}^* D(Q'_{a,i} \| P_i) + \sum_i D(Q'_{B_C-K+2,i} \| P_i).
\end{aligned}$$

517

□

518 **F.4 Proof of Lemma 4**

Proof of Lemma 4.

$$\begin{aligned}
\sum_{i,b \in [B+K-1]} r_{b,i} D(Q_{b,i} \| P_i) & \geq \sum_{i,b \in [B-1]} r_{b,i}^* D(Q'_{b,i} \| P_i) + \sum_i D(Q'_{B,i} \| P_i). \quad \text{(by (15))} \\
& \geq \sum_{i,b \in [B]} r_{b,i}^* D(Q'_{b,i} \| P_i) \\
& \geq \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})} \quad \text{(by definition of } \epsilon\text{-optimal solution).}
\end{aligned}$$

519

□

520 **F.5 Proof of Theorem 5**

521 *Proof of Theorem 5, Bernoulli rewards.* Since the reward is binary, the possible values that
522 $Q_{b,i}$ take lies in a finite set

$$\mathcal{V} = \left\{ \frac{l}{m} : l \in \mathbb{N}, m \in \mathbb{N}^+ \right\},$$

523 where it is easy to prove $|\mathcal{V}| \leq (T/(B+K-1) + 2)^2 \leq (T/B + 2)^2$. We have

$$\begin{aligned}
\mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] & = \sum_{\mathbf{V}_1, \dots, \mathbf{V}_B \in \mathcal{V}^K} \mathbb{P} \left[J(T) \notin \mathcal{I}^*(\mathbf{P}), \bigcap_b \{Q_b = \mathbf{V}_b\} \right] \\
& = \sum_{\mathbf{V}_1, \dots, \mathbf{V}_B \in \mathcal{V}^K : J^*(\mathbf{V}_1, \dots, \mathbf{V}_B) \notin \mathcal{I}^*(\mathbf{P})} \mathbb{P} \left[\bigcap_b \{Q_b = \mathbf{V}_b\} \right].
\end{aligned}$$

524 By using the Chernoff bound, we have

$$\mathbb{P} \left[Q_{b,i} = V_{b,i} \mid \bigcap_{b' \in [b-1]} \{Q_{b'} = \mathbf{V}_{b'}\} \right] \leq e^{-\frac{T'}{B+K-1} r_{b,i} D(V_{b,i} \| P_i)}, \quad (16)$$

525 and thus

$$\begin{aligned}
& \mathbb{P} \left[\bigcap_b \{Q_b = \mathbf{V}_b\} \right] \\
& = \prod_b \mathbb{P} \left[Q_b = \mathbf{V}_b \mid \bigcap_{b'=1}^{b-1} \{Q_{b'} = \mathbf{V}_{b'}\} \right]
\end{aligned}$$

$$\begin{aligned}
&\leq \prod_b e^{-\frac{T'}{B+K-1} \sum_i r_{b,i} D(V_{b,i} || P_i)} \quad (\text{by Eq. (16)}) \\
&= e^{-\frac{T'}{B+K-1} \sum_{b,i} r_{b,i} D(V_{b,i} || P_i)}. \tag{17}
\end{aligned}$$

526 Furthermore,

$$\begin{aligned}
&\mathbb{P} \left[\bigcap_b \{Q_b = \mathbf{V}_b\} \right] \\
&= \mathbb{P} \left[\bigcap_b \{Q_b = \mathbf{V}_b\}, \sum_{i,b \in [B+K-1]} r_{b,i} D(Q_{b,i} || P_i) \geq \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})} \right] \\
&\quad (\text{by Lemma 4}). \\
&= \mathbb{P} \left[\bigcap_b \{Q_b = \mathbf{V}_b\} \right] \mathbb{P} \left[\sum_{i,b \in [B+K-1]} r_{b,i} D(Q_{b,i} || P_i) \geq \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})} \mid \bigcap_b \{Q_b = \mathbf{V}_b\} \right] \\
&= \mathbb{P} \left[\bigcap_b \{Q_b = \mathbf{V}_b\} \right] \mathbb{P} \left[\sum_{i,b \in [B+K-1]} r_{b,i} D(V_{b,i} || P_i) \geq \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})} \right] \\
&= \mathbb{P} \left[\bigcap_b \{Q_b = \mathbf{V}_b\} \right] \mathbb{E} \left[\mathbf{1} \left[\sum_{i,b \in [B+K-1]} r_{b,i} D(V_{b,i} || P_i) \geq \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})} \right] \right] \\
&\leq e^{-\frac{T'}{B+K-1} \sum_{b,i} r_{b,i} D(V_{b,i} || P_i)} \mathbb{E} \left[\mathbf{1} \left[\sum_{i,b \in [B+K-1]} r_{b,i} D(V_{b,i} || P_i) \geq \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})} \right] \right] \\
&\quad (\text{by Eq. (17)}) \\
&= \mathbb{E} \left[e^{-\frac{T'}{B+K-1} \sum_{b,i} r_{b,i} D(V_{b,i} || P_i)} \mathbf{1} \left[\sum_{i,b \in [B+K-1]} r_{b,i} D(V_{b,i} || P_i) \geq \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})} \right] \right] \\
&\leq \mathbb{E} \left[e^{-\frac{T'}{B+K-1} \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})}} \right] \\
&= e^{-\frac{T'}{B+K-1} \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})}}. \tag{18}
\end{aligned}$$

527 Therefore, we have

$$\begin{aligned}
&\mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \\
&\leq \sum_{\mathbf{V}_1, \dots, \mathbf{V}_B \in \mathcal{V}^K} e^{-\frac{B}{B+K-1} \frac{(R_B^{\text{go}} - \epsilon) T'}{H(\mathbf{P})}} \\
&\quad (\text{by Eq. (18)}) \\
&\leq (T/B + 2)^{2KB} e^{-\frac{B}{B+K-1} \frac{(R_B^{\text{go}} - \epsilon) T'}{H(\mathbf{P})}}.
\end{aligned}$$

528 Here, $\log((T/B + 2)^{2KB}) = o(T)$ to T when we consider K, B as constants.

529

□

530 *Proof of Theorem 5, Normal rewards.* Let

$$\mathcal{B} = \bigcup_{i,b} \{|Q_{b,i}| \geq T\}.$$

531 Then, it is easy to see

$$\mathbb{P}[\mathcal{B}] = T^{2KB} O(e^{-T^2/2}),$$

532 which is negligible because $\log(1/\mathbb{P}[\mathcal{B}])/T$ diverges.

533 The PoE is bounded as

$$\mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] = \mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P}), \mathcal{B}^c] + \mathbb{P}[\mathcal{B}]$$

534 We have,

$$\begin{aligned} & \mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P}), \mathcal{B}^c] \\ &= \int_{-T}^T \cdots \int_{-T}^T \mathbf{1}[J(T) \notin \mathcal{I}^*(\mathbf{P})] p(\mathbf{Q}_B | \mathbf{Q}_{B-1} \dots \mathbf{Q}_1) d\mathbf{Q}_B \dots p(\mathbf{Q}_B | \mathbf{Q}_{B-1} \dots \mathbf{Q}_1) d\mathbf{Q}_B \dots p(\mathbf{Q}_1) d\mathbf{Q}_1. \end{aligned} \quad (19)$$

535 Here,

$$\begin{aligned} p(\mathbf{Q}_B | \mathbf{Q}_{B-1} \dots \mathbf{Q}_1) &= \prod_{i \in [K]} \frac{n_{b,i}}{\sqrt{2\pi}} \exp\left(-\frac{n_{b,i}(Q_{b,i} - P_i)^2}{2}\right) \\ &= \prod_{i \in [K]} \frac{n_{b,i}}{\sqrt{2\pi}} \exp(-n_{b,i}D(Q_{b,i} || P_i)) \\ &\leq \prod_{i \in [K]} T \exp(-n_{b,i}D(Q_{b,i} || P_i)). \end{aligned}$$

536 Finally, we have

$$\begin{aligned} (19) &\leq T^{BK} \int_{-T}^T \cdots \int_{-T}^T \mathbf{1}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \prod_{b \in [B]} \prod_{i \in [K]} \exp(-n_{b,i}D(Q_{b,i} || P_i)) d\mathbf{Q}_B \dots d\mathbf{Q}_1 \\ &\leq T^{BK} \int_{-T}^T \cdots \int_{-T}^T \mathbf{1}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \prod_{b \in [B]} \prod_{i \in [K]} \exp\left(-\frac{T' r^{(b,i)}}{B+K-1} D(Q_{b,i} || P_i)\right) d\mathbf{Q}_B \dots d\mathbf{Q}_1 \\ &\leq T^{BK} \int_{-T}^T \cdots \int_{-T}^T \mathbf{1}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \exp\left(-\frac{B}{B+K-1} \frac{(R_B^{\text{go}} - \epsilon)T'}{H(\mathbf{P})}\right) d\mathbf{Q}_B \dots d\mathbf{Q}_1 \quad (\text{by Lemma 4}) \\ &\leq T^{BK} \int_{-T}^T \cdots \int_{-T}^T \exp\left(-\frac{B}{B+K-1} \frac{(R_B^{\text{go}} - \epsilon)T'}{H(\mathbf{P})}\right) d\mathbf{Q}_B \dots d\mathbf{Q}_1 \\ &\leq T^{BK} (2T)^{BK} \exp\left(-\frac{B}{B+K-1} \frac{(R_B^{\text{go}} - \epsilon)T'}{H(\mathbf{P})}\right). \end{aligned}$$

537

□

538 F.6 Proof of Theorem 7

539 *Proof of Theorem 7.* We first show that the limit

$$R_\infty^{\text{go}} = \lim_{B \rightarrow \infty} R_B^{\text{go}}$$

540 exists. Namely, for any $\eta > 0$ there exists $B_0 \in \mathbb{N}$ such that for any $B_1 > B_0$ we have

$$|R_{B_0}^{\text{go}} - R_{B_1}^{\text{go}}| \leq \eta.$$

541 Theorem 5 implies that Algorithm 2 with $B = B_0$ and $\epsilon = \eta/2$ satisfies⁹

$$\liminf_{T \rightarrow \infty} \frac{\log(1/\mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})])}{T} \geq \frac{B_0}{B_0 + K - 1} \frac{R_{B_0}^{\text{go}} - \eta/2}{H(\mathbf{P})},$$

⁹Strictly speaking, Algorithm 2 depends on T , and we take sequence of the algorithm $(\pi_{\text{DOT}, T})_{T=1,2,\dots}$.

542 and thus

$$\inf H(\mathbf{P}) \liminf_{T \rightarrow \infty} \frac{\log(1/\mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})])}{T} \geq \frac{B_0}{B_0 + K - 1} \left(R_{B_0}^{\text{go}} - \frac{\eta}{2} \right). \quad (20)$$

543 Moreover, Theorem 2 implies that any algorithm satisfies

$$\inf H(\mathbf{P}) \limsup_{T \rightarrow \infty} \frac{\log(1/\mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})])}{T} \leq R_{B_1}^{\text{go}}. \quad (21)$$

544 Combining Eq. (20) and Eq. (21), we have

$$\frac{B_0}{B_0 + K - 1} (R_{B_0}^{\text{go}} - \eta/2) \leq R_{B_1}^{\text{go}}$$

545 and thus

$$\begin{aligned} R_{B_0}^{\text{go}} &\leq R_{B_1}^{\text{go}} + \frac{\eta}{2} + \frac{K-1}{B_0 + K - 1} R_{B_0}^{\text{go}} \\ &\leq R_{B_1}^{\text{go}} + \frac{\eta}{2} + \frac{K-1}{B_0 + K - 1} R_{B_0}^{\text{go}} \quad (\text{by Corollary 3}) \\ &\leq R_{B_1}^{\text{go}} + \frac{\eta}{2} + \frac{\eta}{2} \quad (\text{by } K \geq 2, \text{ by taking } B_0 \geq 2KR_{B_0}^{\text{go}}/\eta) \\ &\leq R_{B_1}^{\text{go}} + \eta. \end{aligned}$$

546 By swapping B_0, B_1 , it is easy to show that

$$R_{B_1}^{\text{go}} \leq R_{B_0}^{\text{go}} + \eta,$$

547 and thus

$$|R_{B_0}^{\text{go}} - R_{B_1}^{\text{go}}| \leq \eta,$$

548 which implies that the limit exists. It is easy to confirm that the performance of Algorithm 2
549 with any $B \geq 2KR_{B_1}^{\text{go}}/\eta$ and $\epsilon = \eta/2$ satisfies Eq. (6). \square