

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
\mathcal{P}	hypothesis class of \mathbf{P}
\mathcal{Q}	distribution of estimated parameter of \mathbf{Q}
$\mathbf{P} \in \mathcal{P}^K$	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
$\mathbf{Q}^b \in \mathcal{Q}^{Kb}$	$= (Q_1, Q_2, \dots, Q_b)$
$\mathbf{Q}'_b \in \mathcal{Q}^K$	stored parameters (in Algorithm 2)
$Q'_{b,i} \in \mathcal{Q}$	i -th component of \mathbf{Q}'_b
$D(Q\ P)$	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	$= (\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_b)$
\mathbf{n}_b	numbers of draws of Algorithm 2 at b -th batch
$n_{b,i}$	i -th component of \mathbf{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)
θ	model parameter of the neural network
\mathbf{r}_θ	allocation by a neural network with model parameters θ
$r_{\theta,i}$	i -th component of \mathbf{r}_θ

A Notation table

Table 1 summarizes our notation.

B Uniform optimality in the fixed-confidence setting

For sufficiently small $\delta > 0$, the asymptotic sample complexity for the FC setting is known.

Namely, any fixed-confidence δ -PAC algorithm require at least $C^{\text{conf}}(\mathbf{P}) \log \delta^{-1} + o(\log \delta^{-1})$ samples, where

$$C^{\text{conf}}(\mathbf{P}) = \left(\sup_{\mathbf{r}(\mathbf{P}) \in \Delta_K} \inf_{\mathbf{P}': i^*(\mathbf{P}') \notin \mathcal{I}^*(\mathbf{P})} \sum_{i=1}^K r_i D(P_i \| P'_i) \right)^{-1}. \quad (8)$$

Garivier and Kaufmann (2016) proposed C -Tracking and D -Tracking algorithms that have a sample complexity bound that matches Eq. (8). This achievability bound implies that there is no tradeoff between the performances for different instances \mathbf{P} , and sacrificing the performance for some \mathbf{P} never improves the performance for another \mathbf{P}' . To be more specific, for example, even if we consider a (δ -correct) algorithm that has a suboptimal sample complexity of $2C^{\text{conf}}(\mathbf{P}) \log \delta^{-1} + o(\log \delta^{-1})$ for some instance \mathbf{P} , it is still impossible to achieve sample complexity better than $C^{\text{conf}}(\mathbf{Q}) \log \delta^{-1} + o(\log \delta^{-1})$ for another instance \mathbf{P}' as far as the algorithm is δ -PAC.

C Suboptimal performance of fixed-confidence algorithms in view of fixed-budget setting

This section shows that an optimal algorithm for the FC-BAI can be arbitrarily bad for the FB-BAI.

For a small $\epsilon \in (0, 0.1)$, consider a three-armed Bernoulli bandit instance with $\mathbf{P}^{(1)} = (0.6, 0.5, 0.5 - \epsilon)$ and $\mathbf{P}^{(2)} = (0.4, 0.5, 0.5 - \epsilon)$. Here, the best arm is arm 1 (resp. arm 2) in the instance $\mathbf{P}^{(1)}$ (resp. $\mathbf{P}^{(2)}$).

Let $\mathbf{r}^{\text{conf}}(\mathbf{P}) = (r_1^{\text{conf}}(\mathbf{P}), r_2^{\text{conf}}(\mathbf{P}), r_3^{\text{conf}}(\mathbf{P}))$ be the optimal FC allocation of Eq. (8). The following characterizes the optimal allocation for $\mathbf{P}^{(1)}, \mathbf{P}^{(2)}$:

Lemma 8. The optimal solution of Eq. (8) for instance $\mathbf{P}^{(1)}$ satisfies the following:

$$r_1^{\text{conf}}(\mathbf{P}^{(1)}), r_2^{\text{conf}}(\mathbf{P}^{(1)}), r_3^{\text{conf}}(\mathbf{P}^{(1)}) \geq 0.07 = \Theta(1).$$

Lemma 9. The optimal solution of Eq. (8) for instance $\mathbf{P}^{(2)}$ satisfies the following:

$$r_1^{\text{conf}}(\mathbf{P}^{(2)}), r_2^{\text{conf}}(\mathbf{P}^{(2)}), r_3^{\text{conf}}(\mathbf{P}^{(2)}) = \Theta(\epsilon^2), \Theta(1), \Theta(1).$$

These two lemmas are derived in Section C.1.

Assume that we run an FC algorithm that draws arms according to allocation $\mathbf{r}^{\text{conf}}(\cdot)$ in an FB problem with T rounds. Under the parameters $\mathbf{P}^{(2)}$, it draws arm 1 for $O(\epsilon^2) + o(T)$ times. Letting $\delta = \mathbf{P}^{(1)}[J(T) = 2]$, Lemma 1 in Kaufmann et al. (2016) implies that

$$\begin{aligned} (TO(\epsilon^2) + o(T))D(0.4\|0.6) &\geq d(\mathbf{P}^{(2)}[J(T) = 2], \mathbf{P}^{(1)}[J(T) = 2]) \\ &\geq d(1/2, \mathbf{P}^{(1)}[J(T) = 2]) \quad (\text{assuming the consistency of algorithm}) \\ &= \frac{1}{2} \left(\log \left(\frac{1}{2\delta} \right) + \log \left(\frac{1}{2(1-\delta)} \right) \right) \\ &\geq \frac{1}{2} \log \left(\frac{1}{2\delta} \right), \end{aligned}$$

which implies

$$\mathbf{P}^{(1)}[J(T) = 2] = \delta \geq \frac{1}{2} \exp(-2(TO(\epsilon^2) + o(T))D(0.4\|0.6)). \quad (9)$$

The exponent of Eq.(9) can be arbitrarily small as $\epsilon \rightarrow +0$. In other words, the rate of this algorithm can be arbitrarily close to 0, while the complexity is $H_1(\mathbf{P}^{(1)}) = \Theta(1)$. This fact implies that the optimal algorithm for the FC-BAI has an arbitrarily bad performance in terms of the minimax rate of the FB-BAI.

C.1 Proofs of Lemmas 8 and 9

Proof of Lemma 8. For $\mathbf{r} = (1/3, 1/3, 1/3)$, we have

$$\begin{aligned} \inf_{\mathbf{P}': i^*(\mathbf{P}') \notin \mathcal{I}^*(\mathbf{P}^{(1)})} \sum_{i=1}^K r_i D(P_i^{(1)} \| P'_i) &> \frac{1}{3} \min(D(0.6 \| 0.55), D(0.5 \| 0.55)) \\ &\text{(by } i^*(\mathbf{P}') \notin \mathcal{I}^*(\mathbf{P}^{(1)}) \text{ implies } P'_1 < 0.55 \text{ or } P'_2 > 0.55 \text{ or } P'_3 > 0.55) \\ &\geq 1/600. \end{aligned}$$

We have

$$\begin{aligned} \inf_{\mathbf{P}': i^*(\mathbf{P}') \notin \mathcal{I}^*(\mathbf{P}^{(1)})} \sum_{i=1}^K r_1^{\text{conf}}(\mathbf{P}^{(1)}) D(P_i \| P'_i) &\leq r_1^{\text{conf}}(\mathbf{P}^{(1)}) D(0.6 \| 0.5) \\ &\text{(on instance } \mathbf{P}' = (0.5, 0.5, 0.5 - \epsilon)) \\ &\leq 0.021 r_1, \end{aligned}$$

which implies $r_1^{\text{conf}}(\mathbf{P}^{(1)}) \geq (1/600) \times (1/0.021) \geq 0.07$ for the optimal allocation $r_1^{\text{conf}}(\mathbf{P}^{(1)})$. Similar discussion yields $r_2, r_3 \geq 0.07$. \square

Proof of Lemma 9. For $\mathbf{r} = (1/3, 1/3, 1/3)$, we have

$$\begin{aligned} \inf_{\mathbf{P}': i^*(\mathbf{P}') \notin \mathcal{I}^*(\mathbf{P}^{(2)})} \sum_{i=1}^K r_i D(P_i^{(2)} \| P'_i) &> \frac{1}{3} \min(D(0.5 \| 0.5 - \epsilon/2), D(0.5 - \epsilon \| 0.5 - \epsilon/2)), \\ &\text{(by } \mathbf{P}' \notin \mathcal{I}^*(\mathbf{P}^{(2)}) \text{ implies } P'_2 < 0.5 - \epsilon/2 \text{ or } P'_1 > 0.5 - \epsilon/2 \text{ or } P'_3 > 0.5 - \epsilon/2) \\ &\geq \frac{\epsilon^2}{6}. \\ &\text{(by Pinsker's inequality)} \end{aligned}$$

We have

$$\begin{aligned} \inf_{\mathbf{P}': i^*(\mathbf{P}') \notin \mathcal{I}^*(\mathbf{P}^{(2)})} \sum_{i=1}^K r_i^{\text{conf}}(\mathbf{P}^{(2)}) D(P_i^{(2)} \| P'_i) &\leq r_2^{\text{conf}}(\mathbf{P}^{(2)}) D(0.5 \| 0.5 - \epsilon/2), \\ &\text{(on instance } \mathbf{P}' = (0.4, 0.5 - \epsilon/2, 0.5 - \epsilon/2)) \end{aligned}$$

which implies $r_i^{\text{conf}}(\mathbf{P}^{(2)}) = \Omega(1)$ for the optimal allocation. Similar discussion yields $r_3^{\text{conf}}(\mathbf{P}^{(2)}) = \Omega(1)$.

In the rest of this proof, we show $r_1^{\text{conf}}(\mathbf{P}^{(2)}) = O(\epsilon^2)$. For the ease of exposition, we drop $(\mathbf{P}^{(2)})$ to denote $\mathbf{r}^{\text{conf}} = (r_1^{\text{conf}}, r_2^{\text{conf}}, r_3^{\text{conf}})$. Lemma 4 in Garivier and Kaufmann (2016) states that the optimal solution satisfies:

$$(r_2^{\text{conf}} + r_1^{\text{conf}}) I_{\frac{r_2^{\text{conf}}}{r_2^{\text{conf}} + r_1^{\text{conf}}}}(P_2^{(2)}, P_1^{(2)}) = (r_2^{\text{conf}} + r_3^{\text{conf}}) I_{\frac{r_2^{\text{conf}}}{r_2^{\text{conf}} + r_3^{\text{conf}}}}(P_2^{(2)}, P_3^{(2)}), \quad (10)$$

where

$$I_\alpha(P_2^{(2)}, P_i^{(2)}) = \alpha D(P_2^{(2)}, \alpha P_2^{(2)} + (1 - \alpha) P_i^{(2)}) + (1 - \alpha) D(P_i^{(2)}, \alpha P_2^{(2)} + (1 - \alpha) P_i^{(2)}).$$

We can confirm that

$$(r_2^{\text{conf}} + r_3^{\text{conf}}) I_{\frac{r_2^{\text{conf}}}{r_2^{\text{conf}} + r_3^{\text{conf}}}}(P_2^{(2)}, P_3^{(2)}) = \Theta(1) \times \Theta(\epsilon^2),$$

and

$$(r_2^{\text{conf}} + r_1^{\text{conf}}) \geq r_2^{\text{conf}} = \Theta(1),$$

which, combined with Eq.(10), implies that

$$I_{\frac{r_2^{\text{conf}}}{r_2^{\text{conf}} + r_1^{\text{conf}}}}(P_2^{(2)}, P_1^{(2)}) = \Theta(\epsilon^2),$$

which implies $r_1^{\text{conf}} = \Theta(\epsilon^2)$. □

D Extension to wider models

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

Let us consider an exponential family of form

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

where F is a base measure and $\theta \in \Theta \subset \mathbb{R}^d$ is a natural parameter. We assume that $A'(\theta) = \mathbb{E}_{X \sim F(\cdot|\theta)}[T(X)]$ has the inverse $(A')^{-1} : \text{im}(T) \rightarrow \Theta$, where $\text{im}(T)$ is the image of T .

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

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

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

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

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 weight $(\alpha, 1 - \alpha)$. When \mathcal{P} is the exponential family, the convexity of the KL divergence holds when $\alpha Q + (1 - \alpha)Q'$ is defined as the distribution in this family such that the expectation of the sufficient statistics $T(X)$ is equal to $\alpha \mathbb{E}_{X \sim Q}[T(X)] + (1 - \alpha) \mathbb{E}_{X \sim Q'}[T(X)]$. Note that this corresponds to taking the convex combination of the empirical means when we consider Bernoulli distributions or Gaussian distributions with a known variance.

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

Proposition 10. Theorems 1 and 2, Corollary 3, and Lemma 4 hold under the models \mathcal{P} with the definition of the convex combination in this section.

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

E Computational resources

We used a modern laptop (Macbook Pro) for learning θ . It took less than one hour to learn θ . For conducting a large number of simulations (i.e., Run TNN and existing algorithms for

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

F Implementation details

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

The final model θ of the neural network is chosen as follows. We stored sequence of models $\theta^{(1)}, \theta^{(2)}, \dots$ during training (Algorithm 3). Among these models, we chose the one with 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 minimum is taken over a finite dataset of size $|\mathcal{P}^{\text{emp}}| = 32$ and $|\mathcal{Q}^{\text{emp}}| = 10^5$.

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 computed by the grid search of \mathbf{Q} with each Q_i separated by intervals of 5.0×10^{-3} .

G Proofs

G.1 Proofs of Theorems 1

In this section, we prove Theorem 1. This theorem as well as its proof is a special case of Theorem 2, but we solely prove Theorem 1 here since it is easier to follow.

In this proof, we write candidates of the true distributions and empirical distributions by $\mathbf{P} = (P_1, P_2, \dots, P_K)$ and $\mathbf{Q} = (Q_1, Q_2, \dots, Q_K)$, respectively. In this Sections G.1 and G.2, we write $\mathbf{P}[\mathcal{A}]$ and $\mathbf{Q}[\mathcal{A}]$ to denote the probability of the event \mathcal{A} when the reward of each arm i follows P_i and Q_i , respectively. The entire history of the drawn arms and observed rewards is denoted by $\mathcal{H} = ((I(1), X(1)), (I(2), X(2)), \dots, (I(T), X(T)))$. We write $X_{i,n}$ to denote the reward of the n -th draw of arm i . We define $\mathbf{n} = (n_1, n_2, \dots, n_K)$ and $\mathbf{r} = (r_1, r_2, \dots, r_K) = \mathbf{n}/T$ as the numbers of draws of K arms and their fractions, respectively, for which we write $\mathbf{n}(\mathcal{H})$ and $\mathbf{r}(\mathcal{H})$ when we emphasize the dependence on the history \mathcal{H} .

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

Fix an arbitrary $\epsilon > 0$. We define sets of “typical” rewards under \mathbf{Q} : we write $\mathcal{T}_\epsilon(\mathbf{Q})$ to denote the event such that the rewards (some of which might not be revealed as noted above) satisfy

$$\sum_{i=1}^K \left| \left(n_i D(Q_i \| P_i) - \sum_{m=1}^{n_i} \log \frac{dQ_i}{dP_i}(X_{i,m}) \right) \right| \leq \epsilon T. \quad (11)$$

By the strong law of large numbers, $\lim_{T \rightarrow \infty} \mathbf{Q}[\mathcal{T}_\epsilon(\mathbf{Q})] = 1$.

Let $\mathcal{R}_T \subset \Delta_K$ be the set of all possible $\mathbf{r} = \mathbf{n}/T$. Since $n_i \in \{0, 1, \dots, T\}$ we have

$$|\mathcal{R}_T| \leq (T+1)^K,$$

which is polynomial in T .

Consider an arbitrary algorithm π and define the “typical” allocation $\mathbf{r}(\mathbf{Q}; \pi, \epsilon)$ and decision $J(\mathbf{Q}; \pi, \epsilon)$ of the algorithm for distributions \mathbf{Q} as

$$\begin{aligned} \mathbf{r}(\mathbf{Q}; \pi, \epsilon) &= \arg \max_{\mathbf{r} \in \mathcal{R}_T} \mathbf{Q}[\mathbf{r}(\mathcal{H}) = \mathbf{r} | \mathcal{T}_\epsilon(\mathbf{Q})], \\ J(\mathbf{Q}; \pi, \epsilon) &= \arg \max_{i \in [K]} \mathbf{Q}[J(T) = i | \mathbf{r}(\mathcal{H}) = \mathbf{r}(\mathbf{Q}; \pi, \epsilon), \mathcal{T}_\epsilon(\mathbf{Q})]. \end{aligned}$$

Then we have

$$\mathbf{Q}[\mathbf{r}(\mathcal{H}) = \mathbf{r}(\mathbf{Q}; \pi, \epsilon) | \mathcal{T}_\epsilon(\mathbf{Q})] \geq \frac{1}{|\mathcal{R}_T|}, \quad (12)$$

$$\mathbf{Q} \left[J(T) = J(\mathbf{Q}; \pi, \epsilon) \mid \mathbf{r}(\mathcal{H}) = \mathbf{r}(\mathbf{Q}; \pi, \epsilon), \mathcal{T}_\epsilon(\mathbf{Q}) \right] \geq \frac{1}{K}. \quad (13)$$

Lemma 11. Let $\epsilon > 0$ and algorithm π be arbitrary. Then, for any \mathbf{P}, \mathbf{Q} such that $J(\mathbf{Q}; \pi, \epsilon) \notin \mathcal{I}^*(\mathbf{P})$ it holds that

$$\frac{1}{T} \log \mathbf{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \geq - \sum_{i=1}^K r_i(\mathbf{Q}; \pi, \epsilon) D(Q_i \| P_i) - \epsilon - \delta_{\mathbf{P}, \mathbf{Q}, \epsilon}(T)$$

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

Proof. For arbitrary \mathbf{Q} 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}), \mathbf{r}(\mathcal{H}) = \mathbf{r}(\mathbf{Q}; \pi, \epsilon), J(T) = J(\mathbf{Q}; \pi, \epsilon)] \\ & = \mathbf{P}[\mathcal{T}_\epsilon(\mathbf{Q}), \mathbf{r}(\mathcal{H}) = \mathbf{r}(\mathbf{Q}; \pi, \epsilon)] \mathbf{P}[J(T) = J(\mathbf{Q}; \pi, \epsilon) \mid \mathcal{T}_\epsilon(\mathbf{Q}), \mathbf{r}(\mathcal{H}) = \mathbf{r}(\mathbf{Q}; \pi, \epsilon)] \\ & = \mathbf{P}[\mathcal{T}_\epsilon(\mathbf{Q}), \mathbf{r}(\mathcal{H}) = \mathbf{r}(\mathbf{Q}; \pi, \epsilon)] \mathbf{Q}[J(T) = J(\mathbf{Q}; \pi, \epsilon) \mid \mathcal{T}_\epsilon(\mathbf{Q}), \mathbf{r}(\mathcal{H}) = \mathbf{r}(\mathbf{Q}; \pi, \epsilon)] \quad (14) \\ & \geq \frac{1}{K} \mathbf{P}[\mathcal{T}_\epsilon(\mathbf{Q}), \mathbf{r}(\mathcal{H}) = \mathbf{r}(\mathbf{Q}; \pi, \epsilon)] \quad (\text{by (13)}) \\ & = \frac{1}{K} \mathbb{E}_{\mathbf{P}} [\mathbf{1}[\mathcal{H} \in \mathcal{T}_\epsilon(\mathbf{Q}), \mathbf{r}(\mathcal{H}) = \mathbf{r}(\mathbf{Q}; \pi, \epsilon)]] \\ & = \frac{1}{K} \mathbb{E}_{\mathbf{Q}} \left[\mathbf{1}[\mathcal{T}_\epsilon(\mathbf{Q}), \mathbf{r}(\mathcal{H}) = \mathbf{r}(\mathbf{Q}; \pi, \epsilon)] \prod_{t=1}^T \frac{dP_{I(t)}}{dQ_{I(t)}}(X(t)) \right] \\ & \geq \frac{1}{K} \mathbb{E}_{\mathbf{Q}} [\mathbf{1}[\mathcal{H} \in \mathcal{T}_\epsilon(\mathbf{Q}), \mathbf{r}(\mathcal{H}) = \mathbf{r}(\mathbf{Q}; \pi, \epsilon)]] \exp \left(-T \sum_{i=1}^K r_{b,i}(\mathbf{Q}; \pi, \epsilon) D(Q_i \| P_i) - \epsilon T \right) \\ & \quad (\text{by (11)}) \\ & = \frac{1}{K} \mathbf{Q}[\mathcal{T}_\epsilon(\mathbf{Q}), \mathbf{r}(\mathcal{H}) = \mathbf{r}(\mathbf{Q}; \pi, \epsilon)] \exp \left(-T \sum_{i=1}^K r_i(\mathbf{Q}; \pi, \epsilon) D(Q_i \| P_i) - \epsilon T \right) \\ & \geq \frac{\mathbf{Q}[\mathcal{T}_\epsilon(\mathbf{Q})]}{K |\mathcal{R}_T|} \exp \left(-T \sum_{i=1}^K r_i(\mathbf{Q}; \pi, \epsilon) D(Q_i \| P_i) - \epsilon T \right), \quad (\text{by (12)}) \end{aligned}$$

where (14) holds since $J(T)$ does not depend on the true distribution \mathbf{P} given the history \mathcal{H} . The proof is completed by letting $\delta_{\mathbf{P}, \mathbf{Q}, \epsilon} = \log \frac{\mathbf{Q}[\mathcal{H} \in \mathcal{T}_\epsilon(\mathbf{Q})]}{K |\mathcal{R}_T|}$. \square

Proof of Theorem 1. For each \mathbf{Q} , let $\mathbf{r}(\mathbf{Q}; \{\pi_T\}, \epsilon)$, $J(\mathbf{Q}; \{\pi_T\}, \epsilon)$ be such that there exists a subsequence $\{T_n\}_n \subset \mathbb{N}$ satisfying

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

Such $\mathbf{r}(\mathbf{Q}; \{\pi_T\}, \epsilon) \in \Delta_K$ and $J(\mathbf{Q}; \{\pi_T\}, \epsilon) \in [K]$ exist since Δ_K and $[K]$ are compact. By Lemma 11, for any $J(\mathbf{Q}; \{\pi_T\}, \epsilon) \notin \mathcal{I}^*(\mathbf{P})$ we have

$$\begin{aligned} \liminf_{T \rightarrow \infty} \frac{1}{T} \log 1/\mathbf{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] &\leq \liminf_{n \rightarrow \infty} \frac{1}{T_n} \log 1/\mathbf{P}[J(T_n) \notin \mathcal{I}^*(\mathbf{P})] \\ &\leq \sum_{i=1}^K r_i(\mathbf{Q}; \{\pi_T\}, \epsilon) D(Q_i \| P_i) + \epsilon. \quad (15) \end{aligned}$$

By taking the worst case we have

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

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

$$\begin{aligned}
R(\{\pi_T\}) &\leq \sup_{\{\pi_T\}} \inf_{\mathbf{P} \in \mathcal{P}^K} H(\mathbf{P}) \liminf_{T \rightarrow \infty} \frac{1}{T} \log 1/\mathbf{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \\
&= \sup_{\mathbf{r}(\cdot), J(\cdot)} \sup_{\{\pi_T\}: \mathbf{r}(\cdot; \{\pi_T\}, \epsilon) = \mathbf{r}(\cdot)} \inf_{\mathbf{P} \in \mathcal{P}^K} H(\mathbf{P}) \liminf_{T \rightarrow \infty} \frac{1}{T} \log 1/\mathbf{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \\
&\leq \sup_{\mathbf{r}(\cdot), J(\cdot)} \sup_{\{\pi_T\}: \mathbf{r}(\cdot; \{\pi_T\}, \epsilon) = \mathbf{r}(\cdot)} \inf_{\mathbf{P} \in \mathcal{P}^K, \mathbf{Q} \in \mathcal{Q}^K: J(\mathbf{Q}) \notin \mathcal{I}^*(\mathbf{P})} H(\mathbf{P}) \sum_{i=1}^K r_i(\mathbf{Q}) D(Q_i \| P_i) + \epsilon \\
&\hspace{20em} \text{(by (15))} \\
&\leq \sup_{\mathbf{r}(\cdot), J(\cdot)} \inf_{\mathbf{P} \in \mathcal{P}^K, \mathbf{Q} \in \mathcal{Q}^K: J(\mathbf{Q}) \notin \mathcal{I}^*(\mathbf{P})} H(\mathbf{P}) \sum_{i=1}^K r_i(\mathbf{Q}) D(Q_i \| P_i) + \epsilon.
\end{aligned}$$

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

G.2 Proof of Theorem 2

Theorem 2 is a generalization of Theorem 1, and we consider different candidates of empirical distributions depending on the batch.

As in the case of the proof of Theorem 1, we write $\mathbf{P} = (P_1, P_2, \dots, P_i)$ and $\mathbf{P}[A]$ to denote a candidate of the true distributions and the probability of the event under \mathbf{P} . We divide T rounds into B batches, and the b -th batch corresponds to $(t_b, t_b + 1, \dots, t_{b+1} - 1)$ -th rounds for $b \in [B]$ and $t_b = \lfloor (b-1)T/B \rfloor + 1$. The entire history of the drawn arms and observed rewards is denoted by $\mathcal{H} = ((I(1), X(1)), (I(2), X(2)), \dots, (I(T), X(T)))$. We write $X_{b,i,n}$ to denote the reward of the n -th draw of arm i in the b -th batch. We define $\mathbf{n}_b = (n_{b,1}, n_{b,2}, \dots, n_{b,K})$ and $\mathbf{r} = (r_{b,1}, r_{b,2}, \dots, r_{b,K}) = \mathbf{n}_b/T$ as the numbers of draws of K arms and their fractions in the b -th batch, respectively, for which we write $\mathbf{n}_b(\mathcal{H})$ and $\mathbf{r}_b(\mathcal{H})$ when we emphasize the dependence on the history \mathcal{H} .

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

Fix an arbitrary $\epsilon > 0$. We define sets of ‘‘typical’’ rewards under \mathbf{Q}^B : we write $\mathcal{T}_\epsilon(\mathbf{Q}^B)$ to denote the event such that the rewards (a part of which might be unrevealed as noted above) satisfy

$$\sum_{i=1}^K \left| \left(n_{b,i} D(Q_{b,i} \| P_i) - \sum_{m=1}^{n_{b,i}} \log \frac{dQ_{b,i}}{dP_i}(X_{b,i,m}) \right) \right| \leq \epsilon T/B \quad (16)$$

for any $b \in [B]$. By the 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 under which $X_k(t)$ follows distribution $Q_{b,i}$ for $t \in \{t_b, t_b + 1, \dots, t_{b+1} - 1\}$.

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

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

which is polynomial in T .

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

$$\begin{aligned}
\mathbf{r}_1(\mathbf{Q}^1; \pi, \epsilon) &= \arg \max_{\mathbf{r} \in \mathcal{R}_{T,1}} \mathbf{Q}^1 [\mathbf{r}_1(\mathcal{H}) = \mathbf{r} | \mathcal{T}_\epsilon(\mathbf{Q}^B)], \\
\mathbf{r}_b(\mathbf{Q}^b; \pi, \epsilon) &= \arg \max_{\mathbf{r} \in \mathcal{R}_{T,b}} \mathbf{Q}^b [\mathbf{r}_b(\mathcal{H}) = \mathbf{r} | \mathbf{r}^{b-1}(\mathcal{H}^{b-1}) = \mathbf{r}^{b-1}(\mathbf{Q}^{b-1}; \pi, \epsilon), \mathcal{T}_\epsilon(\mathbf{Q}^B)],
\end{aligned}$$

$$b = 2, 3, \dots, B,$$

$$J(\mathbf{Q}^B; \pi, \epsilon) = \arg \max_{i \in [K]} \mathbf{Q}^B \left[J(T) = i \mid \mathbf{r}^B(\mathcal{H}) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon), \mathcal{T}_\epsilon(\mathbf{Q}^B) \right].$$

Then we have

$$\mathbf{Q}^B \left[\mathbf{r}^B(\mathcal{H}) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon) \mid \mathcal{T}_\epsilon(\mathbf{Q}^B) \right] \geq \frac{1}{|\mathcal{R}_{T,B}|}, \quad (17)$$

$$\mathbf{Q}^B \left[J(T) = J(\mathbf{Q}^B; \pi, \epsilon) \mid \mathbf{r}^B(\mathcal{H}) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon), \mathcal{T}_\epsilon(\mathbf{Q}^B) \right] \geq \frac{1}{K}. \quad (18)$$

Lemma 12. Let $\epsilon > 0$ and algorithm π be arbitrary. Then, for any \mathbf{P}, \mathbf{Q}^B such that $J(\mathbf{Q}^B; \pi, \epsilon) \neq \mathcal{I}^*(\mathbf{P})$ it holds that

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

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

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}) = \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}) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon)] \\ & \quad \times \mathbf{P}[J(T) = J(\mathbf{Q}^B; \pi, \epsilon) \mid \mathcal{T}_\epsilon(\mathbf{Q}^B), \mathbf{r}^B(\mathcal{H}) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon)] \\ & = \mathbf{P}[\mathcal{T}_\epsilon(\mathbf{Q}^B), \mathbf{r}^B(\mathcal{H}) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon)] \\ & \quad \times \mathbf{Q}^B[J(T) = J(\mathbf{Q}^B; \pi, \epsilon) \mid \mathcal{T}_\epsilon(\mathbf{Q}^B), \mathbf{r}^B(\mathcal{H}) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon)] \quad (19) \\ & \geq \frac{1}{K} \mathbf{P}[\mathcal{T}_\epsilon(\mathbf{Q}^B), \mathbf{r}^B(\mathcal{H}) = \mathbf{r}^B(\mathbf{Q}^B; \pi, \epsilon)] \quad (\text{by (18)}) \\ & = \frac{1}{K} \mathbb{E}_{\mathbf{P}} [\mathbf{1}[\mathcal{H} \in \mathcal{T}_\epsilon(\mathbf{Q}^B), \mathbf{r}^B(\mathcal{H}) = \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}) = \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} \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_{i=1}^K r_{b,i}(\mathbf{Q}^b; \pi, \epsilon) D(Q_{b,i} \| P_i) - \epsilon T \right) \quad (\text{by (16)}) \\ & = \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_{i=1}^K r_{b,i}(\mathbf{Q}^b; \pi, \epsilon) D(Q_{b,i} \| P_i) - \epsilon T \right) \\ & \geq \frac{\mathbf{Q}^B[\mathcal{T}_\epsilon(\mathbf{Q}^B)]}{K |\mathcal{R}_{T,B}|} \exp \left(-\frac{T}{B} \sum_{b=1}^B \sum_{i=1}^K r_{b,i}(\mathbf{Q}^b; \pi, \epsilon) D(Q_{b,i} \| P_i) - \epsilon T \right), \quad (\text{by (17)}) \end{aligned}$$

where (19) holds since $J(T)$ does not depend on the true distribution \mathbf{P} given the history \mathcal{H} . The proof is completed by letting $\delta_{\mathbf{P}, \mathbf{Q}^B, \epsilon} = \log \frac{\mathbf{Q}^B[\mathcal{T}_\epsilon(\mathbf{Q}^B)]}{K |\mathcal{R}_{T,B}|}$. \square

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 exists a subsequence $\{T_n\}_n \subset \mathbb{N}$ satisfying

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

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 compact. By Lemma 12, for any $J(\mathbf{Q}^B; \{\pi_T\}, \epsilon) \notin \mathcal{I}^*(\mathbf{P})$ we have

$$\begin{aligned} \liminf_{T \rightarrow \infty} \frac{1}{T} \log 1/\mathbf{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] &\leq \liminf_{n \rightarrow \infty} \frac{1}{T_n} \log 1/\mathbf{P}[J(T_n) \notin \mathcal{I}^*(\mathbf{P})] \\ &\leq \frac{1}{B} \sum_{b=1}^B \sum_{i=1}^K r_{b,i}(\mathbf{Q}^b; \{\pi_T\}, \epsilon) D(Q_{b,i} \| P_i) + \epsilon. \end{aligned} \quad (20)$$

By taking the worst case we have

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

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

$$\begin{aligned} R(\{\pi_T\}) &\leq \sup_{\{\pi_T\}} \inf_{\mathbf{P} \in \mathcal{P}^K} H(\mathbf{P}) \liminf_{T \rightarrow \infty} \frac{1}{T} \log 1/\mathbf{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \\ &= \sup_{\mathbf{r}^B(\cdot), J(\cdot)} \sup_{\{\pi_T\}: \mathbf{r}^B(\cdot; \{\pi_T\}, \epsilon) = \mathbf{r}^B(\cdot)} \inf_{\mathbf{P} \in \mathcal{P}^K} \frac{H(\mathbf{P})}{B} \liminf_{T \rightarrow \infty} \frac{1}{T} \log 1/\mathbf{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \\ &\leq \sup_{\mathbf{r}^B(\cdot), J(\cdot)} \sup_{\{\pi_T\}: \mathbf{r}^B(\cdot; \{\pi_T\}, \epsilon) = \mathbf{r}^B(\cdot)} \inf_{\mathbf{P} \in \mathcal{P}^K, \mathbf{Q}^B \in \mathcal{Q}^{KB}: J(\mathbf{Q}^B) \notin \mathcal{I}^*(\mathbf{P})} \frac{H(\mathbf{P})}{B} \sum_{b=1}^B \sum_{i=1}^K r_{b,i}(\mathbf{Q}^b) D(Q_{b,i} \| P_i) + \epsilon \\ &\hspace{15em} \text{(by (20))} \\ &\leq \sup_{\mathbf{r}^B(\cdot), J(\cdot)} \inf_{\mathbf{P} \in \mathcal{P}^K, \mathbf{Q}^B \in \mathcal{Q}^{KB}: J(\mathbf{Q}^B) \notin \mathcal{I}^*(\mathbf{P})} \frac{H(\mathbf{P})}{B} \sum_{b=1}^B \sum_{i=1}^K r_{b,i}(\mathbf{Q}^b) D(Q_{b,i} \| P_i) + \epsilon. \end{aligned}$$

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

G.3 Proof of Corollary 3

Proof of Corollary 3. We have

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

\square

G.4 Additional lemmas

The following lemma is used to derive the regret bound.

Lemma 13. Assume that we run Algorithm 2. Then, for any $B_C \in K, K+1, \dots, B$, it follows that

$$\sum_{i,b \in [B_C]} r_{b,i} D(Q_{b,i} \| P_i) \geq \sum_{i,a \in [B_C-K]} r_{a,i}^* D(Q'_{a,i} \| P_i) + \sum_{i \in [K]} D(Q'_{B_C-K+1,i} \| P_i). \quad (21)$$

Proof of Lemma 13. We use induction over $B_C \geq K$. (i) It is trivial to derive Eq. (21) for $B_C = K$. (ii) Assume that Eq. (21) holds for B_C . In batch $B_C + 1$, the algorithm draws arms in accordance with allocation $\mathbf{r}_{B_C+1} = \mathbf{r}_{B_C-K+1}^*$. We have,

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

□

G.5 Proof of Lemma 4

Proof of Lemma 4.

$$\begin{aligned} \sum_{i,b \in [B+K-1]} r_{b,i} D(Q_{b,i} \| P_i) & \geq \sum_{i,b \in [B-1]} r_{b,i}^* D(Q'_{b,i} \| P_i) + \sum_i D(Q'_{B,i} \| P_i). \quad (\text{by (21)}) \\ & \geq \sum_{i,b \in [B]} r_{b,i}^* D(Q'_{b,i} \| P_i) \\ & \geq \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})} \quad (\text{by definition of } \epsilon\text{-optimal solution}). \end{aligned}$$

□

G.6 Proof of Theorem 5

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

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

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

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

By using the Chernoff bound, we have

$$\mathbb{P} \left[Q_{b,i} = V_{b,i} \mid \bigcap_{b' \in [b-1]} \{Q_{b'} = \mathbf{V}_{b'}\} \right] \leq e^{-\frac{T'}{B+K-1} r_{b,i} D(V_{b,i} \| P_i)}, \quad (22)$$

and thus

$$\begin{aligned} &\mathbb{P} \left[\bigcap_b \{Q_b = \mathbf{V}_b\} \right] \\ &= \prod_b \mathbb{P} \left[Q_b = \mathbf{V}_b \mid \bigcap_{b'=1}^{b-1} \{Q_{b'} = \mathbf{V}_{b'}\} \right] \\ &\leq \prod_b e^{-\frac{T'}{B+K-1} \sum_i r_{b,i} D(V_{b,i} \| P_i)} \quad (\text{by Eq. (22)}) \\ &= e^{-\frac{T'}{B+K-1} \sum_{b,i} r_{b,i} D(V_{b,i} \| P_i)}. \end{aligned} \quad (23)$$

Furthermore,

$$\begin{aligned} &\mathbb{P} \left[\bigcap_b \{Q_b = \mathbf{V}_b\} \right] \\ &= \mathbb{P} \left[\bigcap_b \{Q_b = \mathbf{V}_b\}, \sum_{i, b \in [B+K-1]} r_{b,i} D(Q_{b,i} \| P_i) \geq \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})} \right] \\ &\quad (\text{by Lemma 4}). \\ &= \mathbb{P} \left[\bigcap_b \{Q_b = \mathbf{V}_b\} \right] \mathbb{P} \left[\sum_{i, b \in [B+K-1]} r_{b,i} D(Q_{b,i} \| P_i) \geq \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})} \mid \bigcap_b \{Q_b = \mathbf{V}_b\} \right] \\ &= \mathbb{P} \left[\bigcap_b \{Q_b = \mathbf{V}_b\} \right] \mathbb{P} \left[\sum_{i, b \in [B+K-1]} r_{b,i} D(V_{b,i} \| P_i) \geq \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})} \right] \\ &= \mathbb{P} \left[\bigcap_b \{Q_b = \mathbf{V}_b\} \right] \mathbb{E} \left[\mathbf{1} \left[\sum_{i, b \in [B+K-1]} r_{b,i} D(V_{b,i} \| P_i) \geq \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})} \right] \right] \end{aligned}$$

$$\begin{aligned}
&\leq e^{-\frac{T'}{B+K-1} \sum_{b,i} r_{b,i} D(V_{b,i}||P_i)} \mathbb{E} \left[\mathbf{1} \left[\sum_{i,b \in [B+K-1]} r_{b,i} D(V_{b,i}||P_i) \geq \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})} \right] \right] \\
&\quad (\text{by Eq. (23)}) \\
&= \mathbb{E} \left[e^{-\frac{T'}{B+K-1} \sum_{b,i} r_{b,i} D(V_{b,i}||P_i)} \mathbf{1} \left[\sum_{i,b \in [B+K-1]} r_{b,i} D(V_{b,i}||P_i) \geq \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})} \right] \right] \\
&\leq \mathbb{E} \left[e^{-\frac{T'}{B+K-1} \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})}} \right] \\
&= e^{-\frac{T'}{B+K-1} \frac{B(R_B^{\text{go}} - \epsilon)}{H(\mathbf{P})}}. \tag{24}
\end{aligned}$$

Therefore, we have

$$\begin{aligned}
&\mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})] \\
&\leq \sum_{\mathbf{V}_1, \dots, \mathbf{V}_B \in \mathcal{V}^K} e^{-\frac{B}{B+K-1} \frac{(R_B^{\text{go}} - \epsilon) T'}{H(\mathbf{P})}} \\
&\quad (\text{by Eq. (24)}) \\
&\leq (T/B + 2)^{2KB} e^{-\frac{B}{B+K-1} \frac{(R_B^{\text{go}} - \epsilon) T'}{H(\mathbf{P})}}.
\end{aligned}$$

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

Proof of Theorem 5, Normal rewards. For the ease of discussion, we assume unit variance $\sigma = 1$. Extending it to the case of common known variance σ is straightforward. Let

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

Then, it is easy to see

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

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

The PoE is bounded as

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

We have,

$$\begin{aligned}
&\mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P}), \mathcal{B}^c] \\
&= \int_{-T}^T \cdots \int_{-T}^T \mathbf{1}[J(T) \notin \mathcal{I}^*(\mathbf{P})] p(\mathbf{Q}_B | \mathbf{Q}_{B-1} \dots \mathbf{Q}_1) d\mathbf{Q}_B \dots p(\mathbf{Q}_B | \mathbf{Q}_{B-1} \dots \mathbf{Q}_1) d\mathbf{Q}_b \dots p(\mathbf{Q}_1) d\mathbf{Q}_1. \tag{25}
\end{aligned}$$

Here,

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

Finally, we have

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

□

G.7 Proof of Theorem 7

Proof of Theorem 7. We first show that the limit

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

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

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

Theorem 5 implies that Algorithm 2 with $B = B_0$ and $\epsilon = \eta/2$ satisfies¹⁵

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

and thus

$$\inf H(\mathbf{P}) \liminf_{T \rightarrow \infty} \frac{\log(1/\mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})])}{T} \geq \frac{B_0}{B_0 + K - 1} \left(R_{B_0}^{\text{go}} - \frac{\eta}{2}\right). \quad (26)$$

Moreover, Theorem 2 implies that any algorithm satisfies

$$\inf H(\mathbf{P}) \limsup_{T \rightarrow \infty} \frac{\log(1/\mathbb{P}[J(T) \notin \mathcal{I}^*(\mathbf{P})])}{T} \leq R_{B_1}^{\text{go}}. \quad (27)$$

Combining Eq. (26) and Eq. (27), we have

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

and thus

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

¹⁵Strictly speaking, Algorithm 2 depends on T , and we take sequence of the algorithm $(\pi_{\text{DOT}, T})_{T=1,2,\dots}$.

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

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

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

and thus

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

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