

---

# Regret Lower Bound and Optimal Algorithm in Finite Stochastic Partial Monitoring

---

**Junpei Komiyama**  
The University of Tokyo  
junpei@komiyama.info

**Junya Honda**  
The University of Tokyo  
honda@stat.t.u-tokyo.ac.jp

**Hiroshi Nakagawa**  
The University of Tokyo  
nakagawa@dl.itc.u-tokyo.ac.jp

## Abstract

Partial monitoring is a general model for sequential learning with limited feedback formalized as a game between two players. In this game, the learner chooses an action and at the same time the opponent chooses an outcome, then the learner suffers a loss and receives a feedback signal. The goal of the learner is to minimize the total loss. In this paper, we study partial monitoring with finite actions and stochastic outcomes. We derive a logarithmic distribution-dependent regret lower bound that defines the hardness of the problem. Inspired by the DMED algorithm (Honda and Takemura, 2010) for the multi-armed bandit problem, we propose PM-DMED, an algorithm that minimizes the distribution-dependent regret. PM-DMED significantly outperforms state-of-the-art algorithms in numerical experiments. To show the optimality of PM-DMED with respect to the regret bound, we slightly modify the algorithm by introducing a hinge function (PM-DMED-Hinge). Then, we derive an asymptotically optimal regret upper bound of PM-DMED-Hinge that matches the lower bound.

## 1 Introduction

Partial monitoring is a general framework for sequential decision making problems with imperfect feedback. Many classes of problems, including prediction with expert advice [1], the multi-armed bandit problem [2], dynamic pricing [3], the dark pool problem [4], label efficient prediction [5], and linear and convex optimization with full or bandit feedback [6, 7] can be modeled as an instance of partial monitoring.

Partial monitoring is formalized as a repeated game played by two players called a learner and an opponent. At each round, the learner chooses an action, and at the same time the opponent chooses an outcome. Then, the learner observes a feedback signal from a given set of symbols and suffers some loss, both of which are deterministic functions of the selected action and outcome.

The goal of the learner is to find the optimal action that minimizes his/her cumulative loss. Alternatively, we can define the regret as the difference between the cumulative losses of the learner and the single optimal action, and minimization of the loss is equivalent to minimization of the regret. A learner with a small regret balances exploration (acquisition of information about the strategy of the opponent) and exploitation (utilization of information). The rate of regret indicates how fast the learner adapts to the problem: a linear regret indicates the inability of the learner to find the optimal action, whereas a sublinear regret indicates that the learner can approach the optimal action given sufficiently large time steps.

The study of partial monitoring is classified into two settings with respect to the assumption on the outcomes. On one hand, in the stochastic setting, the opponent chooses an outcome distribution before the game starts, and an outcome at each round is an i.i.d. sample from the distribution. On the other hand, in the adversarial setting, the opponent chooses the outcomes to maximize the regret of the learner. In this paper, we study the former setting.

## 1.1 Related work

The paper by Piccolboni and Schindelhauer [8] is one of the first to study the regret of the finite partial monitoring problem. They proposed the FeedExp3 algorithm, which attains  $O(T^{3/4})$  minimax regret on some problems. This bound was later improved by Cesa-Bianchi et al. [9] to  $O(T^{2/3})$ , who also showed an instance in which the bound is optimal. Since then, most literature on partial monitoring has dealt with the minimax regret, which is the worst-case regret over all possible opponent's strategies. Bartók et al. [10] classified the partial monitoring problems into four categories in terms of the minimax regret: a trivial problem with zero regret, an easy problem with  $\tilde{\Theta}(\sqrt{T})$  regret<sup>1</sup>, a hard problem with  $\Theta(T^{2/3})$  regret, and a hopeless problem with  $\Theta(T)$  regret. This shows that the class of the partial monitoring problems is not limited to the bandit sort but also includes larger classes of problems, such as dynamic pricing. Since then, several algorithms with a  $\tilde{O}(\sqrt{T})$  regret bound for easy problems have been proposed [11, 12, 13]. Among them, the Bayes-update Partial Monitoring (BPM) algorithm [13] is state-of-the-art in the sense of empirical performance.

**Distribution-dependent and minimax regret:** we focus on the distribution-dependent regret that depends on the strategy of the opponent. While the minimax regret in partial monitoring has been extensively studied, little has been known on distribution-dependent regret in partial monitoring. To the authors' knowledge, the only paper focusing on the distribution-dependent regret in finite discrete partial monitoring is the one by Bartók et al. [11], which derived  $O(\log T)$  distribution-dependent regret for easy problems. In contrast to this situation, much more interest in the distribution-dependent regret has been shown in the field of multi-armed bandit problems. Upper confidence bound (UCB), the most well-known algorithm for the multi-armed bandits, has a distribution-dependent regret bound [2, 14], and algorithms that minimize the distribution-dependent regret (e.g., KL-UCB) has been shown to perform better than ones that minimize the minimax regret (e.g., MOSS), even in instances in which the distributions are hard to distinguish (e.g., Scenario 2 in Garivier et al. [15]). Therefore, in the field of partial monitoring, we can expect that an algorithm that minimizes the distribution-dependent regret would perform better than the existing ones.

**Contribution:** the contributions of this paper lie in the following three aspects. First, we derive the regret lower bound: in some special classes of partial monitoring (e.g., multi-armed bandits), an  $O(\log T)$  regret lower bound is known to be achievable. In this paper, we further extend this lower bound to obtain a regret lower bound for general partial monitoring problems. Second, we propose an algorithm called Partial Monitoring DMED (PM-DMED). We also introduce a slightly modified version of this algorithm (PM-DMED-Hinge) and derive its regret bound. PM-DMED-Hinge is the first algorithm with a logarithmic regret bound for hard problems. Moreover, for both easy and hard problems, it is the first algorithm with the optimal constant factor on the leading logarithmic term. Third, performances of PM-DMED and existing algorithms are compared in numerical experiments. Here, the partial monitoring problems consisted of three specific instances of varying difficulty. In all instances, PM-DMED significantly outperformed the existing methods when a number of rounds is large. The regret of PM-DMED on these problems quickly approached the theoretical lower bound.

## 2 Problem Setup

This paper studies the finite stochastic partial monitoring problem with  $N$  actions,  $M$  outcomes, and  $A$  symbols. An instance of the partial monitoring game is defined by a loss matrix  $L = (l_{i,j}) \in \mathbb{R}^{N \times M}$  and a feedback matrix  $H = (h_{i,j}) \in [A]^{N \times M}$ , where  $[A] = \{1, 2, \dots, A\}$ . At the beginning, the learner is informed of  $L$  and  $H$ . At each round  $t = 1, 2, \dots, T$ , a learner selects an action  $i(t) \in [N]$ , and at the same time an opponent selects an outcome  $j(t) \in [M]$ . The learner

<sup>1</sup>Note that  $\tilde{\Theta}$  ignores a polylog factor.

suffers loss  $l_{i(t),j(t)}$ , which he/she cannot observe: the only information the learner receives is the signal  $h_{i(t),j(t)} \in [A]$ . We consider a stochastic opponent whose strategy for selecting outcomes is governed by the opponent's strategy  $p^* \in \mathcal{P}_M$ , where  $\mathcal{P}_M$  is a set of probability distributions over an  $M$ -ary outcome. The outcome  $j(t)$  of each round is an i.i.d. sample from  $p^*$ .

The goal of the learner is to minimize the cumulative loss over  $T$  rounds. Let the optimal action be the one that minimizes the loss in expectation, that is,  $i^* = \arg \min_{i \in [N]} L_i^\top p^*$ , where  $L_i$  is the  $i$ -th row of  $L$ . Assume that  $i^*$  is unique. Without loss of generality, we can assume that  $i^* = 1$ . Let  $\Delta_i = (L_i - L_1)^\top p^* \in [0, \infty)$  and  $N_i(t)$  be the number of rounds before the  $t$ -th in which action  $i$  is selected. The performance of the algorithm is measured by the (pseudo) regret,

$$\text{Regret}(T) = \sum_{t=1}^T \Delta_{i(t)} = \sum_{i \in [N]} \Delta_i N_i(T+1),$$

which is the difference between the expected loss of the learner and the optimal action. It is easy to see that minimizing the loss is equivalent to minimizing the regret. The expectation of the regret measures the performance of an algorithm that the learner uses.

For each action  $i \in [N]$ , let  $\mathcal{C}_i$  be the set of opponent strategies for which action  $i$  is optimal:

$$\mathcal{C}_i = \{q \in \mathcal{P}_M : \forall_{j \neq i} (L_i - L_j)^\top q \leq 0\}.$$

We call  $\mathcal{C}_i$  the optimality cell of action  $i$ . Each optimality cell is a convex closed polytope. Furthermore, we call the set of optimality cells  $\{\mathcal{C}_1, \dots, \mathcal{C}_N\}$  the cell decomposition as shown in Figure 1. Let  $\mathcal{C}_i^c = \mathcal{P}_M \setminus \mathcal{C}_i$  be the set of strategies with which action  $i$  is not optimal.

The signal matrix  $S_i \in \{0, 1\}^{A \times M}$  of action  $i$  is defined as  $(S_i)_{k,j} = \mathbb{1}[h_{i,j} = k]$ , where  $\mathbb{1}[X] = 1$  if  $X$  is true and 0 otherwise. The signal matrix defined here is slightly different from the one in the previous papers (e.g., Bartók et al. [10]) in which the number of rows of  $S_i$  is the number of the different symbols in the  $i$ -th row of  $H$ . The advantage in using the definition here is that,  $S_i p^* \in \mathbb{R}^A$  is a probability distribution over symbols that the algorithm observes when it selects an action  $i$ . Examples of signal matrices are shown in Section 5. An instance of partial monitoring is *globally observable* if for all pairs  $i, j$  of actions,  $L_i - L_j \in \oplus_{k \in [N]} \text{Im} S_k^\top$ . In this paper, we exclusively deal with globally observable instances: in view of the minimax regret, this includes trivial, easy, and hard problems.

### 3 Regret Lower Bound

A good algorithm should work well against any opponent's strategy. We extend this idea by introducing the notion of strong consistency: a partial monitoring algorithm is strongly consistent if it satisfies  $\mathbb{E}[\text{Regret}(T)] = o(T^a)$  for any  $a > 0$  and  $p \in \mathcal{P}_M$  given  $L$  and  $H$ .

In the context of the multi-armed bandit problem, Lai and Robbins [2] derived the regret lower bound of a strongly consistent algorithm: an algorithm must select each arm  $i$  until its number of draws  $N_i(t)$  satisfies  $\log t \lesssim N_i(t) d(\theta_i \| \theta_1)$ , where  $d(\theta_i \| \theta_1)$  is the KL divergence between the two one-parameter distributions from which the rewards of action  $i$  and the optimal action are generated. Analogously, in the partial monitoring problem, we can define the minimum number of observations.

**Lemma 1.** *For sufficiently large  $T$ , a strongly consistent algorithm satisfies:*

$$\forall_{q \in \mathcal{C}_1^c} \sum_{i \in [N]} \mathbb{E}[N_i(T)] D(p_i^* \| S_i q) \geq \log T - o(\log T),$$

where  $p_i^* = S_i p^*$  and  $D(p \| q) = \sum_i (p)_i \log((p)_i / (q)_i)$  is the KL divergence between two discrete distributions, in which we define  $0 \log 0 / 0 = 0$ .

Lemma 1 can be interpreted as follows: for each round  $t$ , consistency requires the algorithm to make sure that the possible risk that action  $i \neq 1$  is optimal is smaller than  $1/t$ . Large deviation principle [16] states that, the probability that an opponent with strategy  $q$  behaves like  $p^*$  is

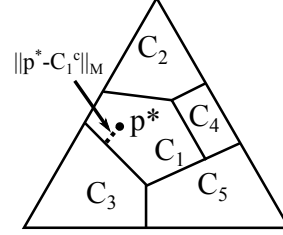


Figure 1: Cell decomposition of a partial monitoring instance with  $M = 3$ .

roughly  $\exp(-\sum_i N_i(t)D(p_i^* \| S_i q))$ . Therefore, we need to continue exploration of the actions until  $\sum_i N_i(t)D(p_i^* \| S_i q) \sim \log t$  holds for any  $q \in \mathcal{C}_1^c$  to reduce the risk to  $\exp(-\log t) = 1/t$ .

The proof of Lemma 1 is in Appendix B in the supplementary material. Based on the technique used in Lai and Robbins [2], the proof considers a modified game in which another action  $i \neq 1$  is optimal. The difficulty in proving the lower bound in partial monitoring lies in that, the feedback structure can be quite complex: for example, to confirm the superiority of action 1 over 2, one might need to use the feedback from action 3  $\notin \{1, 2\}$ . Still, we can derive the lower bound by utilizing the consistency of the algorithm in the original and modified games.

We next derive a lower bound on the regret based on Lemma 1. Note that, the expectation of the regret can be expressed as  $\mathbb{E}[\text{Regret}(T)] = \sum_{i \neq 1} \mathbb{E}[N_i(t)](L_i - L_1)^\top p^*$ . Let

$$\mathcal{R}_j(\{p_i\}) = \left\{ \{r_i\}_{i \neq j} \in [0, \infty)^{N-1} : \inf_{q \in \text{cl}(\mathcal{C}_j^c): p_j = S_j q} \sum_i r_i D(p_i \| S_i q) \geq 1 \right\},$$

where  $\text{cl}(\cdot)$  denotes a closure. Moreover, let

$$C_j^*(p, \{p_i\}) = \inf_{r_i \in \mathcal{R}_j(\{p_i\})} \sum_{i \neq j} r_i (L_i - L_j)^\top p,$$

the optimal solution of which is

$$\mathcal{R}_j^*(p, \{p_i\}) = \left\{ \{r_i\}_{i \neq j} \in \mathcal{R}_j(\{p_i\}) : \sum_{i \neq j} r_i (L_i - L_j)^\top p = C_j^*(p, \{p_i\}) \right\}.$$

The value  $C_1^*(p^*, \{p_i^*\}) \log T$  is the possible minimum regret for observations such that the minimum divergence of  $p^*$  from any  $q \in \mathcal{C}_1^c$  is larger than  $\log T$ . Using Lemma 1 yields the following regret lower bound:

**Theorem 2.** *The regret of a strongly consistent algorithm is lower bounded as:*

$$\mathbb{E}[\text{Regret}(T)] \geq C_1^*(p^*, \{p_i^*\}) \log T - o(\log T).$$

From this theorem, we can naturally measure the harshness of the instance by  $C_1^*(p^*, \{p_i^*\})$ , whereas the past studies (e.g., Vanchinathan et al. [13]) ambiguously define the harshness as the closeness to the boundary of the cells. Furthermore, we show in Lemma 5 in the Appendix that  $C_1^*(p^*, \{p_i^*\}) = O(N/\|p^* - \mathcal{C}_1^c\|_M^2)$ : the regret bound has at most quadratic dependence on  $\|p^* - \mathcal{C}_1^c\|_M$ , which is defined in Appendix D as the closeness of  $p^*$  to the boundary of the optimal cell.

## 4 PM-DMED Algorithm

In this section, we describe the partial monitoring deterministic minimum empirical divergence (PM-DMED) algorithm, which is inspired by DMED [17] for solving the multi-armed bandit problem. Let  $\hat{p}_i(t) \in [0, 1]^A$  be the empirical distribution of the symbols under the selection of action  $i$ . Namely, the  $k$ -th element of  $\hat{p}_i(t)$  is  $(\sum_{t'=1}^{t-1} \mathbb{1}[i(t') = i \cap h_{i(t'), j(t')} = k]) / (\sum_{t'=1}^{t-1} \mathbb{1}[i(t') = i])$ . We sometimes omit  $t$  from  $\hat{p}_i$  when it is clear from the context. Let the empirical divergence of  $q \in \mathcal{P}_M$  be  $\sum_{i \in [N]} N_i(t)D(\hat{p}_i(t) \| S_i q)$ , the exponential of which can be considered as a likelihood that  $q$  is the opponent's strategy.

The main routine of PM-DMED is in Algorithm 1. At each loop, the actions in the current list  $Z_C$  are selected once. The list for the actions in the next loop  $Z_N$  is determined by the subroutine in Algorithm 2. The subroutine checks whether the empirical divergence of each point  $q \in \mathcal{C}_1^c$  is larger than  $\log t$  or not (Eq. (3)). If it is large enough, it exploits the current information by selecting  $\hat{i}(t)$ , the optimal action based on the estimation  $\hat{p}(t)$  that minimizes the empirical divergence. Otherwise, it selects the actions with the number of observations below the minimum requirement for making the empirical divergence of each suboptimal point  $q \in \mathcal{C}_1^c$  larger than  $\log t$ .

Unlike the  $N$ -armed bandit problem in which a reward is associated with an action, in the partial monitoring problem, actions, outcomes, and feedback signals can be intricately related. Therefore, we need to solve a non-trivial optimization to run PM-DMED. Later in Section 5, we discuss a practical implementation of the optimization.

---

**Algorithm 1** Main routine of PM-DMED and PM-DMED-Hinge

---

```

1: Initialization: select each action once.
2:  $Z_C, Z_R \leftarrow [N], Z_N \leftarrow \emptyset$ .
3: while  $t \leq T$  do
4:   for  $i(t) \in Z_C$  in an arbitrarily fixed order do
5:     Select  $i(t)$ , and receive feedback.
6:      $Z_R \leftarrow Z_R \setminus \{i(t)\}$ .
7:     Add actions to  $Z_N$  in accordance with
        $\begin{cases} \text{Algorithm 2 (PM-DMED)} \\ \text{Algorithm 3 (PM-DMED-Hinge)} \end{cases}$ .
8:      $t \leftarrow t + 1$ .
9:   end for
10:   $Z_C, Z_R \leftarrow Z_N, Z_N \leftarrow \emptyset$ .
11: end while

```

---



---

**Algorithm 2** PM-DMED subroutine for adding actions to  $Z_N$  (without duplication).

---

```

1: Parameter:  $c > 0$ .
2: Compute an arbitrary  $\hat{p}(t)$  such that
   
$$\hat{p}(t) \in \arg \min_q \sum_i N_i(t) D(\hat{p}_i(t) \| S_i q) \quad (1)$$

   and let  $\hat{i}(t) = \arg \min_i L_i^\top \hat{p}(t)$ .
3: If  $\hat{i}(t) \notin Z_R$  then put  $\hat{i}(t)$  into  $Z_N$ .
4: If there are actions  $i \notin Z_R$  such that
   
$$N_i(t) < c\sqrt{\log t} \quad (2)$$

   then put them into  $Z_N$ .
5: If
   
$$\{N_i(t)/\log t\}_{i \neq \hat{i}(t)} \notin \mathcal{R}_{\hat{i}(t)}(\{\hat{p}_i(t)\}) \quad (3)$$

   then compute some
   
$$\{r_i^*\}_{i \neq \hat{i}(t)} \in \mathcal{R}_{\hat{i}(t)}^*(\hat{p}(t), \{\hat{p}_i(t)\}) \quad (4)$$

   and put all actions  $i$  such that  $i \notin Z_R$  and
    $r_i^* > N_i(t)/\log t$  into  $Z_N$ .

```

---

**Necessity of  $\sqrt{\log T}$  exploration:** PM-DMED tries to observe each action to some extent (Eq. (2)), which is necessary for the following reason: consider a four-state game characterized by

$$L = \begin{pmatrix} 0 & 1 & 1 & 0 \\ 10 & 1 & 0 & 0 \\ 10 & 0 & 1 & 0 \\ 11 & 11 & 11 & 11 \end{pmatrix}, \quad H = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 2 & 3 \\ 1 & 2 & 2 & 3 \\ 1 & 1 & 2 & 2 \end{pmatrix}, \text{ and } p^* = (0.1, 0.2, 0.3, 0.4)^\top.$$

The optimal action here is action 1, which does not yield any useful information. By using action 2, one receives three kinds of symbols from which one can estimate  $(p^*)_1$ ,  $(p^*)_2 + (p^*)_3$ , and  $(p^*)_4$ , where  $(p^*)_j$  is the  $j$ -th component of  $p^*$ . From this, an algorithm can find that  $(p^*)_1$  is not very small and thus the expected loss of actions 2 and 3 is larger than that of action 1. Since the feedback of actions 2 and 3 are the same, one may also use action 3 in the same manner. However, the loss per observation is 1.2 and 1.3 for actions 2 and 3, respectively, and thus it is better to use action 2. This difference comes from the fact that  $(p^*)_2 = 0.2 < 0.3 = (p^*)_3$ . Since an algorithm does not know  $p^*$  beforehand, it needs to observe action 4, the only source for distinguishing  $(p^*)_2$  from  $(p^*)_3$ . Yet, an optimal algorithm cannot select it more than  $\Omega(\log T)$  times because it affects the  $O(\log T)$  factor in the regret. In fact,  $O((\log T)^a)$  observations of action 4 with some  $a > 0$  are sufficient to be convinced that  $(p^*)_2 < (p^*)_3$  with probability  $1 - o(1/T^{\text{poly}(a)})$ . For this reason, PM-DMED selects each action  $\sqrt{\log t}$  times.

## 5 Experiment

Following Bartók et al. [11], we compared the performances of algorithms in three different games: the four-state game (Section 4), a three-state game and dynamic pricing. Experiments on the  $N$ -armed bandit game was also done, and the result is shown in Appendix C.1.

The three-state game, which is classified as easy in terms of the minimax regret, is characterized by:

$$L = \begin{pmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \end{pmatrix} \quad \text{and} \quad H = \begin{pmatrix} 1 & 2 & 2 \\ 2 & 1 & 2 \\ 2 & 2 & 1 \end{pmatrix}.$$

The signal matrices of this game are,

$$S_1 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \end{pmatrix}, \quad S_2 = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix}, \quad \text{and } S_3 = \begin{pmatrix} 0 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix}.$$

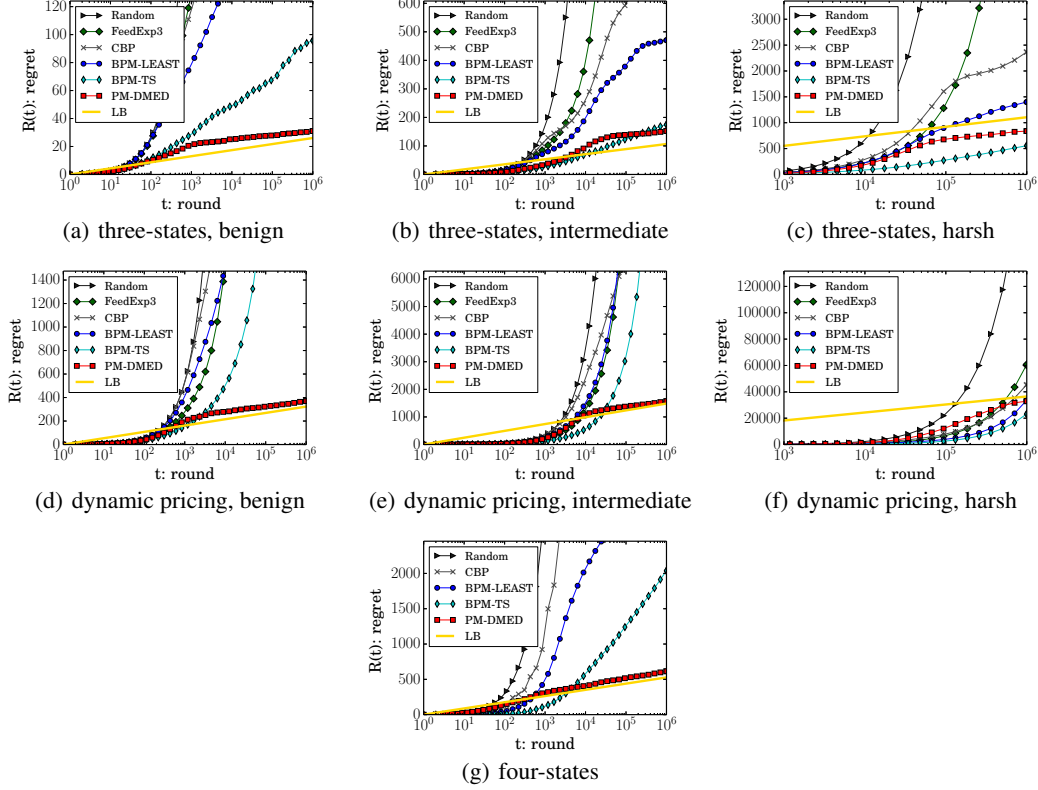


Figure 2: Regret-round semilog plots of algorithms. The regrets are averaged over 100 runs. LB is the asymptotic regret lower bound of Theorem 2.

Dynamic pricing, which is classified as hard in terms of the minimax regret, is a game that models a repeated auction between a seller (learner) and a buyer (opponent). At each round, the seller sets a price for a product, and at the same time, the buyer secretly sets a maximum price he is willing to pay. The signal is “buy” or “no-buy”, and the seller’s loss is either a given constant (no-buy) or the difference between the buyer’s and the seller’s prices (buy). The loss and feedback matrices are:

$$L = \begin{pmatrix} 0 & 1 & \dots & N-1 \\ c & 0 & \dots & N-2 \\ \vdots & \ddots & \ddots & \vdots \\ c & \dots & c & 0 \end{pmatrix} \quad \text{and} \quad H = \begin{pmatrix} 2 & 2 & \dots & 2 \\ 1 & 2 & \dots & 2 \\ \vdots & \ddots & \ddots & \vdots \\ 1 & \dots & 1 & 2 \end{pmatrix},$$

where signals 1 and 2 correspond to no-buy and buy. The signal matrix of action  $i$  is

$$S_i = \begin{pmatrix} \overbrace{1 \dots 1}^{i-1} & \overbrace{0 \dots 0}^{M-i+1} \\ 0 & \dots & 1 \end{pmatrix}.$$

Following Bartók et al. [11], we set  $N = 5$ ,  $M = 5$ , and  $c = 2$ .

In our experiments with the three-state game and dynamic pricing, we tested three settings regarding the harshness of the opponent: at the beginning of a simulation, we sampled 1,000 points uniformly at random from  $\mathcal{P}_M$ , then sorted them by  $C_1^*(p^*, \{p_i^*\})$ . We chose the top 10%, 50%, and 90% harshest ones as the opponent’s strategy in the harsh, intermediate, and benign settings, respectively.

We compared Random, FeedExp3 [8], CBP [11] with  $\alpha = 1.01$ , BPM-LEAST, BPM-TS [13], and PM-DMED with  $c = 1$ . Random is a naive algorithm that selects an action uniformly random. FeedExp3 requires a matrix  $G$  such that  $H^\top G = L^\top$ , and thus one cannot apply it to the four-state game. CBP is an algorithm of logarithmic regret for easy games. The parameters  $\eta$  and  $f(t)$  of CBP were set in accordance with Theorem 1 in their paper. BPM-LEAST is a Bayesian algorithm with  $\tilde{O}(\sqrt{T})$  regret for easy games, and BPM-TS is a heuristic of state-of-the-art performance. The priors of two BPMs were set to be uninformative to avoid a misspecification, as recommended in their paper.

---

**Algorithm 3** PM-DMED-Hinge subroutine for adding actions to  $Z_N$  (without duplication).

---

- 1: **Parameters:**  $c > 0$ ,  $f(n) = bn^{-1/2}$  for  $b > 0$ ,  $\alpha(t) = a/(\log \log t)$  for  $a > 0$ .
- 2: Compute arbitrary  $\hat{p}(t)$  which satisfies

$$\hat{p}(t) \in \arg \min_q \sum_i N_i(t)(D(\hat{p}_i(t) \| S_i q) - f(N_i(t)))_+ \quad (5)$$

and let  $\hat{i}(t) = \arg \min_i L_i^\top \hat{p}(t)$ .

- 3: If  $\hat{i}(t) \notin Z_R$  then put  $\hat{i}(t)$  into  $Z_N$ .

- 4: If

$$\hat{p}(t) \notin \mathcal{C}_{\hat{i}(t), \alpha(t)} \quad (6)$$

or there exists an action  $i$  such that

$$D(\hat{p}_i(t) \| S_i \hat{p}(t)) > f(N_i(t)) \quad (7)$$

then put all actions  $i \notin Z_R$  into  $Z_N$ .

- 5: If there are actions  $i$  such that

$$N_i(t) < c\sqrt{\log t} \quad (8)$$

then put the actions not in  $Z_R$  into  $Z_N$ .

- 6: If

$$\{N_i(t)/\log t\}_{i \neq \hat{i}(t)} \notin \mathcal{R}_{\hat{i}(t)}(\{\hat{p}_i(t), f(N_i(t))\}) \quad (9)$$

then compute some

$$\{r_i^*\}_{i \neq \hat{i}(t)} \in \mathcal{R}_{\hat{i}(t)}^*(\hat{p}(t), \{\hat{p}_i(t), f(N_i(t))\}) \quad (10)$$

and put all actions such that  $i \notin Z_R$  and  $r_i^* > N_i(t)/\log t$  into  $Z_N$ . If such  $r_i^*$  is infeasible then put all action  $i \notin Z_R$  into  $Z_N$ .

---

The computation of  $\hat{p}(t)$  in (1) and the evaluation of the condition in (3) involve convex optimizations, which were done with Ipopt [18]. Moreover, obtaining  $\{r_i^*\}$  in (4) is classified as a linear semi-infinite programming (LSIP) problem, a linear programming (LP) with finitely many variables and infinitely many constraints. Following the optimization of BPM-LEAST [13], we resorted to a finite sample approximation and used the Gurobi LP solver [19] in computing  $\{r_i^*\}$ : at each round, we sampled 1,000 points from  $\mathcal{P}_M$ , and relaxed the constraints on the samples. To speed up the computation, we skipped these optimizations in most rounds with large  $t$  and used the result of the last computation. The computation of the coefficient  $C_1^*(p^*, \{p_i^*\})$  of the regret lower bound (Theorem 2) is also an LSIP, which was approximated by 100,000 sample points from  $\mathcal{C}_1^c$ .

The experimental results are shown in Figure 2. In the four-state game and the other two games with an easy or intermediate opponent, PM-DMED outperforms the other algorithms when the number of rounds is large. In particular, in the dynamic pricing game with an intermediate opponent, the regret of PM-DMED at  $T = 10^6$  is ten times smaller than those of the other algorithms. Even in the harsh setting in which the minimax regret matters, PM-DMED has some advantage over all algorithms except for BPM-TS. With sufficiently large  $T$ , the slope of an optimal algorithm should converge to LB. In all games and settings, the slope of PM-DMED converges to LB, which is empirical evidence of the optimality of PM-DMED.

## 6 Theoretical Analysis

Section 5 shows that the empirical performance of PM-DMED is very close to the regret lower bound in Theorem 2. Although the authors conjecture that PM-DMED is optimal, it is hard to analyze PM-DMED. The technically hardest part arises from the case in which the divergence of each action is small but not yet fully converged. To circumvent this difficulty, we can introduce a discount factor. Let

$$\mathcal{R}_j(\{p_i, \delta_i\}) = \left\{ \{r_i\}_{i \neq j} \in [0, \infty)^{N-1} : \inf_{q \in \text{cl}(\mathcal{C}_j^c): D(p_j \| S_j q) \leq \delta_j} \sum_i r_i (D(p_i \| S_i q) - \delta_i)_+ \geq 1 \right\}, \quad (11)$$

where  $(X)_+ = \max(X, 0)$ . Note that  $\mathcal{R}_j(\{p_i, \delta_i\})$  in (11) is a natural generalization of  $\mathcal{R}_j(\{p_i\})$  in Section 4 in the sense that  $\mathcal{R}_j(\{p_i, 0\}) = \mathcal{R}_j(\{p_i\})$ . Event  $\{N_i(t)/\log t\}_{i \neq 1} \in \mathcal{R}_1(\{\hat{p}_i(t), \delta_i\})$  means that the number of observations  $\{N_i(t)\}$  is enough to ensure that the “ $\{\delta_i\}$ -discounted” empirical divergence of each  $q \in \mathcal{C}_1^c$  is larger than  $\log t$ . Analogous to  $\mathcal{R}_j(\{p_i, \delta_i\})$ , we define

$$C_j^*(p, \{p_i, \delta_i\}) = \inf_{\{r_i\}_{i \neq j} \in \mathcal{R}_j(\{p_i, \delta_i\})} \sum_{i \neq j} r_i (L_j - L_i)^\top p$$

and its optimal solution by

$$\mathcal{R}_j^*(p, \{p_i, \delta_i\}) = \left\{ \{r_i\}_{i \neq j} \in \mathcal{R}_j(\{p_i, \delta_i\}) : \sum_{i \neq j} r_i (L_j - L_i)^\top p = C_j^*(p, \{p_i, \delta_i\}) \right\}.$$

We also define  $\mathcal{C}_{i,\alpha} = \{p \in \mathcal{P}_M : L_i^\top p + \alpha \leq \min_{j \neq i} L_j^\top p\}$ , the optimal region of action  $i$  with margin. PM-DMED-Hinge shares the main routine of Algorithm 1 with PM-DMED and lists the next actions by Algorithm 3. Unlike PM-DMED, it (i) discounts  $f(N_i(t))$  from the empirical divergence  $D(\hat{p}_i(t) \| S_i q)$ . Moreover, (ii) when  $\hat{p}(t)$  is close to the cell boundary, it encourages more exploration to identify the cell it belongs to by Eq. (6).

**Theorem 3.** *Assume that the following regularity conditions hold for  $p^*$ . (1)  $\mathcal{R}_1^*(p, \{p_i, \delta_i\})$  is unique at  $p = p^*, p_i = S_i p^*, \delta_i = 0$ . Moreover, (2) for  $\mathcal{S}_\delta = \{q : D(p_1^* \| S_1 q) \leq \delta\}$ , it holds that  $\text{cl}(\text{int}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta) = \text{cl}(\text{cl}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta)$  for all  $\delta \geq 0$  in some neighborhood of  $\delta = 0$ , where  $\text{cl}(\cdot)$  and  $\text{int}(\cdot)$  denote the closure and the interior, respectively. Then,*

$$\mathbb{E}[\text{Regret}(T)] \leq C_1^*(p^*, \{p_i^*\}) \log T + o(\log T).$$

We prove this theorem in Appendix D. Recall that  $\mathcal{R}_1^*(p, \{\hat{p}_i(t), \delta_i\})$  is the set of optimal solutions of an LSIP. In this problem, KKT conditions and the duality theorem apply as in the case of finite constraints; thus, we can check whether Condition 1 holds or not for each  $p^*$  (see, e.g., Ito et al. [20] and references therein). Condition 2 holds in most cases, and an example of an exceptional case is shown in Appendix A.

Theorem 3 states that PM-DMED-Hinge has a regret upper bound that matches the lower bound of Theorem 2.

**Corollary 4.** *(Optimality in the  $N$ -armed bandit problem) In the  $N$ -armed Bernoulli bandit problem, the regularity conditions in Theorem 3 always hold. Moreover, the coefficient of the leading logarithmic term in the regret bound of the partial monitoring problem is equal to the bound given in Lai and Robbins [2]. Namely,  $C_1^*(p^*, \{p_i^*\}) = \sum_{i \neq 1}^N (\Delta_i / d(\mu_i \| \mu_1))$ , where  $d(p \| q) = p \log(p/q) + (1-p) \log((1-p)/(1-q))$  is the KL-divergence between Bernoulli distributions.*

Corollary 4, which is proven in Appendix C, states that PM-DMED-Hinge attains the optimal regret of the  $N$ -armed bandit if we run it on an  $N$ -armed bandit game represented as partial monitoring.

**Asymptotic analysis:** it is Theorem 6 where we lose the finite-time property. This theorem shows the continuity of the optimal solution set  $\mathcal{R}_1^*(p, \{p_i, \delta_i\})$  of  $C_j^*(p, \{p_j\})$ , which does not mention how close  $\mathcal{R}_1^*(p, \{p_i, \delta_i\})$  is to  $\mathcal{R}_1^*(p^*, \{p_i^*, 0\})$  if  $\max\{\|p - p^*\|_M, \max_i \|p_i - p_i^*\|_M, \max_i \delta_i\} \leq \delta$  for given  $\delta$ . To obtain an explicit bound, we need *sensitivity analysis*, the theory of the robustness of the optimal value and the solution for small deviations of its parameters (see e.g., Fiacco [21]). In particular, the optimal solution of partial monitoring involves an infinite number of constraints, which makes the analysis quite hard. For this reason, we will not perform a finite-time analysis. Note that, the  $N$ -armed bandit problem is a special instance in which we can avoid solving the above optimization and a finite-time optimal bound is known.

**Necessity of the discount factor:** we are not sure whether discount factor  $f(n)$  in PM-DMED-Hinge is necessary or not. We also empirically tested PM-DMED-Hinge: although it is better than the other algorithms in many settings, such as dynamic pricing with an intermediate opponent, it is far worse than PM-DMED. We found that our implementation, which uses the Ipopt nonlinear optimization solver, was sometimes inaccurate at optimizing (5): there were some cases in which the true  $p^*$  satisfies  $\forall_{i \in [N]} D(\hat{p}_i(t) \| S_i p^*) - f(N_i(t)) = 0$ , while the solution  $\hat{p}(t)$  we obtained had non-zero hinge values. In this case, the algorithm lists all actions from (7), which degrades performance. Determining whether the discount factor is essential or not is our future work.

## Acknowledgements

The authors gratefully acknowledge the advice of Kentaro Minami and sincerely thank the anonymous reviewers for their useful comments. This work was supported in part by JSPS KAKENHI Grant Number 15J09850 and 26106506.



## References

- [1] Nick Littlestone and Manfred K. Warmuth. The weighted majority algorithm. *Inf. Comput.*, 108(2):212–261, February 1994.
- [2] T. L. Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1):4–22, 1985.
- [3] Robert D. Kleinberg and Frank Thomson Leighton. The value of knowing a demand curve: Bounds on regret for online posted-price auctions. In *FOCS*, pages 594–605, 2003.
- [4] Alekh Agarwal, Peter L. Bartlett, and Max Dama. Optimal allocation strategies for the dark pool problem. In *AISTATS*, pages 9–16, 2010.
- [5] Nicolò Cesa-Bianchi, Gábor Lugosi, and Gilles Stoltz. Minimizing regret with label efficient prediction. *IEEE Transactions on Information Theory*, 51(6):2152–2162, 2005.
- [6] Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *ICML*, pages 928–936, 2003.
- [7] Varsha Dani, Thomas P. Hayes, and Sham M. Kakade. Stochastic linear optimization under bandit feedback. In *COLT*, pages 355–366, 2008.
- [8] Antonio Piccolboni and Christian Schindelhauer. Discrete prediction games with arbitrary feedback and loss. In *COLT*, pages 208–223, 2001.
- [9] Nicolò Cesa-Bianchi, Gábor Lugosi, and Gilles Stoltz. Regret minimization under partial monitoring. *Math. Oper. Res.*, 31(3):562–580, 2006.
- [10] Gábor Bartók, Dávid Pál, and Csaba Szepesvári. Minimax regret of finite partial-monitoring games in stochastic environments. In *COLT*, pages 133–154, 2011.
- [11] Gábor Bartók, Navid Zolghadr, and Csaba Szepesvári. An adaptive algorithm for finite stochastic partial monitoring. In *ICML*, 2012.
- [12] Gábor Bartók. A near-optimal algorithm for finite partial-monitoring games against adversarial opponents. In *COLT*, pages 696–710, 2013.
- [13] Hastagiri P. Vanchinathan, Gábor Bartók, and Andreas Krause. Efficient partial monitoring with prior information. In *NIPS*, pages 1691–1699, 2014.
- [14] Peter Auer, Nicolò Cesa-bianchi, and Paul Fischer. Finite-time Analysis of the Multiarmed Bandit Problem. *Machine Learning*, 47:235–256, 2002.
- [15] Aurélien Garivier and Olivier Cappé. The KL-UCB algorithm for bounded stochastic bandits and beyond. In *COLT*, pages 359–376, 2011.
- [16] Amir Dembo and Ofer Zeitouni. *Large deviations techniques and applications*. Applications of mathematics. Springer, New York, Berlin, Heidelberg, 1998.
- [17] Junya Honda and Akimichi Takemura. An Asymptotically Optimal Bandit Algorithm for Bounded Support Models. In *COLT*, pages 67–79, 2010.
- [18] Andreas Wächter and Carl D. Laird. Interior point optimizer (IPOPT).
- [19] Gurobi Optimization Inc. Gurobi optimizer.
- [20] S. Ito, Y. Liu, and K. L. Teo. A dual parametrization method for convex semi-infinite programming. *Annals of Operations Research*, 98(1-4):189–213, 2000.
- [21] Anthony V. Fiacco. *Introduction to sensitivity and stability analysis in nonlinear programming*. Academic Press, New York, 1983.
- [22] T.L. Graves and T.L. Lai. Asymptotically efficient adaptive choice of control laws in controlled Markov chains. *SIAM J. Contr. and Opt.*, 35(3):715–743, 1997.
- [23] Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of A/B testing. In *COLT*, pages 461–481, 2014.
- [24] Sébastien Bubeck. *Bandits Games and Clustering Foundations*. Theses, Université des Sciences et Technologie de Lille - Lille I, June 2010.
- [25] William W. Hogan. Point-to-set maps in mathematical programming. *SIAM Review*, 15(3):591–603, 1973.
- [26] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory, Second Edition (Wiley Series in Telecommunications and Signal Processing)*. Wiley-Interscience, 2006.

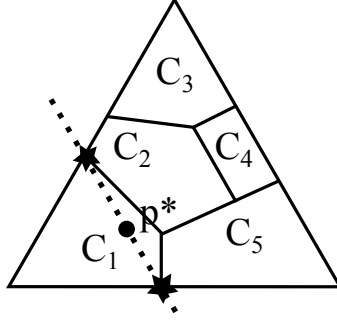


Figure 3: A corner case.

## A Case in which Condition 2 Does Not Hold

Figure 3 is an example that Theorem 3 does not cover. The dotted line is  $\{q : p_1^* = S_1 q\}$ , which (accidentally) coincides with a line that makes the convex polytope of  $\mathcal{C}_1^c$ . In this case, Condition 2 does not hold because  $\text{int}(\mathcal{C}_1^c) \cap S_0 = \emptyset$  whereas  $\text{cl}(\mathcal{C}_1^c) \cap S_0 \neq \emptyset$  (two starred points), which means that a slight modification of  $p^*$  changes the set of cells that intersects with the dotted line discontinuously. We exclude these unusual cases for the ease of analysis.

The authors consider that it is quite hard to give the optimal regret bound without such regularity conditions. In fact, many regularity conditions are assumed in Graves and Lai [22], where another generalization of the bandit problem is considered and the regret lower bound is expressed in terms of LSIP. In this paper, the regularity conditions are much simplified by the continuity argument in Theorem 6 but it remains an open problem to fully remove them.

## B Proof: Regret Lower Bound

In this section, we prove Lemma 1 and Theorem 2.

*Proof of Lemma 1.* The technique here is mostly inspired from Theorem 1 in Lai and Robbins [2]. The use of a  $\sqrt{T}$  term is inspired from Kaufmann et al. [23]. Let  $p' \in \text{int}(\mathcal{C}_1^c)$  and  $i' \neq 1$  be the optimal action under the opponent's strategy  $p'$ . We consider a modified partial monitoring game with its opponent's strategy is  $p'$ .

**Notation:** Let  $\hat{X}_i^m \in [A]$  is the signal of the  $m$ -th observation of action  $i$ . Let

$$\widehat{\text{KL}}_i(n) = \sum_{m=1}^n \log \left( \frac{(S_i p^*)_{\hat{X}_i^m}}{(S_i p')_{\hat{X}_i^m}} \right),$$

and  $\widehat{\text{KL}} = \sum_{i \in [N]} \widehat{\text{KL}}_i(N_i(T))$ . Let  $\mathbb{P}'$  and  $\mathbb{E}'$  be the probability and the expectation with respect to the modified game, respectively. Then, for any event  $\mathcal{E}$ ,

$$\mathbb{P}'[\mathcal{E}] = \mathbb{E} \left[ \mathbb{1}[\mathcal{E}] \exp \left( -\widehat{\text{KL}} \right) \right] \quad (12)$$

holds. Now, let us define the following events:

$$\begin{aligned} \mathcal{D}_1 &= \left\{ \sum_{i \in [N]} N_i(T) D(S_i p^* \| S_i p') < (1 - \epsilon) \log T, N_{i'}(T) < \sqrt{T} \right\}, \\ \mathcal{D}_2 &= \left\{ \widehat{\text{KL}} \leq \left( 1 - \frac{\epsilon}{2} \right) \log T \right\}, \\ \mathcal{D}_{12} &= \mathcal{D}_1 \cap \mathcal{D}_2, \\ \mathcal{D}_{1 \setminus 2} &= \mathcal{D}_1 \cap \mathcal{D}_2^c. \end{aligned}$$

**First step** ( $\Pr[\mathcal{D}_{12}] = o(1)$ ): from (12),

$$\mathbb{P}'[\mathcal{D}_{12}] \geq \mathbb{E} \left[ \mathbb{1}[\mathcal{D}_{12}] \exp \left( - \left( 1 - \frac{\epsilon}{2} \right) \log T \right) \right] = T^{-(1-\epsilon/2)} \Pr[\mathcal{D}_{12}].$$

By using this we have

$$\begin{aligned} \Pr[\mathcal{D}_{12}] &\leq T^{(1-\epsilon/2)} \mathbb{P}'[\mathcal{D}_{12}] \\ &\leq T^{(1-\epsilon/2)} \mathbb{P}' \left[ N_{i'}(T) < \sqrt{T} \right] \\ &= T^{(1-\epsilon/2)} \mathbb{P}' \left[ T - N_{i'}(T) > T - \sqrt{T} \right] \\ &\leq T^{(1-\epsilon/2)} \frac{\mathbb{E}'[T - N_{i'}(T)]}{T - \sqrt{T}} \quad (\text{by the Markov inequality}). \end{aligned} \quad (13)$$

Since this algorithm is strongly consistent,  $\mathbb{E}'[T - N_{i'}(T)] \rightarrow o(T^a)$  for any  $a > 0$ . Therefore, the RHS of the last line of (13) is  $o(T^{a-\epsilon/2})$ , which, by choosing sufficiently small  $a$ , converges to zero as  $T \rightarrow \infty$ . In summary,  $\Pr[\mathcal{D}_{12}] = o(1)$ .

**Second step** ( $\Pr[\mathcal{D}_{1 \setminus 2}] = o(1)$ ): we have

$$\begin{aligned} &\Pr[\mathcal{D}_{1 \setminus 2}] \\ &= \Pr \left[ \sum_{i \in [N]} N_i(T) D(S_i p^* \| S_i p') < (1 - \epsilon) \log T, N_{i'}(T) < \sqrt{T}, \sum_{i \in [N]} \widehat{\text{KL}}_i(N_i(T)) > \left( 1 - \frac{\epsilon}{2} \right) \log T \right]. \end{aligned}$$

Note that

$$\max_{1 \leq n \leq N} \widehat{\text{KL}}_i(n) = \max_{1 \leq n \leq N} \sum_{m=1}^n \log \left( \frac{(S_i p^*)_{\hat{X}_i^m}}{(S_i p')_{\hat{X}_i^m}} \right),$$

is the maximum of the sum of positive-mean random variables, and thus converges to its average (c.f., Lemma 10.5 in [24]). Namely,

$$\lim_{N \rightarrow \infty} \max_{1 \leq n \leq N} \frac{\widehat{\text{KL}}_i(n)}{N} \rightarrow D(S_i p^* \| S_i p')$$

almost surely. Therefore,

$$\lim_{T \rightarrow \infty} \frac{\max_{\{N_i(T)\} \in \mathbb{N}^N, \sum_{i \in [N]} N_i(T) D(S_i p^* \| S_i p') < (1-\epsilon) \log T} \sum_{i \in [N]} \widehat{\text{KL}}_i(N_i(T))}{\log T} \rightarrow 1 - \epsilon$$

almost surely. By using this fact and  $1 - \epsilon/2 > 1 - \epsilon$ , we have

$$\Pr \left[ \max_{\{N_i(T)\} \in \mathbb{N}^N, \sum_{i \in [N]} N_i(T) D(S_i p^* \| S_i p') < (1-\epsilon) \log T} \sum_{i \in [N]} \widehat{\text{KL}}_i(N_i(T)) > \left( 1 - \frac{\epsilon}{2} \right) \log T \right] = o(1).$$

In summary, we obtain  $\Pr[\mathcal{D}_{1 \setminus 2}] = o(1)$ .

**Last step:** we here have

$$\begin{aligned} \mathcal{D}_1 &= \left\{ \sum_{i \in [N]} N_i(T) D(S_i p^* \| S_i p') < (1 - \epsilon) \log T \right\} \cap \left\{ N_{i'}(T) < \sqrt{T} \right\} \\ &\supseteq \left\{ \sum_{i \in [N]} N_i(T) D(S_i p^* \| S_i p') + \frac{(1 - \epsilon) \log T}{\sqrt{T}} N_{i'}(T) < (1 - \epsilon) \log T \right\}, \end{aligned}$$

where we used the fact that  $\{A < C\} \cap \{B < C\} \supseteq \{A + B < C\}$  for  $A, B > 0$  in the last line. Note that, by using the result of the previous steps,  $\Pr[\mathcal{D}_1] = \Pr[\mathcal{D}_{12}] + \Pr[\mathcal{D}_{1 \setminus 2}] = o(1)$ . By using the complementarity of this fact,

$$\Pr \left[ \sum_{i \in [N]} N_i(T) D(S_i p^* \| S_i p') + \frac{(1 - \epsilon) \log T}{\sqrt{T}} N_{i'}(T) \geq (1 - \epsilon) \log T \right] \geq \Pr[\mathcal{D}_1^c] = 1 - o(1).$$

Using the Markov inequality yields

$$\mathbb{E} \left[ \sum_{i \in [N]} N_i(T) D(S_i p^* \| S_i p') + \frac{(1 - \epsilon) \log T}{\sqrt{T}} N_{i'}(T) \right] \geq (1 - \epsilon)(1 - o(1)) \log T. \quad (14)$$

Because  $\mathbb{E}[N_{i'}(T)]$  is subpolynomial as a function of  $T$  due to the consistency, the second term in LHS of (14) is  $o(1)$  and thus negligible. Lemma 1 follows from the fact that (14) holds for sufficiently small  $\epsilon$  and arbitrary  $p' \in \text{int}(\mathcal{C}_1^c)$ .  $\square$

*Proof of Theorem 2.* Assume that there exists  $\delta > 0$  and a sequence  $T_1 < T_2 < T_3 < \dots$  such that for all  $t$

$$\mathbb{E}[\text{Regret}(T_t)] < (1 - \delta) C_1^*(p^*, \{p_i^*\}) \log T_t,$$

that is,

$$\sum_{i \neq 1} \frac{\mathbb{E}[N_i(T_t)]}{(1 - \delta) \log T_t} (L_i - L_1)^\top p^* < C_1^*(p^*, \{p_i^*\}).$$

From the definition of  $C_1^*$ , there exists  $q'_t \in \{q \in \text{cl}(\mathcal{C}_1^c) : p_1^* = S_j q\} =: \mathcal{S}$  such that

$$\sum_{i \neq 1} \frac{\mathbb{E}[N_i(T_t)]}{(1 - \delta) \log T_t} D(p_i^* \| S_i q'_t) < 1.$$

Since  $\mathcal{S}$  is compact, there exists a subsequence  $t_0 < t_1 < \dots$  such that  $\lim_{u \rightarrow \infty} q'_{t_u} = q'$  for some  $q' \in \mathcal{S}$ . Therefore from the lower semicontinuity of the divergence we obtain

$$\begin{aligned} 1 &\geq \sum_{i \neq 1} \liminf_{u \rightarrow \infty} \frac{\mathbb{E}[N_i(T_t)]}{(1 - \delta) \log T_t} D(p_i^* \| S_i q'_{t_u}) \\ &\geq \sum_{i \neq 1} \liminf_{t \rightarrow \infty} \frac{\mathbb{E}[N_i(T_t)]}{(1 - \delta) \log T_t} D(p_i^* \| S_i q') \\ &= \sum_i \liminf_{t \rightarrow \infty} \frac{\mathbb{E}[N_i(T_t)]}{(1 - \delta) \log T_t} D(p_i^* \| S_i q'), \end{aligned}$$

which contradicts Lemma 1.  $\square$

## C The $N$ -armed Bandit Problem as Partial Monitoring

In Section 6, we have introduced PM-DMED-Hinge, an asymptotically optimal algorithm for partial monitoring. In this appendix, we prove that this algorithm also has an optimal regret bound of the  $N$ -armed bandit problem when we run it on an  $N$ -armed bandit game represented as an instance of partial monitoring.

In the  $N$ -armed bandit problem, the learner selects one of  $N$  actions (arms) and receives a corresponding reward. This problem can be considered as a special case of partial monitoring in which the learner directly observes the loss matrix. For example, three-armed Bernoulli bandit can be represented by the following loss and feedback matrices, and the strategy:

$$L = H = \begin{pmatrix} 2 & 1 & 2 & 1 & 2 & 1 & 2 & 1 \\ 2 & 2 & 1 & 1 & 2 & 2 & 1 & 1 \\ 2 & 2 & 2 & 2 & 1 & 1 & 1 & 1 \end{pmatrix}, \text{ and } p^* = \begin{pmatrix} (1 - \mu_1)(1 - \mu_2)(1 - \mu_3) \\ \mu_1(1 - \mu_2)(1 - \mu_3) \\ (1 - \mu_1)\mu_2(1 - \mu_3) \\ \mu_1\mu_2(1 - \mu_3) \\ (1 - \mu_1)(1 - \mu_2)\mu_3 \\ \mu_1(1 - \mu_2)\mu_3 \\ (1 - \mu_1)\mu_2\mu_3 \\ \mu_1\mu_2\mu_3 \end{pmatrix}, \quad (15)$$

where  $\mu_1, \mu_2$ , and  $\mu_3$  are the expected rewards of the actions. Signals 1 and 2 correspond to the rewards of 1 and 0 generated by the selected arm, respectively. More generally,  $N$ -armed Bernoulli

bandit is represented as an instance of partial monitoring in which the loss and feedback matrices are the same  $N \times 2^N$  matrix

$$l_{i,j} = h_{i,j} = \mathbb{1}[(j-1 \bmod 2^i) < 2^{i-1}] + 1,$$

where  $\bmod$  denotes the modulo operation. This problem is associated with  $N$  parameters  $\mu_1, \mu_2, \dots, \mu_N$  that correspond to the expected rewards of the actions. For the ease of analysis, we assume  $\{\mu_i\}$  are in  $(0, 1)$  and different from each other. Without loss of generality, we assume  $1 > \mu_1 > \mu_2 > \dots > \mu_N > 0$ , and thus action 1 is the optimal action. The opponent's strategy is

$$p_j^* = \prod_{i \in [N]} (\mu_i + (1 - 2\mu_i) \mathbb{1}[(j-1 \bmod 2^i) < 2^{i-1}]).$$

Note that  $\mu_i = (S_i p^*)_1$ .

*Proof of Corollary 4.* In the following, we prove that the regularity conditions in Theorem 3 are always satisfied in the case of the  $N$ -armed bandit. During the proof we also show that  $C_1^*(p^*, \{p_i^*\})$  is equal to the optimal constant factor of Lai and Robbins [2].

Because signal 1 corresponds to the reward of 1, we can define  $\hat{\mu}_i(q) = (S_i q)_1$ , and thus

$$\mathcal{C}_i = \{q \in \mathcal{P}_M : \forall i' \neq i, \hat{\mu}_i(q) \geq \hat{\mu}_{i'}(q)\}.$$

First, we show the uniqueness of  $\mathcal{R}_1^*(p, \{p_i, \delta_i\})$  at  $p = p^*, \{p_i\} = S_i p^*, \delta_i = 0$ . It is easy to check

$$D(p_i^* \| S_i q) = d(\hat{\mu}_i(p^*) \| \hat{\mu}_i(q)) = d(\mu_i \| \hat{\mu}_i(q)),$$

where  $d(a \| b)$  is the KL divergence between two Bernoulli distributions with parameters  $a$  and  $b$ . Then

$$\begin{aligned} \mathcal{R}_1(\{p_i^*\}) &= \left\{ \{r_i\}_{i \neq 1} \in [0, \infty)^{N-1} : \inf_{q \in \text{cl}(\mathcal{C}_1^c) : p_i^* = S_i q} \sum_i r_i D(p_i^* \| S_i q) \geq 1 \right\} \\ &= \left\{ \{r_i\}_{i \neq 1} \in [0, \infty)^{N-1} : \inf_{q \in \text{cl}(\mathcal{C}_1^c) : \mu_1 = \hat{\mu}_1(q)} \sum_i r_i D(p_i^* \| S_i q) \geq 1 \right\} \\ &= \left\{ \{r_i\}_{i \neq 1} : r_i \geq \frac{1}{d(\mu_i \| \mu_1)} \right\}, \end{aligned} \quad (16)$$

where the last inequality follows from the fact that

$$\{q \in \text{cl}(\mathcal{C}_1^c) : \hat{\mu}_1(q) = \mu_1\} = \{q \in \mathcal{P}_M : \hat{\mu}_1(q) = \mu_1, \exists i \neq 1, \hat{\mu}_i(q) \geq \mu_1\}.$$

By Eq. (16), the regret minimizing solution is

$$C_1^*(p^*, \{p_i^*\}) = \sum_{i \neq 1} \frac{\Delta_i}{d(\mu_i \| \mu_1)},$$

and

$$\mathcal{R}_1^*(p^*, \{p_i^*\}) = \left\{ \{r_i\}_{i \neq 1} : r_i = \frac{1}{d(\mu_i \| \mu_1)} \right\},$$

which is unique.

Second, we show that  $\text{cl}(\text{int}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta) = \text{cl}(\text{cl}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta)$  for sufficiently small  $\delta \geq 0$ . Note that,

$$\text{cl}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta = \{q \in \mathcal{P}_M : \exists i' \neq 1, \hat{\mu}_1(q) \leq \hat{\mu}_{i'}(q), d(\mu_1 \| \hat{\mu}_1(q)) \leq \delta\}$$

and

$$\text{int}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta = \{q \in \mathcal{P}_M : \exists i' \neq 1, \hat{\mu}_1(q) < \hat{\mu}_{i'}(q), d(\mu_1 \| \hat{\mu}_1(q)) \leq \delta\}.$$

To prove

$$\text{cl}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta \subset \text{cl}(\text{int}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta), \quad (17)$$

it suffices to show that, an open ball centered at any position in

$$\{q \in \mathcal{P}_M : \exists i' \neq 1, \hat{\mu}_1(q) = \hat{\mu}_{i'}(q), d(\mu_1 \| \hat{\mu}_1(q)) \leq \delta\} \supset (\text{cl}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta) \setminus (\text{int}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta)$$

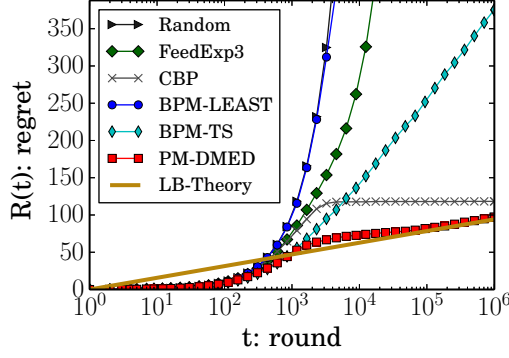


Figure 4: Regret-round semilog plots of algorithms. The regrets are averaged over 100 runs. LB-Theory is the asymptotic regret lower bound of Lai and Robbins [2].

contains a point in  $\text{int}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta$ . This holds because we can make a slight move towards the direction of increasing  $\hat{\mu}_{i'}$ : we can always find  $q'$  in an open ball centered at  $q$  such that  $\hat{\mu}_{i'}(q') > \hat{\mu}_{i'}(q)$  and  $\hat{\mu}_1(q') = \hat{\mu}_1(q)$  because of (i) the fact that there always exists  $q \in \mathcal{P}_M$  such that  $\{q \in \mathcal{P}_M, \forall i \in [N] \hat{\mu}_i(q) = \mu_i\}$  for arbitrary  $\{\mu_i\} \in (0, 1)^N$  and (ii) the continuity of the  $\hat{\mu}_i$  operator. Therefore, any open ball centered at  $q \in \text{cl}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta$  contains an element of  $\text{int}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta$ , by which we obtain (17). By using (17), we have

$$\text{cl}(\text{cl}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta) \subset \text{cl}(\text{cl}(\text{int}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta)) = \text{cl}(\text{int}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta), \quad (18)$$

where we used the fact that  $\text{cl}(\text{cl}(X)) = \text{cl}(X)$ . Combining (18) with the fact that  $\text{cl}(\text{cl}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta) \supset \text{cl}(\text{int}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta)$  yields  $\text{cl}(\text{cl}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta) = \text{cl}(\text{int}(\mathcal{C}_1^c) \cap \mathcal{S}_\delta)$ .

Therefore, in the  $N$ -armed Bernoulli bandit problem, the regularity conditions are always satisfied and  $C_1^*(p^*, \{p_i^*\})$  matches the optimal coefficient of the logarithmic regret bound. From Theorem 3, if we run PM-DMED-Hinge in this game, its expected regret is asymptotically optimal in view of the  $N$ -armed bandit problem.  $\square$

### C.1 Experiment

We also assessed the performance of PM-DMED and other algorithms in solving the three-armed Bernoulli bandit game defined by (15) with parameters  $\mu_1 = 0.4$ ,  $\mu_2 = 0.3$ , and  $\mu_3 = 0.2$ . The settings of the algorithms are the same as that of the main paper. The results of simulations are shown in Figure 4. LB-Theory is the regret lower bound of Lai and Robbins [2], that is,  $\sum_{i \neq 1} \frac{\Delta_i \log t}{d(\mu_i \| \mu_1)}$ . The slope of PM-DMED quickly approaches that of LB-Theory, which is empirical evidence that PM-DMED has optimal performance in  $N$ -armed bandits.

## D Optimality of PM-DMED-Hinge

In this appendix we prove Theorem 3. First we define distances among distributions. For distributions  $p_i, p'_i \in \mathcal{P}_A$  of symbols we use the total variation distance

$$\|p_i - p'_i\| = \frac{1}{2} \sum_{a=1}^A |(p_i)_a - (p'_i)_a|.$$

For distributions  $p, p' \in \mathcal{P}_M$  of outcomes, we identify  $p$  with the set  $\{p' : \forall i, S_i p = S_i p'\}$  and define

$$\|p - p'\|_M = \max_i \|S_i p - S_i p'\|.$$

For  $\mathcal{Q} \subset \mathcal{P}_M$  we define

$$\|p - \mathcal{Q}\|_M = \inf_{p' \in \mathcal{Q}} \|p - p'\|_M.$$

In the following, we use Pinsker's inequality given below many times.

$$D(p_i \| q_i) \geq 2 \| p_i - q_i \|^2.$$

Let

$$\begin{aligned}\rho_{i,L} &= \sup_{\lambda > 0} \frac{1}{\lambda} \min_{x \in \mathcal{C}_{i,\lambda}} \|x - \mathcal{C}_i^c\|_M, \\ \nu_{i,L} &= \sup_{\lambda > 0} \frac{1}{\lambda} \max_{x \in \mathcal{C}_{i,\lambda}^c} \|x - \mathcal{C}_i^c\|_M.\end{aligned}$$

Note that these two constants are positive from the global observability.

### D.1 Properties of regret lower bound

In this section, we give Lemma 5 and Theorem 6 that are about the functions  $C_j^*(p, \{p_i, \delta_i\})$  and  $\mathcal{R}_j^*(p, \{p_i, \delta_i\})$ . In the following, we always consider these functions on  $p \in \mathcal{P}_M$ ,  $p_i \in \{S_i p : \text{supp}(p) \subset \text{supp}(p^*)\}$  and  $\delta_i \geq 0$ , where  $\text{supp}(\cdot)$  denotes the support of the distribution.

We define

$$L_{\max} = \max_{i', j'} l_{i', j'}.$$

**Lemma 5.** *Let  $p \in \mathcal{C}_{j,\alpha}$  and  $\{p_i, \delta_i\}$  be satisfying  $\|p_i - S_i p\| \leq \alpha \rho_{j,L}/2$  and  $\delta_i \leq (\alpha \rho_{j,L})^2/4$  for all  $i$ . Then*

$$C_j^*(p, \{p_i, \delta_i\}) \leq \frac{4NL_{\max}}{(\alpha \rho_{j,L})^2}. \quad (19)$$

Furthermore,  $\mathcal{R}_j(\{p_i, \delta_i\})$  is nonempty and

$$\mathcal{R}_j^*(p, \{p_i, \delta_i\}) \subset \left[0, \frac{4NL_{\max}}{(\rho_{j,L})^2 \alpha^3}\right]^{N-1}.$$

*Proof of Lemma 5.* Since  $\|p - \mathcal{C}_1^c\|_M \geq \alpha \rho_{j,L}$ , there exists  $i = i(q)$  for any  $q \in \mathcal{C}_1^c$  such that

$$\|S_i q - S_i p\| \geq \alpha \rho_{j,L}.$$

For this  $i$  we have

$$\begin{aligned}D(p_i \| S_i q) - \delta_i &\geq 2 \|p_i - S_i q\|^2 - \delta_i \\ &\geq 2(\|S_i q - S_i p\| - \|p_i - S_i p\|_+)^2 - \delta_i \\ &\geq (\alpha \rho_{j,L})^2/2 - \delta_i \\ &\geq (\alpha \rho_{j,L})^2/4.\end{aligned}$$

Thus, by letting  $r_i = 4/(\alpha \rho_{j,L} \alpha)^2$  for all  $i \neq j$  we have

$$\{r_i\}_{i \neq j} \in \mathcal{R}_j(\{p_i, \delta_i\}),$$

which implies (19). On the other hand it holds for any  $\{r_i^*\}_{i \neq j} \in \mathcal{R}_j^*(p, \{p_i, \delta_i\})$  from  $p \in \mathcal{C}_{j,\alpha}$  that

$$C_j^*(p, \{p_i, \delta_i\}) = \sum_{i \neq j} r_i^* L_i^\top p \geq \max_{i \neq j} r_i^* \alpha$$

and therefore we have

$$\max_{i \neq j} r_i^* \leq \frac{4NL_{\max}}{(\rho_{j,L})^2 \alpha^3}.$$

□

**Theorem 6.** Assume that the regularity conditions in Theorem 3 hold. Then the point-to-set map  $\mathcal{R}_1^*(p, \{p_i, \delta_i\})$  is (i) nonempty near  $p = p^*, p_i = S_i p^*, \delta_i = 0$  and (ii) continuous at  $p = p^*, p_i = S_i p^*, \delta_i = 0$ .

See Hogan [25] for definitions of terms such as continuity of point-to-set maps.

*Proof of Theorem 6.* Define

$$\bar{\mathcal{R}}_1(\{p_i, \delta_i\}) = \left\{ \{r_i\}_{i \neq 1} \in [0, \xi]^{N-1} : \inf_{q \in \text{cl}(\mathcal{C}_1^c) : D(p_1 \| S_1 q) \leq \delta_1} \sum_{i \neq 1} r_i (D(p_i \| S_i q) - \delta_i)_+ \geq 1 \right\}$$

for

$$\xi = \frac{4NL_{\max}}{(\rho_{1,L})^2 (\max_{i \neq 1} L_i^\top p^* - L_1^\top p^*)^3}.$$

Note that  $p^* \in \mathcal{C}_{1,\alpha}$  for  $\alpha \leq \max_{i \neq 1} L_i^\top p^* - L_1^\top p^*$ . From Lemma 5, near  $p = p^*, p_i = S_i p^*, \delta_i = 0$ ,

$$\bar{\mathcal{R}}_1(\{p_i, \delta_i\}) \supset \mathcal{R}_1^*(p, \{p_i, \delta_i\})$$

and

$$C_1^*(p, \{p_i, \delta_i\}) = \inf_{\{r_i\}_{i \neq 1} \in \bar{\mathcal{R}}_1(\{p_i, \delta_i\})} \sum_{i \neq 1} r_i (L_i - L_1)^\top p$$

hold. Since the function

$$\sum_i r_i (D(p_i \| S_i q) - \delta_i)_+$$

is continuous in  $\{r_i\}$ ,  $\bar{\mathcal{R}}_1(\{p_i, \delta_i\})$  is a closed set and therefore  $\mathcal{R}_1^*(p, \{p_i, \delta_i\})$  is nonempty near  $p = p^*, p_i = S_i p^*, \delta_i = 0$ .

From the continuity of  $D(p_i \| S_i q)$  at any  $q$  such that  $D(p_i \| S_i q) < \infty$ , we have

$$\begin{aligned} \inf_{q \in \text{cl}(\mathcal{C}_1^c) \cap \mathcal{S}_{\delta_1}} \sum_{i \neq 1} r_i (D(p_i \| S_i q) - \delta_i)_+ &= \inf_{q \in \text{cl}(\text{cl}(\mathcal{C}_1^c) \cap \mathcal{S}_{\delta_1})} \sum_{i \neq 1} r_i (D(p_i \| S_i q) - \delta_i)_+ \\ &= \inf_{q \in \text{cl}(\text{int}(\mathcal{C}_1^c) \cap \mathcal{S}_{\delta_1})} \sum_{i \neq 1} r_i (D(p_i \| S_i q) - \delta_i)_+ \\ &= \inf_{q \in \text{int}(\mathcal{C}_1^c) \cap \mathcal{S}_{\delta_1}} \sum_{i \neq 1} r_i (D(p_i \| S_i q) - \delta_i)_+. \end{aligned}$$

Thus, we have

$$\bar{\mathcal{R}}_1(\{p_i, \delta_i\}) = \left\{ \{r_i\}_{i \neq 1} \in [0, \xi]^{N-1} : \inf_{q \in \text{int}(\mathcal{C}_1^c) : D(p_1 \| S_1 q) \leq \delta_1} \sum_i r_i (D(p_i \| S_i q) - \delta_i)_+ \geq 1 \right\}. \quad (20)$$

Since the objective function  $\sum_{i \neq j} r_i (L_i - L_j)^\top p$  is continuous in  $\{r_i\}$  and  $p$ , and (20) is compact, now it suffices to show that (20) is continuous in  $\{p_i, \delta_i\}$  at  $\{S_i p^*, 0\}$  to prove the theorem from [25, Corollary 8.1].

First we show that  $\bar{\mathcal{R}}_1(\{p_i, \delta_i\})$  is closed at  $\{S_i p^*, 0\}$ . Consider  $\{r_i^{(m)}\}_{i \neq 1} \in \bar{\mathcal{R}}_1(\{p_i^{(m)}, \delta_i^{(m)}\})$  for a sequence  $\{p_i^{(m)}, \delta_i^{(m)}\}_i$  which converges to  $\{S_i p^*, 0\}_i$  as  $m \rightarrow \infty$ . We show that  $\{r_i\}_{i \neq 1} \in \bar{\mathcal{R}}_1(\{S_i p^*, 0\})$  if  $r_i^{(m)} \rightarrow r_i$  as  $m \rightarrow \infty$ .

Take an arbitrary  $q \in \text{int}(\mathcal{C}_1^c)$  such that  $D(S_1 p^* \| S_1 q) = 0$ . Since  $\|S_1 p^* - p_1^{(m)}\| \rightarrow 0$  and  $p_1 \in \{S_1 p : \text{supp}(p) \subset \text{supp}(p^*)\}$ , there exists  $\tilde{p}^{(m)}$  such that  $p_1^{(m)} = S_1 \tilde{p}^{(m)}$  and  $\|p^* - \tilde{p}^{(m)}\|_M \rightarrow 0$ .



Thus, from  $q \in \text{int}(\mathcal{C}_1^c)$ , it holds for sufficiently large  $m$  that  $q^{(m)} = q - p^* + \tilde{p}^{(m)} \in \text{int}(\mathcal{C}_1^c)$ . For this  $q^{(m)}$  we have

$$D(p_1^{(m)} \| S_1 q^{(m)}) \leq D(S_1 \tilde{p}^{(m)} \| S_1 (q - p^* + \tilde{p}^{(m)})) = 0 \leq \delta_1.$$

That is,  $q^{(m)} \in \text{int}(\mathcal{C}_1^c) \cap \mathcal{S}_{\delta_1}$ . Therefore, for sufficiently large  $m$  we have

$$\sum_i r_i^{(m)} (D(p_i \| S_i q^{(m)}) - \delta_i^{(m)})_+ \geq 0$$

and, letting  $m \rightarrow \infty$ ,

$$\sum_i r_i D(p_i \| S_i q) \geq 0.$$

This means that  $\{r_i\}_{i \neq 1} \in \bar{\mathcal{R}}_1(\{p_i, \delta_i\})$ , that is,  $\bar{\mathcal{R}}_1(\{p_i, \delta_i\})$  is closed at  $\{S_i p^*, 0\}$ .

Next we show that  $\bar{\mathcal{R}}_1(\{p_i, \delta_i\})$  is open at  $\{S_i p^*, 0\}$ . Consider  $\{r_i\}_{i \neq 1} \in \bar{\mathcal{R}}_1(\{S_i p^*, 0\})$  and a sequence  $\{p_i^{(m)}, \delta_i^{(m)}\}_i$  which converges to  $\{S_i p^*, 0\}_i$  as  $m \rightarrow \infty$ . We show that there exists a sequence  $\{r_i^{(m)}\}_{i \neq 1} \in \bar{\mathcal{R}}_1(\{p_i^{(m)}, \delta_i^{(m)}\})$  such that  $r_i^{(m)} \rightarrow r_i$ .

Consider the optimal value function

$$v(\{p_i^{(m)}, \delta_i^{(m)}\}) = \inf_{q \in \text{cl}(\mathcal{C}_1^c) \cap \mathcal{S}_{\delta_1}} \sum_i r_i (D(p_i^{(m)} \| S_i q) - \delta_i^{(m)})_+. \quad (21)$$

Since the feasible region of (21) is closed at  $p_i = S_i p^*, \delta_i = 0$  and the objective function of (21) is lower semicontinuous in  $q, \{p_i, \delta_i\}$  we see that  $v(\{p_i^{(m)}, \delta_i^{(m)}\})$  is lower semicontinuous from [25, Theorem 2]. Therefore, for any  $\epsilon > 0$  there exists  $m_0 > 0$  such that for all  $m \geq m_0$

$$v(\{p_i^{(m)}, \delta_i^{(m)}\}) \geq (1 - \epsilon)v(\{S_i p^*, 0\}) \geq 1$$

since  $v(\{S_i p^*, 0\}) \geq 1$  from  $r_i \in \bar{\mathcal{R}}_1(\{S_i p^*, 0\})$ . Thus, by letting  $r_i^{(m)} := r_i / (1 - \epsilon)$  we have

$$\inf_{v \in \text{cl}(\mathcal{C}_1^c) \cap \mathcal{S}_{\delta_1}} \sum_i r_i^{(m)} (D(p_i^{(m)} \| S_i q^{(m)}) - \delta_i^{(m)})_+ \geq 1,$$

that is,  $\{r_i^{(m)}\}_{i \neq 1} \in \bar{\mathcal{R}}_1(\{p_i^{(m)}, \delta_i^{(m)}\})$ . □

## D.2 Regret analysis of PM-DMED-Hinge

Let  $\hat{p}_{i,n} \in [0, 1]^A$  be the empirical distribution of the symbols from the action  $i$  when the action  $i$  is selected  $n$  times. Then we have  $\hat{p}_i(t) = \hat{p}_{i,N_i(t)}$ . Let  $P_{i,n_i}(u) = \Pr[D(\hat{p}_{i,n_i} \| S_i p^*) \geq u]$ . Then, from the large deviation bound on discrete distributions (Theorem 11.2.1 in Cover and Thomas [26]), we have

$$P_{i,n_i}(u) \leq (n_i + 1)^A e^{-n_i u}. \quad (22)$$

We also define

$$\mathcal{H}(\{p_i, n_i\}) = \{i \in [N] : D(p_i \| S_i p^*) - f(n_i) > 0\}.$$

For

$$0 < \delta \leq \|p^* - \mathcal{C}_1^c\|_M^2 / 8 \quad (23)$$

define events

$$\begin{aligned}
\mathcal{A}(t) &= \{\hat{p}(t) \in \mathcal{C}_1\} \\
\mathcal{A}'(t) &= \{\hat{p}(t) \in \mathcal{C}_{1,\alpha(t)}\} \\
\mathcal{B}(t) &= \bigcap_i \{\|\hat{p}_i(t) - S_i p^*\| \leq \sqrt{\delta}\} \\
\mathcal{C}(t) &= \{\hat{i}(t) \notin \mathcal{H}(\{\hat{p}_i(t), N_i(t)\}), \mathcal{H}(\{\hat{p}_i(t), N_i(t)\}) \neq \emptyset\} \\
&= \left\{ D(\hat{p}_{\hat{i}(t)}(t) \| S_{\hat{i}(t)} p^*) \leq f(N_{\hat{i}(t)}(t)), \bigcup_i \{D(\hat{p}_i(t) \| S_i p^*) > f(N_i(t))\} \right\} \\
\mathcal{D}(t) &= \bigcap_i \{D(\hat{p}_i(t) \| S_i \hat{p}(t)) \leq f(N_i(t))\} \\
\mathcal{E}(t) &= \left\{ \max_i f(N_i(t)) \leq \min \{2\delta, (\rho_{1,L}\alpha(t))^2/4\}, \right. \\
&\quad \left. \min_i N_i(t) \geq \max\{c\sqrt{\log t}, (\log \log T)^{1/3}\}, 2\nu_{1,L}\alpha(t) \leq \|p^* - \mathcal{C}_1^c\|_M \right\}, \tag{24}
\end{aligned}$$

where we write  $\{\mathcal{T}, \mathcal{U}\}$  instead of  $\{\mathcal{T} \cap \mathcal{U}\}$  for events  $\mathcal{T}$  and  $\mathcal{U}$ .

*Proof of theorem 3.* Since  $\mathcal{A}'(t) \subset \mathcal{A}(t)$ , the whole sample space is covered by

$$\begin{aligned}
&\{\mathcal{A}'(t), \mathcal{B}(t)\} \cup \{\mathcal{A}'(t), \mathcal{B}^c(t)\} \cup \{\mathcal{A}(t), (\mathcal{A}'(t))^c\} \cup \{\mathcal{A}^c(t), \mathcal{C}(t)\} \cup \{\mathcal{A}^c(t), \mathcal{C}^c(t)\} \\
&\subset \{\mathcal{A}'(t), \mathcal{B}(t), \mathcal{D}(t), \mathcal{E}(t)\} \cup \{\mathcal{A}'(t), \mathcal{B}^c(t), \mathcal{D}(t), \mathcal{E}(t)\} \cup \{\mathcal{A}(t), (\mathcal{A}'(t))^c, \mathcal{D}(t), \mathcal{E}(t)\} \cup \{\mathcal{A}^c(t), \mathcal{C}(t)\} \\
&\quad \cup \{\mathcal{A}^c(t), \mathcal{C}^c(t), \mathcal{D}(t), \mathcal{E}(t)\} \cup \mathcal{D}^c(t) \cup \mathcal{E}^c(t). \tag{25}
\end{aligned}$$

Let  $J_i(t)$  denote the event that action  $i$  is newly added into the list  $L_N$  at the  $t$ -th round and  $J'_i(t) \subset J_i(t)$  denote the event that  $J_i(t)$  occurred by Step 6 of Algorithm 3. Note that if  $\{\mathcal{A}'(t), \mathcal{D}(t), \mathcal{E}(t)\}$  occurred then  $J_i(t)$  is equivalent to  $J'_i(t)$ . Combining this fact with (25) we can bound the regret as

$$\begin{aligned}
\text{Regret}(T) &\leq \sum_{i \neq 1} \Delta_i \sum_{t=1}^T \mathbb{1}[J_i(t)] + N \\
&\leq \sum_{i \neq 1} \Delta_i \sum_{t=1}^T \left( \mathbb{1}[J'_i(t), \mathcal{A}'(t), \mathcal{B}(t), \mathcal{D}(t), \mathcal{E}(t)] + \mathbb{1}[J'_i(t), \mathcal{A}(t), \mathcal{B}^c(t), \mathcal{D}(t), \mathcal{E}(t)] \right. \\
&\quad \left. + \mathbb{1}[J_i(t), \mathcal{A}(t), (\mathcal{A}'(t))^c, \mathcal{D}(t), \mathcal{E}(t)] + \mathbb{1}[J_i(t), \mathcal{A}^c(t), \mathcal{C}^c(t), \mathcal{D}(t), \mathcal{E}(t)] \right. \\
&\quad \left. + \mathbb{1}[J_i(t), \mathcal{D}^c(t) \cup \mathcal{E}^c(t)] \right) + \left( \sum_{i \neq 1} \Delta_i \right) \sum_{t=1}^T \mathbb{1}[\mathcal{A}^c(t), \mathcal{C}(t)] + N.
\end{aligned}$$

The following Lemmas 7–13 bound the expectation of each term and complete the proof.  $\square$

**Lemma 7.** Let  $\{r_i^*\}_{i \neq 1}$  be the unique member of  $\mathcal{R}_j^*(p^*, \{S_i p^*, 0\})$ . Then there exists  $\epsilon_\delta > 0$  such that  $\lim_{\delta \rightarrow 0} \epsilon_\delta = 0$  and for all  $i \neq 1$

$$\sum_{t=1}^T \mathbb{1}[J'_i(t), \mathcal{A}'(t), \mathcal{B}(t), \mathcal{D}(t), \mathcal{E}(t)] \leq (1 + \epsilon_\delta) r_i^* \log T + 1.$$

**Lemma 8.**

$$\mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}[J'_i(t), \mathcal{A}'(t), \mathcal{B}^c(t), \mathcal{D}(t), \mathcal{E}(t)] \right] = o(\log T).$$

**Lemma 9.**

$$\mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}[\mathcal{A}^c(t), \mathcal{C}(t)] \right] = O(1).$$

**Lemma 10.**

$$\mathbb{E} \left[ \sum_{t=1}^T \mathbb{1} [J_i(t), \mathcal{A}(t), (\mathcal{A}'(t))^c, \mathcal{D}(t), \mathcal{E}(t)] \right] = O(1).$$

**Lemma 11.**

$$\mathbb{E} \left[ \sum_{t=1}^T \mathbb{1} [J_i(t), \mathcal{A}^c(t), \mathcal{C}^c(t), \mathcal{D}(t), \mathcal{E}(t)] \right] = O(1).$$

**Lemma 12.**

$$\mathbb{E} \left[ \sum_{t=1}^T \mathbb{1} [J_i(t), \mathcal{D}^c(t)] \right] = O(1).$$

**Lemma 13.**

$$\sum_{t=1}^T \mathbb{1} [J_i(t), \mathcal{E}^c(t)] = o(\log T).$$

*Proof of Lemma 7.* From  $\mathcal{D}(t)$  we have

$$\sum_i N_i(t) (D(\hat{p}_i(t) \| S_i \hat{p}(t)) - f(N_i(t)))_+ = 0. \quad (26)$$

Here assume that  $\|\hat{p}(t) - p^*\|_M > 2\sqrt{\delta}$ . Then

$$\begin{aligned} \max_i D(\hat{p}_i(t) \| S_i \hat{p}(t)) &\geq 2 \max_i \|\hat{p}_i(t) - S_i \hat{p}(t)\|^2 \quad (\text{by Pinsker's inequality}) \\ &\geq 2 \max_i (\|S_i p^* - S_i \hat{p}(t)\| - \|S_i p^* - \hat{p}_i(t)\|)_+^2 \\ &\geq 2 \max_i (\|S_i p^* - S_i \hat{p}(t)\| - \sqrt{\delta})_+^2 \quad (\text{by definition of } \mathcal{B}(t)) \\ &> 2\delta \\ &\geq f(N_i(t)), \quad (\text{by definition of } \mathcal{E}(t)) \end{aligned}$$

which contradicts (26) and we obtain  $\|\hat{p}(t) - p^*\|_M \leq 2\sqrt{\delta}$ . Furthermore, from  $\mathcal{B}(t)$  and  $\mathcal{E}(t)$  we have

$$\bigcap_i \{\|\hat{p}_i(t) - S_i p^*\| \leq \sqrt{\delta}\} \text{ and } \bigcap_i \{f(N_i(t)) \leq 2\delta\},$$

respectively. Since  $\mathcal{R}_1^*(p, \{p_i, \delta_i\})$  is continuous at  $p = p^*$ ,  $p_i = S_i p^*$ ,  $\delta_i = 0$  from Theorem 6,  $r_i \leq (1 + \epsilon_\delta) r_i^*$  for all  $\{r_i\}_{i \neq 1} \in \mathcal{R}_{\hat{i}(t)}^*(\hat{p}(t), \{\hat{p}_i(t), f(N_i(t))\})$  where  $r_i^*$  is the unique member of  $\mathcal{R}_1^*(p^*, \{S_i p^*, 0\})$  and we used the fact that  $\mathcal{A}'(t)$  implies  $\hat{i}(t) = 1$ .

We complete the proof by

$$\begin{aligned} &\sum_{t=1}^T \mathbb{1} [J'_i(t), \mathcal{A}'(t), \mathcal{B}(t), \mathcal{D}(t), \mathcal{E}(t)] \\ &= \sum_{n=1}^T \mathbb{1} \left[ \bigcup_{t=1}^T \{J'_i(t), \mathcal{A}'(t), \mathcal{B}(t), \mathcal{D}(t), \mathcal{E}(t), N_i(t) = n\} \right] \\ &\leq \sum_{n=1}^T \mathbb{1} \left[ \bigcup_{t=1}^T \{n / \log t \leq (1 + \epsilon_\delta) r_i^*\} \right] \\ &\leq (1 + \epsilon_\delta) r_i^* \log T + 1. \end{aligned}$$

□

*Proof of Lemma 8.* First, we obtain from  $\mathcal{D}(t)$  and  $\mathcal{E}(t)$  that  $f(N_i(t)) \leq (\rho_{1,L}\alpha(t))^2/4$  and

$$\begin{aligned}\|\hat{p}_i(t) - S_i\hat{p}(t)\| &\leq \sqrt{D(\hat{p}_i(t)\|S_i\hat{p}(t))/2} \\ &\leq \sqrt{f(N_i(t))/2} \\ &\leq \rho_{1,L}\alpha(t)/\sqrt{8}.\end{aligned}$$

Therefore, from Lemma 5, it holds for any  $\{r_i^*\}_{i \neq 1} \in \mathcal{R}_j^*(\hat{p}(t), \{\hat{p}_i(t), f(N_i(t))\})$  that

$$\begin{aligned}r_i^* &\leq \frac{4NL_{\max}}{(\rho_{1,L})^2(\alpha(t))^3} \\ &\leq \frac{4NL_{\max}}{(\rho_{1,L})^2(\alpha(T))^3}.\end{aligned}$$

Now we have

$$\begin{aligned}&\mathbb{E} \left[ \sum_{t=1}^{\infty} \mathbb{1} [J'_i(t), \mathcal{A}'(t), \mathcal{B}^c(t), \mathcal{D}(t), \mathcal{E}(t)] \right] \\ &\leq \mathbb{E} \left[ \sum_{t=1}^{\infty} \mathbb{1} \left[ \frac{N_i(t)}{\log T} < \frac{4NL_{\max}}{(\rho_{1,L})^2(\alpha(T))^3}, \mathcal{B}^c(t), \mathcal{E}(t) \right] \right] \\ &\leq \left( \frac{4NL_{\max} \log T}{(\rho_{1,L})^2(\alpha(T))^3} + 1 \right) \Pr \left[ \bigcup_{t=1}^T \{\mathcal{B}^c(t), \mathcal{E}(t)\} \right].\end{aligned}\tag{27}$$

Here, note that

$$\begin{aligned}\mathcal{B}^c(t) &\subset \bigcup_i \{\|\hat{p}_i(t) - S_i p^*\| \geq \sqrt{\delta}\} \\ &\subset \bigcup_i \{D(\hat{p}_i(t)\|S_i p^*) \geq 2\delta\}.\end{aligned}$$

Since  $N_i(t) \geq (\log \log T)^{1/3}$  holds under event  $\mathcal{E}(t)$ , we can bound the probability in (27) as

$$\begin{aligned}&\Pr \left[ \bigcup_{t=1}^T \{\mathcal{B}^c(t), \mathcal{E}(t)\} \right] \\ &\leq \sum_i \sum_{n_i=(\log \log T)^{1/3}}^{\infty} \Pr[D(\hat{p}_{i,n_i}\|S_i p^*) \geq 2\delta] \\ &\leq N \sum_{n=(\log \log T)^{1/3}}^{\infty} (n+1)^A e^{-2n\delta} \quad (\text{by (22)}) \\ &= e^{-\Theta((\log \log T)^{1/3})}\end{aligned}$$

and combining this with (27) we have

$$\begin{aligned}&\mathbb{E} \left[ \sum_{t=1}^{\infty} \mathbb{1} [J'_i(t), \mathcal{A}'(t), \mathcal{B}^c(t), \mathcal{D}(t), \mathcal{E}(t)] \right] \\ &\leq O((\log T)(\log \log T)^3) e^{-\Theta((\log \log T)^{1/3})} \\ &= o(\log T).\end{aligned}$$

□

*Proof of Lemma 9.* Let  $\mathcal{G} \in 2^{[N]} \setminus \emptyset$  and  $\{n_i\}_{i \in \mathcal{G}} \in \mathbb{N}^{|\mathcal{G}|}$  be arbitrary. Consider the case that

$$\sum_{i \in \mathcal{G}} n_i (D(\hat{p}_{i,n_i}\|S_i p^*) - f(n_i))_+ < x.\tag{28}$$

for some  $x > 0$ . Then under events  $t \geq e^x$ ,  $\bigcap_{i \in \mathcal{G}} \{N_i(t) = n_i\}$ ,  $\mathcal{A}^c(t)$ ,  $\mathcal{C}(t)$  and  $\mathcal{H}(\{\hat{p}_i(t), f(n_i)\}) = \mathcal{G}$  we have

$$\begin{aligned} \min_{p \in \mathcal{C}_{i(t)}^c : D(\hat{p}_{i(t)}(t) \| S_{i(t)} p) \leq f(N_{i(t)}(t))} \sum_i N_i(t) (D(\hat{p}_i(t) \| S_i p) - f(N_i(t)))_+ \\ \leq \sum_i n_i (D(\hat{p}_i(t) \| S_i p^*) - f(n_i))_+ < x \leq \log t, \end{aligned}$$

which implies that the condition (9) is satisfied. On the other hand from (10),  $\{r_i^*\}$  satisfies

$$\sum_{i \in \mathcal{G}} (r_i^* \log t) (D(\hat{p}_i(t) \| S_i p) - f(n_i))_+ \geq \log t. \quad (29)$$

Eqs. (28) and (29) imply that there exists at least one  $i \in \mathcal{G}$  such that  $r_i^* \log t > N_i(t) = n_i$ . This action is selected within  $N$  rounds and therefore  $N_i(t') = n_i$  never holds for all  $t' \geq t + N$ . Thus, under the condition (28) it holds that

$$\sum_t \mathbb{1} \left[ \mathcal{A}^c(t), \mathcal{C}(t), \mathcal{H}(\{\hat{p}_i(t), N_i(t)\}) = \mathcal{G}, \bigcap_{i \in \mathcal{G}} \{N_i(t) = n_i\} \right] \leq e^x + N.$$

By using this inequality we have

$$\begin{aligned} \sum_{t=1}^{\infty} \mathbb{1} [\mathcal{A}^c(t), \mathcal{C}(t)] \\ \leq \sum_{\mathcal{G} \in 2^{[N]} \setminus \emptyset} \sum_{\{n_i\}_{i \in \mathcal{G}} \in \mathbb{N}^{|\mathcal{G}|}} \sum_{t=1}^{\infty} \mathbb{1} \left[ \mathcal{A}^c(t), \mathcal{C}(t), \mathcal{H}(\{\hat{p}_i(t), N_i(t)\}) = \mathcal{G}, \bigcap_{i \in \mathcal{G}} \{N_i(t) = n_i\} \right] \\ \leq \sum_{\mathcal{G} \in 2^{[N]} \setminus \emptyset} \sum_{\{n_i\}_{i \in \mathcal{G}} \in \mathbb{N}^{|\mathcal{G}|}} \mathbb{1} \left[ \bigcap_{i \in \mathcal{G}} \{D(\hat{p}_{i,n_i} \| S_i p^*) \geq f(n_i)\} \right] \left( \exp \left( \sum_{i \in \mathcal{G}} n_i (D(\hat{p}_{i,n_i} \| S_i p^*) - f(n_i)) \right) + N \right). \end{aligned} \quad (30)$$

Let

$$D_i = \sup_n \{\text{ess sup } D(\hat{p}_{i,n} \| S_i p^*)\} = -\log \min_{j: (S_i p^*)_j > 0} (S_i p^*)_j,$$

where  $(S_i p^*)_j$  is the  $j$ -th component of  $S_i p^*$ . Then,

$$\begin{aligned} \mathbb{E} \left[ \sum_{\{n_i\}_{i \in \mathcal{G}} \in \mathbb{N}^{|\mathcal{G}|}} \mathbb{1} \left[ \bigcap_{i \in \mathcal{G}} \{D(\hat{p}_{i,n_i} \| S_i p^*) \geq f(n_i)\} \right] \left( \exp \left( \sum_{i \in \mathcal{G}} n_i (D(\hat{p}_{i,n_i} \| S_i p^*) - f(n_i)) \right) + N \right) \right] \\ \leq \sum_{\{n_i\}_{i \in \mathcal{G}} \in \mathbb{N}^{|\mathcal{G}|}} \left( \prod_{i \in \mathcal{G}} \int_{f(n_i)}^{D_i} e^{n_i(u_i - f(n_i))} d(-P_{i,n_i}(u_i)) + N \prod_{i \in \mathcal{G}} (n_i + 1)^A e^{-n_i f(n_i)} \right). \end{aligned}$$

The first integral is bounded as

$$\begin{aligned} \int_{f(n_i)}^{D_i} e^{n_i(u_i - f(n_i))} d(-P_{i,n_i}(u_i)) \\ = \left[ -e^{n_i(u_i - f(n_i))} P_{i,n_i}(u_i) \right]_{f(n_i)}^{D_i} + \int_{f(n_i)}^{D_i} n_i e^{n_i(u_i - f(n_i))} P_{i,n_i}(u_i) du_i \quad (\text{integration by parts}) \\ \leq (n_i + 1)^A e^{-n_i f(n_i)} + \int_{f(n_i)}^{D_i} n_i (n_i + 1)^A e^{-n_i f(n_i)} du_i \\ \leq (1 + n_i D_i) (n_i + 1)^A e^{-n_i f(n_i)}. \end{aligned} \quad (31)$$

Putting (30)–(31) together we have

$$\begin{aligned}
& \mathbb{E} \left[ \sum_{t=1}^{\infty} \mathbb{1} [\mathcal{A}^c(t), \mathcal{C}(t)] \right] \\
& \leq \sum_{\mathcal{G} \in 2^{[N]} \setminus \emptyset} \sum_{\{n_i\}_{i \in \mathcal{G}} \in \mathbb{N}^{|\mathcal{G}|}} \left( \prod_{i \in \mathcal{G}} (1 + n_i D_i) (n_i + 1)^A e^{-n_i f(n_i)} + N \prod_{i \in \mathcal{G}} (n_i + 1)^{|A|} e^{-n_i f(n_i)} \right) \\
& \leq (N + 1) \sum_{\mathcal{G} \in 2^{[N]} \setminus \emptyset} \prod_{i \in \mathcal{G}} \sum_{n_i \in \mathbb{N}} (1 + n_i D_i) (n_i + 1)^A e^{-n_i f(n_i)} \\
& = O(1). \quad (\text{by } n_i f(n_i) = \Theta(n_i^{1/2}))
\end{aligned}$$

□

*Proof of Lemma 10.* Because  $\mathcal{A}(t)$ ,  $(\mathcal{A}'(t))^c$  and  $\mathcal{E}(t)$  imply

$$\begin{aligned}
\|p^* - \hat{p}(t)\|_M & \geq \sup_{p \in \mathcal{C}_1^c} \{\|p^* - p\|_M - \|\hat{p}(t) - p\|_M\} \\
& \geq \sup_{p \in \mathcal{C}_1^c} \{\|p^* - \mathcal{C}_1^c\|_M - \|\hat{p}(t) - p\|_M\} \\
& \geq \|p^* - \mathcal{C}_1^c\|_M - \nu_{1,L} \alpha(t) \\
& \geq \|p^* - \mathcal{C}_1^c\|_M / 2,
\end{aligned}$$

$(\mathcal{A}'(t))^c$ ,  $\mathcal{D}(t)$  and  $\mathcal{E}(t)$  imply

$$\begin{aligned}
\max_i D(\hat{p}_i(t) \| S_i p^*) & \geq 2 \max_i \|\hat{p}_i(t) - S_i p^*\|^2 \\
& \geq 2 \max_i (\|S_i p^* - S_i \hat{p}(t)\| - \|\hat{p}_i(t) - S_i \hat{p}(t)\|)_+^2 \\
& \geq 2 \max_i \left( \|S_i p^* - S_i \hat{p}(t)\| - \sqrt{f(N_i(t))/4} \right)_+^2 \\
& \geq 2(\|p^* - \mathcal{C}_1^c\|_M / 2 - \sqrt{\delta/2})_+^2 \\
& \geq \|p^* - \mathcal{C}_1^c\|_M^2 / 8. \quad (\text{by (23)})
\end{aligned}$$

On the other hand, event  $\{J_i(t), \mathcal{A}^c(t), \mathcal{A}'(t), \min_j N_j(t) = n\}$  occurs for at most twice since all actions are put into the list if  $\{\mathcal{A}^c(t), \mathcal{A}'(t)\}$  occurred. Thus, we have

$$\begin{aligned}
& \mathbb{E} \left[ \sum_n \mathbb{1} [J_i(t), \mathcal{A}(t), (\mathcal{A}'(t))^c, \mathcal{D}(t), \mathcal{E}(t)] \right] \\
& \leq 2 \mathbb{E} \left[ \sum_n \mathbb{1} \left[ \bigcup_t \{ \mathcal{A}(t), (\mathcal{A}'(t))^c, \mathcal{D}(t), \mathcal{E}(t), \min_j N_j(t) = n \} \right] \right] \\
& \leq 2 \sum_n \Pr \left[ \max_j D(\hat{p}_j(t) \| S_j p^*) \geq \|p^* - \mathcal{C}_1^c\|_M^2 / 8, \bigcap_j \{N_j(t) \geq n\} \right] \\
& \leq 2N \sum_n (n + 1)^A e^{-n \|p^* - \mathcal{C}_1^c\|_M^2 / 8} \quad (\text{by (22)}) \\
& = O(1).
\end{aligned}$$

□

*Proof of Lemma 11.* Recall that

$$\mathcal{C}^c(t) = \left\{ \{D(\hat{p}_{i(t)}(t) \| S_{i(t)} p^*) > f(N_{i(t)}(t))\} \cup \bigcap_j \{D(\hat{p}_j(t) \| S_j p^*) \leq f(N_j(t))\} \right\}.$$

Here

$$\left\{ \mathcal{A}^c(t), \mathcal{D}(t), \mathcal{E}(t), \bigcap_j \{D(\hat{p}_j(t)\|S_j p^*) \leq f(N_j(t))\} \right\} \quad (32)$$

cannot occur since (32) implies that

$$\begin{aligned} \|\hat{p}(t) - p^*\|_M &= \max_j \|S_j \hat{p}(t) - S_j p^*\| \\ &\leq \max_j (\|S_j \hat{p}(t) - \hat{p}_j(t)\| + \|\hat{p}_j(t) - S_j p^*\|) \\ &\leq \max_j \left( \sqrt{D(\hat{p}_j(t)\|S_j \hat{p}(t))/2} + \sqrt{D(\hat{p}_j(t)\|S_j p^*)/2} \right) \\ &\leq \sqrt{2 \max_j f(N_j(t))} \quad (\text{by } \mathcal{D}(t)) \\ &\leq 2\sqrt{\delta} \quad (\text{by (24)}) \\ &\leq \|p^* - \mathcal{C}_1^c\|_M / \sqrt{2} \quad (\text{by (23)}), \end{aligned}$$

which contradicts  $\hat{p}(t) \in \mathcal{C}_1^c$ .

On the other hand,  $\mathbb{1} \left[ J_i(n), \hat{i}(t) = j, D(\hat{p}_j(t)\|S_j p^*) > f(N_j(t)), N_j(t) = n_j \right]$  occurs for at most twice since  $\hat{i}(t)$  is put into the list under this event. Thus, we have

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=1}^{\infty} \mathbb{1} [J_i(t), \mathcal{A}^c(t), \mathcal{C}^c(t), \mathcal{D}(t), \mathcal{E}(t)] \right] &\leq 2 \sum_j \sum_{n=1}^{\infty} \Pr[D(\hat{p}_{j,n}\|S_j p^*) > f(n)] \\ &\leq 2 \sum_j \sum_{n=1}^{\infty} (n+1)^A e^{-nf(n)} \quad (\text{by (22)}) \\ &= O(1). \end{aligned}$$

□

*Proof of Lemma 12.*  $\mathcal{D}^c(t)$  implies

$$\begin{aligned} 0 &< \min_p \sum_j N_j(t) (D(\hat{p}_j(t)\|S_j p) - f(N_j(t)))_+ \\ &\leq \sum_j N_j(t) (D(\hat{p}_j(t)\|S_j p^*) - f(N_j(t)))_+ \end{aligned}$$

and therefore

$$\bigcup_j \{D(\hat{p}_j(t)\|S_j p^*) \geq f(N_j(t))\}.$$

Note that  $\{J_i(t), \mathcal{D}^c(t), N_j(t) = n\}$  occurs for at most twice because all actions are put into the list if  $\mathcal{D}^c(t)$  occurred. Thus, we have

$$\begin{aligned} &\mathbb{E} \left[ \sum_n \mathbb{1} [J_i(t), \mathcal{D}^c(t)] \right] \\ &\leq 2 \mathbb{E} \left[ \sum_j \sum_n \mathbb{1} [D(\hat{p}_{j,n}\|S_j p^*) \geq f(n)] \right] \\ &\leq 2 \sum_j \sum_n (n+1)^A e^{-nf(n)} \quad (\text{by (22)}) \\ &= O(1). \end{aligned}$$

□

*Proof of Lemma 13.* First, we have

$$\begin{aligned}
\mathcal{E}^c(t) &= \left\{ \max_i f(N_i(t)) > \min \{2\delta, (\rho_{1,L}\alpha(t))^2/4\} \right. \\
&\quad \left. \cup \min_i N_i(t) < \max\{c\sqrt{\log t}, (\log \log T)^{1/3}\} \cup 2\nu_{1,L}\alpha(t) > \|p^* - \mathcal{C}_1^c\|_M \right\} \\
&\subset \left\{ f(\min_i N_i(t)) > \min \{2\delta, (\rho_{1,L}\alpha(t))^2/4\} \right. \\
&\quad \left. \cup \min_i N_i(t) < c\sqrt{\log t} \cup c\sqrt{\log t} < (\log \log T)^{1/3} \cup 2\nu_{1,L}\alpha(t) > \|p^* - \mathcal{C}_1^c\|_M \right\} \\
&\subset \left\{ \frac{b}{\sqrt{c\sqrt{\log t}}} > \min \{2\delta, (\rho_{1,L}\alpha(t))^2/4\} \cup \min_i N_i(t) < c\sqrt{\log t} \right. \\
&\quad \left. \cup t < e^{\frac{(\log \log T)^{2/3}}{c}} \cup 2a/\log \log t > \|p^* - \mathcal{C}_1^c\|_M/\rho_{1,L} \right\} \\
&= \left\{ t < e^{\frac{b^4}{16\delta^4 c^2}} \cup \frac{(\log t)^{1/4}}{\log \log t} < \frac{b}{a\sqrt{c}\rho_{1,L}} \cup \min_i N_i(t) < c\sqrt{\log t} \right. \\
&\quad \left. \cup t < e^{\frac{(\log \log T)^{2/3}}{c}} \cup t < e^{2a\rho_{1,L}/\|p^* - \mathcal{C}_1^c\|_M} \right\}.
\end{aligned}$$

From  $\lim_{t \rightarrow \infty} (\log t)^{1/4} / \log \log t = \infty$  and  $e^{\frac{(\log \log T)^{2/3}}{c}} = o(e^{\log \log T}) = o(\log T)$  we have

$$\begin{aligned}
&\sum_{t=1}^T \mathbb{1}[J_i(t), \mathcal{E}^c(t)] \\
&= \sum_j \sum_{t=1}^T \mathbb{1}[J_i(t), N_j(t) < c\sqrt{\log t}] + o(\log T).
\end{aligned}$$

By (8), event  $\{J_i(t), N_j(t) < c\sqrt{\log t}, N_j(t) = n\}$  occurs for at most twice and therefore

$$\begin{aligned}
&\sum_{t=1}^T \mathbb{1}[J_i(t), \mathcal{E}^c(t)] \\
&\leq 2 \sum_j \sum_{n=1}^T \mathbb{1}\left[\bigcup_{t=1}^T \{n < c\sqrt{\log t}\}\right] + o(\log T) \\
&= o(\log T).
\end{aligned}$$

□