
Analysis and Design of Thompson Sampling for Stochastic Partial Monitoring

Taira Tsuchiya
The University of Tokyo
RIKEN AIP
tsuchiya@ms.k.u-tokyo.ac.jp

Junya Honda
The University of Tokyo
RIKEN AIP
honda@edu.k.u-tokyo.ac.jp

Masashi Sugiyama
RIKEN AIP
The University of Tokyo
sugi@k.u-tokyo.ac.jp

Abstract

We investigate finite stochastic partial monitoring, which is a general model for sequential learning with limited feedback. While Thompson sampling is one of the most promising algorithms on a variety of online decision-making problems, its properties for stochastic partial monitoring have not been theoretically investigated, and the existing algorithm relies on a heuristic approximation of the posterior distribution. To mitigate these problems, we present a novel Thompson-sampling-based algorithm, which enables us to exactly sample the target parameter from the posterior distribution. Besides, we prove that the new algorithm achieves the logarithmic *problem-dependent expected pseudo-regret* $O(\log T)$ for a linearized variant of the problem with local observability. This result is the first regret bound of Thompson sampling for partial monitoring, which also becomes the first logarithmic regret bound of Thompson sampling for linear bandits.

1 Introduction

Partial monitoring (PM) is a general sequential decision-making problem with limited feedback (Rus-tichini, 1999; Piccolboni and Schindelhauer, 2001). PM is attracting broad interest because it includes a wide range of problems such as the multi-armed bandit problem (Lai and Robbins, 1985), a linear optimization problem with full or bandit feedback (Zinkevich, 2003; Dani et al., 2008), dynamic pricing (Kleinberg and Leighton, 2003), and label efficient prediction (Cesa-Bianchi et al., 2005).

A PM game can be seen as a sequential game that is played by two players: a learner and an opponent. At every round, the learner chooses an action, while the opponent chooses an outcome. Then, the learner suffers an unobserved loss and receives a feedback symbol, both of which are determined from the selected action and outcome. The main characteristic of this game is that the learner cannot directly observe the outcome and loss. The goal of the learner is to minimize his/her cumulative loss over all rounds. The performance of the learner is evaluated by the regret, which is defined as the difference between the cumulative losses of the learner and the optimal action (*i.e.*, the action whose expected loss is the smallest).

There are mainly two types of PM games, which are the *stochastic* and *adversarial* settings (Piccolboni and Schindelhauer, 2001; Bartók et al., 2011). In the stochastic setting, the outcome at each round is determined from the *opponent's strategy*, which is a probability vector over the opponent's possible choices. On the other hand, in the adversarial setting, the outcomes are arbitrarily decided by the

opponent. We refer to the PM game with finite actions and finite outcomes as a *finite* PM game. In this paper, we focus on the stochastic finite game.

One of the first algorithms for PM was considered by Piccolboni and Schindelhauer (2001). They proposed the FeedExp3 algorithm, the key idea of which is to use an unbiased estimator of the losses. They showed that the FeedExp3 algorithm attains $\tilde{O}(T^{3=4})$ minimax regret for a certain class of PM games, and showed that any algorithm suffers linear minimax regret $\Omega(T)$ for the other class. Here T is the time horizon and the notation $\tilde{O}(\cdot)$ hides polylogarithmic factors. The upper bound $\tilde{O}(T^{3=4})$ is later improved by Cesa-Bianchi et al. (2006) to $O(T^{2=3})$, and they also provided a game with a matching lower bound.

In the seminal paper by Bartók et al. (2011), they classified PM games into four classes based on their minimax regrets. To be more specific, they classified games into trivial, easy, hard, and hopeless games, where their minimax regrets are 0 , $\Theta(\sqrt{T})$, $\Theta(T^{2=3})$, and $\Theta(T)$, respectively. Note that the easy game is also called a *locally observable* game. After their work, several algorithms have been proposed for the finite PM problem (Bartók et al., 2012; Vanchinathan et al., 2014; Komiyama et al., 2015). For the problem-dependent regret analysis, Komiyama et al. (2015) proposed an algorithm that achieves $O(\log T)$ regret with the optimal constant factor. However, it requires to solve a time-consuming optimization problem with infinitely many constraints at each round. In addition, this algorithm relies on the forced exploration to achieve the optimality, which makes the empirical performance near-optimal only after prohibitively many rounds, say, 10^5 or 10^6 .

Thompson sampling (TS, Thompson, 1933) is one of the most promising algorithms on a variety of online decision-making problems such as the multi-armed bandit (Lai and Robbins, 1985) and the linear bandit (Agrawal and Goyal, 2013b), and the effectiveness of TS has been investigated both empirically (Chapelle and Li, 2011) and theoretically (Kaufmann et al., 2012; Agrawal and Goyal, 2013a; Honda and Takemura, 2014). In the literature on PM, Vanchinathan et al. (2014) proposed a TS-based algorithm called BPM-TS (Bayes-update for PM based on TS) for stochastic PM, which empirically achieved state-of-the-art performance. Their algorithm uses Gaussian approximation to handle the complicated posterior distribution of the opponent’s strategy. However, this approximation is somewhat heuristic and can degrade the empirical performance due to the discrepancy from the exact posterior distribution. Furthermore, no theoretical guarantee is provided for BPM-TS.

Our goals are to establish a new TS-based algorithm for stochastic PM, which allows us to sample the opponent’s strategy parameter from the exact posterior distribution, and investigate whether the TS-based algorithm can achieve sub-linear regret in stochastic PM. Using the accept-reject sampling, we propose a new TS-based algorithm for PM (TSPM), which is equipped with a numerical scheme to obtain a posterior sample from the complicated posterior distribution. We derive a logarithmic regret upper bound $O(\log T)$ for the proposed algorithm on the locally observable game under a linearized variant of the problem. This is the first regret bound for TS on the locally observable game. Moreover, our setting includes the linear bandit problem and our result is also the first logarithmic expected regret bound of TS for the linear bandit, whereas a high-probability bound was provided, for example, in Agrawal and Goyal (2013b). Finally, we compare the performance of TSPM with existing algorithms in numerical experiments, and show that TSPM outperforms existing algorithms.

2 Preliminaries

This paper studies finite stochastic PM games (Bartók et al., 2011). A PM game with N actions and M outcomes is defined by a pair of a loss matrix $\mathbf{L} = (\ell_{i,j}) \in \mathbb{R}^{N \times M}$ and feedback matrix $\mathbf{H} = (h_{i,j}) \in [A]^{N \times M}$, where A is the number of feedback symbols and $[A] = \{1, 2, \dots, A\}$.

A PM game can be seen as a sequential game that is played by two players: the learner and the opponent. At each round $t = 1, 2, \dots, T$, the learner selects action $i(t) \in [N]$, and at the same time the opponent selects an outcome based on the opponent’s strategy $p \in \mathcal{P}_M$, where $\mathcal{P}_n = \{p \in \mathbb{R}^n : p_k \geq 0, \sum_{k=1}^n p_k = 1\}$ is the $(n - 1)$ -dimensional probability simplex. The outcome $j(t)$ of each round is an independent and identically distributed sample from p , and then, the learner suffers loss $\ell_{i(t),j(t)}$ at time t . The learner cannot directly observe the value of this loss, but instead observes the *feedback symbol* $y(t) = h_{i(t),j(t)} \in [A]$. The setting explained above has been widely studied in the literature of stochastic PM (Bartók et al., 2011; Komiyama et al., 2015), and we call this the *discrete* setting. In Section 4, we also introduce a *linear* setting for theoretical analysis, which is a slightly different setting from the discrete one.

The learner aims to minimize the cumulative loss over T rounds. The expected loss of action i is given by $L_i^\top p$, where L_i is the i -th column of L^\top . We say action i is *optimal* under strategy p if $(L_i - L_j)^\top p \geq 0$ for any $j \neq i$. We assume that the optimal action is unique, and without loss of generality that the optimal action is action 1. Let $\Delta_i = (L_i - L_1)^\top p \geq 0$ for $i \in [N]$ and $N_i(t)$ be the number of times action i is selected before the t -th round. When the time step is clear from the context, we use n_i instead of $N_i(t)$. We adopt the pseudo-regret to measure the performance: $\text{Reg}(T) = \sum_{t=1}^T \Delta_{I(t)} = \sum_{i \in [N]} \Delta_i N_i(T+1)$. This is the relative performance of the algorithm against the *oracle*, which knows the optimal action 1 before the game starts.

We introduce the following definitions to clarify the class of PM games, for which we develop an algorithm and derive a regret upper bound. The following cell decomposition is the concept to divide the simplex \mathcal{P}_M based on the loss matrix to identify the optimal action, which depends on the opponent's strategy p .

Definition 1 (Cell decomposition and Pareto-optimality (Bartók et al., 2011)). For every action $i \in [N]$, cell $C_i := \{p \in \mathcal{P}_M : (L_i - L_j)^\top p \geq 0, \forall j \neq i\}$ is the set of opponent's strategies for which action i is optimal. Action i is *Pareto-optimal* if there exists an opponent's strategy p under which action i is optimal.

Each cell is a convex closed polytope. Next, we define *neighbors* between two Pareto-optimal actions, which intuitively means that the two actions "touch" each other in their surfaces.

Definition 2 (Neighbors and neighborhood action (Bartók et al., 2011)). Two Pareto-optimal actions i and j are *neighbors* if $C_i \cap C_j$ is an $(M-2)$ -dimensional polytope. For two neighboring actions $i, j \in [N]$, the *neighborhood action set* is defined as $N_{i,j}^+ = \{k \in [N] : C_i \cap C_j \cap C_k \neq \emptyset\}$.

Note that the neighborhood action set $N_{i,j}^+$ includes actions i and j from its definition. Next, we define the *signal matrix*, which encodes the information of the feedback matrix \mathbf{H} so that we can utilize the feedback information.

Definition 3 (Signal matrix (Komiyama et al., 2015)). The signal matrix $S_i \in \mathbb{R}^{M \times M}$ of action i is defined as $(S_i)_{y,j} = \mathbb{P}[h_{i,j} = y]$, where $\mathbb{P}[X] = 1$ if the event X is true and 0 otherwise.

Note that if we define the signal matrix as above, $S_i p \in \mathbb{R}^M$ is a probability vector over feedback symbols of action i . The following *local observability* condition separates easy and hard games, this condition intuitively means that the information obtained by taking actions in the neighborhood action set $N_{i,j}^+$ is sufficient to distinguish the loss difference between actions i and j .

Definition 4 (Local observability (Bartók et al., 2011)). A partial monitoring game is said to be *locally observable* if for all pairs i, j of neighboring actions, $L_i - L_j \in \text{Im} S_k^\top$, where $\text{Im} V$ is the image of the linear map V , and $V \oplus W$ is the direct sum between the vector spaces V and W .

We also consider the concept of the *strong local observability* condition, which implies the above local observability condition.

Definition 5 (Strong local observability). A partial monitoring game is said to be *strongly locally observable* if for all pairs $i, j \in [N]$, $L_i - L_j \in \text{Im} S_i^\top \oplus \text{Im} S_j^\top$.

This condition was assumed in the theoretical analysis in Vanchinathan et al. (2014), and we also assume this condition in theoretical analysis in Section 4. Note that the strong local observability means that, for any $j \neq k$, there exists $z_{j,k} \in \mathbb{R}^{2A}$ such that $L_j - L_k = (S_j^\top, S_k^\top) z_{j,k}$.

Notation. Let $\|\cdot\|$ and $\|\cdot\|_p$ be the Euclidian norm and p -norm, and let $\|x\|_A = \sqrt{x^\top A x}$ be the norm induced by the positive semidefinite matrix $A \succeq 0$. Let $D_{\text{KL}}(p||q) = \sum_{a=1}^A p_a \log(p_a/q_a)$ be the Kullback-Leibler divergence of p from q . The vector $e_y \in \mathbb{R}^M$ is the y -th orthonormal basis of \mathbb{R}^M , and $\mathbf{1}_n = [1, \dots, 1]^\top$ is the n -dimensional all-one vector. Let $q_i^{(t)}$ be the empirical feedback distribution of action i at time t , i.e., $q_i^{(t)} = [n_{i1}/n_i, \dots, n_{iA}/n_i]^\top \in \mathcal{P}_A$, where $n_{iy} = \sum_{s=1}^t \mathbb{1}[i(s) = i, y(s) = y]$ and $n_i = \sum_{y=1}^A n_{iy}$. The notation is summarized in Appendix A.

Methods for Sampling from Posterior Distribution. We briefly review the methods to draw a sample from the posterior distribution. While TS is one of the most promising algorithms, the posterior distribution can be in a quite complicated form, which makes obtaining a sample from it

Algorithm 1: TSPM Algorithm

Input: prior parameter $\lambda > 0$

- 1 Set $B_0 = \lambda I_M, b_0 = 0$.
 - 2 Take each action for $n = 1$ times.
 - 3 **for** $t = 1, 2, \dots, T$ **do**
 - 4 Sample $\tilde{p}_t \sim \pi(p \mid \tilde{r}(s), y(s)g_{s=1}^t)$ based on the accept-reject sampling (Algorithm 2).
 - 5 Take action $i(t) = \arg \max_{i \in [N]} L_i(\tilde{p}_t)$ and observe feedback $y(t)$.
 - 6 Update $B_t = B_{t-1} + S_{i(t)}^>, b_t = b_{t-1} + S_{i(t)}^> e_{y(t)}$.
-

Algorithm 2: Accept-Reject Sampling

Input: constant $R \geq [0, 1]$

- 1 **while** true **do**
 - 2 Sample $\tilde{p}_t \sim g_t(p)$ (Algorithm 3).
 - 3 Sample $\tilde{u} \sim \mathcal{U}([0, 1])$.
 - 4 **if** $R\tilde{u} < F_t(\tilde{p}_t)/G_t(\tilde{p}_t)$ **then**
 - 5 **return** \tilde{p}_t .
-

Algorithm 3: Sampling from $g_t(p)$

- 1 Compute \tilde{B}_t, \tilde{b}_t from B_t, b_t .
 - 2 **repeat**
 - 3 | Sample $p^{(\cdot)} \sim \mathcal{N}(\tilde{B}_t^{-1}\tilde{b}_t, \tilde{B}_t^{-1})$.
 - 4 **until** $p^{(\cdot)} \geq P_{M-1}$;
 - 5 **return** $\tilde{p} = [p^{(\cdot) >}, 1 - \sum_{i=1}^{M-1} p^{(\cdot)}_i]^>$.
-

computationally hard. To overcome this issue, a variety of approximate posterior sampling methods have been considered, such as Gibbs sampling, Langevin Monte Carlo, Laplace approximation, and the bootstrap (Russo et al., 2018, Section 5). Recent work (Lu and Van Roy, 2017) proposed a flexible approximation method, which can even efficiently be applied to quite complex models such as neural networks. However, more recent work revealed that algorithms based on such an approximation procedure *can* suffer a linear regret (Phan et al., 2019), even if the approximation error in terms of the α -divergence is small enough.

Although BPM-TS is one of the best methods for stochastic PM, it approximates the posterior by a Gaussian distribution in a heuristic way, which can degrade the empirical performance due to the distributional discrepancy from the exact posterior distribution. Furthermore, no theoretical guarantee is provided for BPM-TS. In this paper, we mitigate these problems by providing a new algorithm for stochastic PM, which allows us to exactly draw samples from the posterior distribution. We also give theoretical analysis for the proposed algorithm.

3 Thompson-sampling-based Algorithm for Partial Monitoring

In this section, we present a new algorithm for stochastic PM games, where we name the algorithm TSPM (TS-based algorithm for PM). The algorithm is given in Algorithm 1, and we will explain the subroutines in the following.

3.1 Accept-Reject Sampling

We adopt the accept-reject sampling (Casella et al., 2004) to *exactly* draw samples from the posterior distribution. The accept-reject sampling is a technique to draw samples from a specific distribution f , and a key feature is to use a *proposal distribution* g , from which we can easily draw a sample and whose ratio to f , that is f/g , is bounded by a constant value R . To obtain samples from f , (i) we generate samples $X \sim g$; (ii) accept X with probability $f(X)/Rg(X)$. Note that f and g do not have to be normalized when the acceptance probability is calculated.

Let $\pi(p)$ be a prior distribution for p . Then an unnormalized density of the posterior distribution for p can be expressed as

$$F_t(p) = \pi(p) \prod_{i=1}^W \exp(-n_i D_{\text{KL}}(q_i^{(t)} \parallel S_i p)), \quad (1)$$

the detailed derivation of which is given in Appendix B. We use the proposal distribution with unnormalized density

$$G_t(p) = \pi(p) \prod_{i=1}^{\mathcal{N}} \exp\left(-\frac{1}{2}n_i k q_i^{(t)} S_i p^2\right). \quad (2)$$

Based on these distributions, we use Algorithm 2 for exact sampling from the posterior distribution, where $U([0, 1])$ is the uniform distribution over $[0, 1]$ and $g_t(p)$ is the distribution corresponding to the unnormalized density $G_t(p)$ in (2). The following proposition shows that setting $R = 1$ realizes the exact sampling.

Proposition 1. *Let $f_t(p)$ be the distribution corresponding to the unnormalized density $F_t(p)$ in (1). Then, the output of Algorithm 2 with $R = 1$ follows $f_t(p)$.*

This proposition can easily be proved by Pinsker’s inequality, which is detailed in Appendix B.

In practice, $R \in [0, 1]$ is a parameter to balance the amount of over-exploration and the computational efficiency. As R decreases from 1, the algorithm tends to accept a point p far from the mode. The case $R = 0$ corresponds the TSPM algorithm where the proposal distribution is used without the accept-reject sampling, which we call *TSPM-Gaussian*. As we will see in Section 4, TSPM-Gaussian corresponds to exact sampling of the posterior distribution when the feedback follows a Gaussian distribution rather than a multinomial distribution.

TSPM-Gaussian can be related to BPM-TS (Vanchinathan et al., 2014) in the sense that both of them use samples from Gaussian distributions. Nevertheless, they use different Gaussians and TSPM-Gaussian performs much better than BPM-TS as we will see in the experiments. Details on the relation between TSPM-Gaussian and BPM-TS are described in Appendix D.

In general, we can realize efficient sampling with a small number of rejections if the proposal distribution and the target distribution are close to each other. On the other hand, in our problem, the densities in (1) and (2) for each fixed point p exponentially decay with the number of samples n_i if the empirical feedback distribution $q_i^{(t)}$ converges. This means that $F_t(p)$ and $G_t(p)$ have an exponentially large relative gap in most rounds. Nevertheless, the number of rejections does not increase with t as we will see in the experiments, which suggests that the proposal distribution approximates the target distribution well with high probability.

3.2 Sampling from Proposal Distribution

When we consider Gaussian density $N(0, \lambda I_M)$ truncated over P_M as a prior, the proposal distribution also has the Gaussian density $N(B_t^{-1} b_t, B_t^{-1})$ over P_M , where

$$B_t = \lambda I_M + \sum_{i=1}^{\mathcal{N}} n_i S_i^{\succ} S_i = B_{t-1} + S_{i(t)}^{\succ} S_{i(t)}, \quad b_t = \sum_{i=1}^{\mathcal{N}} n_i S_i^{\succ} q_i^{(t)} = b_{t-1} + S_{i(t)}^{\succ} e_{y(t)}. \quad (3)$$

Here note that the probability simplex P_M is in an $(M - 1)$ -dimensional space and a sample from $N(0, \lambda I_M)$ is not contained in P_M with probability one. In the literature, e.g., Altmann et al. (2014), sampling methods for Gaussian distributions truncated on a simplex have been discussed. We use one of these procedures summarized in Algorithm 3, where we first sample $M - 1$ elements of p from another Gaussian distribution and determine the remaining element by the constraint $\sum_{i=1}^M p_i = 1$.

Proposition 2. *Sampling from $g_t(p)$ is equivalent to Algorithm 3 with*

$$\tilde{B}_t = C_t - 2D_t + f_t \mathbf{1}_M - \mathbf{1}_M^{\succ} \mathbf{1}_M^{-1}, \quad \tilde{b}_t = f_t \mathbf{1}_M - d_t + b_t^{(\cdot)} - b^{(M)} \mathbf{1}_M^{-1},$$

where $B_t = \begin{bmatrix} C_t & d_t \\ d_t^{\succ} & f_t \end{bmatrix}$ for $C_t \in \mathbb{R}^{(M-1) \times (M-1)}$, $d_t \in \mathbb{R}^{M-1}$, $f_t \in \mathbb{R}$, $b_t = [b_t^{(\cdot)\succ}, b_t^{(M)}] \in \mathbb{R}^{M-1}$, and $D_t = \frac{1}{2}(d_t \mathbf{1}_M^{-1} + \mathbf{1}_M^{-1} d_t^{\succ})$.

We give the proof of this proposition for self-containedness in Appendix C.

4 Theoretical Analysis

This section considers a regret upper bound of the TSPM algorithm.

In the theoretical analysis, we consider a *linear* setting of PM. In the linear PM, the learner suffers the expected loss $L_{i(t)}^\top p$ as in the discrete setting, and receives feedback vector $y(t) = S_i p + \epsilon_t$ for $\epsilon_t \sim \mathcal{N}(0, I_M)$ whereas the one-hot representation of $y(t)$ is distributed by the probability vector $S_i p$ in the discrete setting. Therefore, if ϵ_t can be regarded as a sub-Gaussian random variable as in [Kirschner et al. \(2020\)](#) then the linear PM includes the discrete PM, though our theoretical analysis requires ϵ_t to be Gaussian. The relation between discrete and linear settings can also be seen from the observation that bandit problems with Bernoulli and Gaussian rewards can be expressed as discrete and linear PM, respectively. The linear PM also includes the linear bandit problem, where the feedback vector is expressed as $L_i^\top p + \epsilon_t$.

In the linear PM, $G_t(p)$ in (2) becomes the exact posterior distribution rather than a proposal distribution. The definition of the cell decomposition for this setting is largely the same as that of discrete setting and detailed in Appendix F. Therefore, TS with exact posterior sampling in the linear PM corresponds to TSPM-Gaussian. In the linear PM, the unknown parameter p is in \mathbb{R}^M rather than in \mathbb{P}_M , and therefore we consider the prior $\pi(p) = \mathcal{N}(0, \lambda I_M)$ over \mathbb{R}^M , where the posterior distribution becomes $\mathcal{N}(B_t^{-1} b_t, B_t^{-1})$.

There are a few works that analyze TS for the PM because of its difficulty. For example in [Vanchinathan et al. \(2014\)](#), an analysis of the TS-based algorithm (BPM-TS) is not given despite the fact that its performance is better than the algorithm based on a confidence ellipsoid (BPM-LEAST). [Zimmert and Lattimore \(2019\)](#) considered the theoretical aspect of a variant of TS for the linear PM in view of the Bayes regret, but this algorithm is based on the knowledge on the time horizon and different from the family of TS used in practice. More specifically, their algorithm considers the posterior distribution for *regret* (not pseudo-regret), and an action is chosen according to the posterior probability that each arm minimizes the *cumulative* regret. Thus, the time horizon also needs to be known.

Types of Regret Bounds. We focus on the (a) *problem-dependent* (b) *expected pseudo-regret*. (a) In the literature, a *minimax* (or *problem-independent*) regret bound has mainly been considered, for example, to classify difficulties of the PM problem ([Bartók et al., 2010](#); [Bartók et al., 2011](#)). On the other hand, a *problem-dependent* regret bound often reflects the empirical performance more clearly than the minimax regret ([Bartók et al., 2012](#); [Vanchinathan et al., 2014](#); [Komiya et al., 2015](#)). For this reason, we consider this problem-dependent regret bound. (b) In complicated settings of bandit problems, a *high-probability regret bound* has mainly been considered ([Abbasi-Yadkori et al., 2011](#); [Agrawal and Goyal, 2013b](#)), which bounds the pseudo-regret with high probability $1 - \delta$. Though such a bound can be transformed to an expected regret bound, this type of analysis often sacrifices the tightness since a linear regret might be suffered with small probability δ . This is why the analysis in [Vanchinathan et al. \(2014\)](#) for BPM-LEAST finally yielded an $\tilde{O}(\sqrt{T})$ expected regret bound whereas their high-probability bound is $O(\log T)$.

4.1 Regret Upper Bound

In the following theorem, we show that logarithmic problem-dependent expected regret is achievable by the TSPM-Gaussian algorithm.

Theorem 3 (Regret upper bound). *Consider any finite stochastic linear partial monitoring game. Assume that the game is strongly locally observable and $\Delta_i = (L_i - L_1)^\top p > 0$ for any $i \neq 1$. Then, the regret of TSPM-Gaussian satisfies for sufficiently large T that*

$$\mathbb{E} [\text{Reg}(T)] = O \left(\frac{AN^2 M \max_{i \in [N]} \Delta_i}{\Lambda^2} \log T \right), \quad (4)$$

where $\Lambda := \min_{i \neq 1} \Lambda_i$ for $\Lambda_i = \Delta_i / k z_{1:i} k$ with $z_{1:i}$ defined after Definition 5.

Remark. In the proof of Theorem 3, it is sufficient to assume that $L_1 - L_i \succeq \text{Im } S_1^\top - \text{Im } S_i^\top$ for $i \in [N]$, which is weaker than the strong local observability, though it is still sometimes stronger than the local observability condition.

The proof of Theorem 3 is given in Appendix F. This result is the first problem-dependent bound of TS for PM, which also becomes the first logarithmic regret bound of TS for linear bandits.

The norm of $z_{j:k}$ in Λ intuitively indicates the difficulty of the problem. Whereas we can estimate $(S_j p, S_k p)$ with noise through taking actions j and k , the actual interest is the gap of the losses $p^\top (L_j - L_k) = (S_j p, S_k p)^\top z_{j:k}$. Thus, if $k z_{j:k} k$ is large, the gap estimation becomes difficult since the noise is enhanced through $z_{j:k}$.

Unfortunately, the derived bound in Theorem 3 has quadratic dependence on N , which seems to be not tight. This quadratic dependence comes from the difficulty of the *expected* regret analysis. In general, we evaluate the regret before and after the convergence of the statistics separately. Whereas the latter one usually becomes dominant, the main difficulty comes from the analysis of the former one, which might become large with low probability (Agrawal and Goyal, 2012; Kaufmann et al., 2012; Agrawal and Goyal, 2013a).

In our analysis, we were not able to bound the former one within a non-dominant order, though it is still logarithmic in T . In fact, our analysis shows that the regret after convergence is $O(\sum_{i \neq 1} \Delta_i \frac{\Lambda}{2} \log T)$ as shown in Lemma 18 in Appendix F, which will become the regret with high probability. In particular, if we consider the classic bandit problem as a PM game, we can confirm that the derived bound after convergence becomes the best possible bound

$$O \sum_{i \neq 1} \frac{\log T}{\Delta_i}$$

by considering Λ_i depending on each suboptimal arm i as the difficulty measure instead of Λ . Still, deriving a regret bound for the term before convergence within an non-dominant order is an important future work.

4.2 Technical Difficulties of the Analysis

The main difficulty of this regret analysis is that PM requires to consider the statistics of *all* actions when the number of selections $N_i(t)$ of some action i is evaluated. This is in stark contrast to the analysis of the classic bandit problems, where it becomes sufficient to evaluate statistics of action i and the best action 1. This makes the analysis remarkably complicated in TS, where we need to separately consider the randomness caused by the feedback and TS.

To overcome this difficulty, we handle the effect of actions of no interest in two different novel ways depending on each decomposed regret. The first one is to evaluate the worst-case effect of these actions based on an argument (Lemma 10) related to the law of the iterated logarithm (LIL), which is sometimes used in the best-arm identification literature to improve the performance (Jamieson et al., 2014). The second one is to bound the action-selection probability of TS using an argument of (super-)martingale (Theorem 16), which is of independent interest. Whereas such a technique is often used for the construction of confidence bounds (Abbasi-Yadkori et al., 2011), we reveal that it is also useful for evaluation of the regret of TS.

We only focused on the Gaussian noise $\epsilon_t \sim N(0, I_M)$, rather than the more general sub-Gaussian noise. This restriction to the Gaussian noise comes from the essential difficulty of the problem-dependent analysis of TS, where lower bounds for some probabilities are needed whereas the sub-Gaussian assumption is suited for obtaining upper bounds. To the best of our knowledge, the problem-dependent regret analysis for TS on the sub-Gaussian case has never been investigated even for the multi-armed bandit setting, which is quite simple compared to that of PM. In the literature of the problem-dependent regret analysis, the noise distribution is restricted to distributions with explicitly given forms, e.g., Bernoulli, Gaussian, or more generally a one-dimensional canonical exponential family (Kaufmann et al., 2012; Agrawal and Goyal, 2013a; Korda et al., 2013). Their analysis relies on the specific characteristic of the distribution to bound the problem-dependent regret.

5 Experiments

In this section, we numerically compare the performance of TSPM and TSPM-Gaussian against existing methods, which are RandomPM (the algorithm which selects action randomly), FeedExp3 (Piccolboni and Schindelhauer, 2001), and BPM-TS (Vanchinathan et al., 2014). Recently, Lattimore and Szepesvári (2019) considered the sampling-based algorithm called Mario sampling for easy games. Mario sampling coincides with TS (except for the difference between pseudo-regret and regret with known time horizon) mentioned in the last section when any pair of actions is a neighbor. As shown in Appendix G, this property is indeed satisfied for dp-easy games defined in the following. Therefore, the performance is essentially the same between TSPM with $R = 1$ and Mario sampling. To compare the performance, we consider a dynamic pricing problem, which is a typical example of PM games. We conducted experiments on the discrete setting because the experiments for PM has been mainly focused on the discrete setting.

(a) dp-easy, N = M = 3

(b) dp-easy, N = M = 5

(c) dp-easy, N = M = 7

(d) dp-hard, N = M = 3

(e) dp-hard, N = M = 5

(f) dp-hard, N = M = 7

Figure 1: Regret-round plots of algorithms. The solid lines indicate the average over independent trials. The thin fillings are the standard error.

In the dynamic pricing game, the player corresponds to a seller, and the opponent corresponds to a buyer. At each round, the seller sells an item for a specific price $p(t)$ and the buyer comes with an evaluation price $j(t)$ for the item, where the selling price and the evaluation price correspond to the action and outcome, respectively. The buyer buys the item if the selling price is smaller than or equal to $j(t)$ and not otherwise. The seller can only know if the buyer bought the item (denoted as feedback 1) or did not buy the item (denoted as 0). The seller aims to minimize the cumulative loss, and there are two types of definitions for the loss, where each induced game falls into the easy and hard games. We call them dp-easy and dp-hard games, respectively.

In both cases, the seller incurs the constant loss c when the item is not bought due to the loss of opportunity to sell the item. In contrast, when the item is bought, the loss incurred to the seller is different between these settings. The seller in the dp-easy game does not take the buyer's evaluation price into account. In other words, the seller gains the selling price as a reward (equivalently incurs $-i(t)$ as a loss). Therefore, the loss for the selling price and the evaluation price is

$$l_{i(t);j(t)} = -i(t)1[i(t) \leq j(t)] + c1[i(t) > j(t)] :$$

This setting can be regarded as a generalized version of the online posted price mechanism, which was addressed in e.g., [Blum et al. \(2004\)](#) and [Cesa-Bianchi et al. \(2006\)](#), and an example of strongly locally observable games.

On the other hand, the seller in dp-hard game takes the buyer's evaluation price into account when the item is bought. In other words, the seller incurs the difference between the opponent evaluation and the selling price $j(t) - i(t)$ as a loss because the seller could have made more profit if the seller had sold at the price $j(t)$. Therefore, the loss incurred at time t is

$$l_{i(t);j(t)} = (j(t) - i(t))1[i(t) \leq j(t)] + c1[i(t) > j(t)] :$$

This setting is also addressed in [Cesa-Bianchi et al. \(2006\)](#), and belongs to the class of hard games. Note that our algorithm can also be applied to a hard game, though there is no theoretical guarantee.

Setup. In the both dp-easy and dp-hard games, we set $N = M = 2^f - 3; 5; 7$ and $c = 2$. We fixed the time horizon T to 10000 and simulated 100 times. For FeedExp3 and BPM-TS, the setup of hyperparameters follows their original papers. For TSPM, we set $\epsilon = 0.001$, and R was selected from $\{0.01; 1; 10\}$. Here, recall that TSPM with $R = 1$ and $R = 0$ correspond to the exact sampling and TSPM-Gaussian, respectively, and a smaller value of R gives the higher acceptance probability in the accept-reject sampling. Therefore, using smaller R makes the algorithm time-efficient, although it can worsen the performance since it over-explores the tail of the posterior distributions. To stabilize

(a) dp-easyN = M = 3 (b) dp-easyN = M = 5 (c) dp-easyN = M = 7

Figure 2: The number of rejected times by the accept-reject sampling. The solid lines indicate the average over 100 independent trials after taking moving average with window size 10.

sampling from the proposal distribution in Algorithm 3, we used an initialization that takes each action $n = 10A$ times. The detailed settings of the experiments with more results are given in Appendix H.

Results. Figure 1 is the empirical comparison of the proposed algorithms against the benchmark methods. This result shows that, in all cases, the TSPM with exact sampling gives the best performance. TSPM-Gaussian also outperforms BPM-TS even though both of them use Gaussian distributions as posteriors. Besides, the experimental results suggest that our algorithm performs reasonably well even for a hard game. It can be observed that the proposed methods outperform BPM-TS more significantly for a larger number of outcomes. Further discussion for this observation is given in Appendix D.

Figure 2 shows the number of rejections at each time step in the accept-reject sampling. We counted the number of times that either Line 4 in Algorithm 2 or Line 4 in Algorithm 3 was not satisfied. In the accept-reject sampling, it is desirable that the frequency of rejection does not increase as the time-step and does not increase rapidly with the number of outcomes. We can see that the former one is indeed satisfied. For the latter property, the frequency of rejection becomes unfortunately large when exact sampling ($R = 1$) is conducted. Still, we can substantially improve this frequency by setting R to be a small value or zero, which still keeps regret tremendously better than that of BPM with almost the same time efficiency as BPM-TS.

6 Conclusion and Discussion

This paper investigated Thompson sampling (TS) for stochastic partial monitoring from the algorithmic and theoretical viewpoints. We provided a new algorithm that enables exact sampling from the posterior distribution, and numerically showed that the proposed algorithm outperforms existing methods. Besides, we provided an upper bound for the problem-dependent logarithmic expected pseudo-regret for the linearized version of the partial monitoring. To our knowledge, this bound is the first logarithmic problem-dependent expected pseudo-regret bound of a TS-based algorithm for linear bandit problems and strongly locally observable partial monitoring games.

There are several remaining questions. As mentioned in Section 4, Kirschner et al. (2020) considered linear partial monitoring with the feedback structure $y_t = S_{i(t)}p + \xi_t$, where $(\xi_t)_{t=1}^T$ is a sequence of independent sub-Gaussian noise vectors. This setting is the generalization of our linear setting, where $(\xi_t)_{t=1}^T$ are i.i.d. Gaussian vectors. Therefore, a natural question that arises is whether we can extend our analysis on TSPM-Gaussian to the sub-Gaussian case, although we believe it would be not straightforward as discussed in Section 4. It is also an important open problem to derive a regret bound on TSPM using the exact posterior sampling for the discrete partial monitoring. Although we conjecture that the algorithm also achieves logarithmic regret for the setting, there still remain some difficulties in the analysis. In particular, we have to handle the KL divergence $D(p \| q)$ and consider the restriction of the support of the opponent's strategy \mathcal{P}_M , which make the analysis much more complicated. Besides, it is worth noting that the theoretical analysis of TS for hard games has never been theoretically investigated. We believe that in general TS suffers linear regret in the minimax sense due to its greediness. However, we conjecture that TS can achieve the sub-linear regret for some specific instances of hard games in the sense of the problem-dependent regret, as empirically observed in the experiments. Finally, it is an important open problem to derive the minimax regret for anytime TS-based algorithms. This needs more detailed analysis in terms of the regret bound, which were dropped in our main result.

Broader Impact

Application. Partial monitoring (PM) includes various online decision-making problems such as multi-armed bandits, linear bandits, dynamic pricing, and label efficient prediction. Not only can PM handles them, the dueling bandits, combinatorial bandits, transductive bandits, and many other problems can be seen as a partial monitoring game, as discussed in [Kirschner et al. \(2020\)](#). Therefore, our analysis of Thompson sampling (TS) for PM games pushes the application of TS to a more wide range of online decision-making problems forward. Moreover, PM has the potential that novel online-decision making problems are newly discovered, where we have to handle the limited feedback in an online fashion.

Practical Use. The obvious advantage of using TS is that the users can easily apply the algorithm to their problems. They do not have to solve mathematical optimization problems, which are often required to solve when using non-sampling-based algorithms ([Bartók et al., 2012](#); [Komiya et al., 2015](#)). For the negative side, the theoretical analysis for the regret upper bound might make the users become overconfident when the users use their algorithms. For example, they might use the TSPM algorithm to the linear PM game with heavy-tailed noise, such as sub-exponential noise, without noticing it. Nevertheless, this is not an TS-specific problem, but one that can be found in many theoretical studies, and TS is still one of the most promising policies.

Acknowledgements

The authors would like to thank the meta-reviewer and reviewers for a lot of helpful comments. The authors would like to thank Kento Nozawa and Ikko Yamane for maintaining servers for our experiments, and Kenny Song for helpful discussion on the writing. TT was supported by Toyota-Dwango AI Scholarship, and RIKEN Junior Research Associate Program for the final part of the project. JH was supported by KAKENHI 18K17998, and MS was supported by KAKENHI 17H00757.

References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320, 2011.
- Shipra Agrawal and Navin Goyal. Analysis of Thompson sampling for the multi-armed bandit problem. In the 25th Annual Conference on Learning Theory, volume 23, pages 39.1–39.26, 2012.
- Shipra Agrawal and Navin Goyal. Further optimal regret bounds for Thompson sampling. In *Sixteenth International Conference on Artificial Intelligence and Statistics*, volume 31, pages 99–107, 2013a.
- Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear payoffs. In the 30th International Conference on Machine Learning, pages 127–135, 2013b.
- Yoann Altmann, Steve McLaughlin, and Nicolas Dobigeon. Sampling from a multivariate gaussian distribution truncated on a simplex: A review. In *2014 IEEE Workshop on Statistical Signal Processing (SSP)*, pages 113–116, 2014.
- Gábor Bartók, Dávid Pál, and Csaba Szepesvári. Toward a classification of finite partial-monitoring games. In *Algorithmic Learning Theory*, pages 224–238, 2010.
- Gábor Bartók, Navid Zolghadr, and Csaba Szepesvári. An adaptive algorithm for finite stochastic partial monitoring. In the 29th International Conference on Machine Learning, pages 1–20, 2012.
- Gábor Bartók, Dávid Pál, and Csaba Szepesvári. Minimax regret of finite partial-monitoring games in stochastic environments. In the 24th Annual Conference on Learning Theory, volume 19, pages 133–154, 2011.
- Avrim Blum, Vijay Kumar, Atri Rudra, and Felix Wu. Online learning in online auctions. *Theoretical Computer Science*, 324(2):137–146, 2004.
- George Casella, Christian P. Robert, and Martin T. Wells. Generalized accept-reject sampling schemes volume 45 of *Lecture Notes Monograph Series*, pages 342–347. Institute of Mathematical Statistics, 2004.
- Nicolò Cesa-Bianchi, Gábor Lugosi, and Gilles Stoltz. Minimizing regret with label efficient prediction. *IEEE Transactions on Information Theory*, 51(6):2152–2162, 2005.

- Nicolò Cesa-Bianchi, Gábor Lugosi, and Gilles Stoltz. Regret minimization under partial monitoring. *Mathematics of Operations Research* 31(3):562–580, 2006.
- Olivier Chapelle and Lihong Li. An empirical evaluation of Thompson sampling. *Advances in Neural Information Processing Systems* 24:2249–2257, 2011.
- Varsha Dani, Thomas P. Hayes, and Sham M. Kakade. Stochastic linear optimization under bandit feedback. In *21st Annual Conference on Learning Theory*, pages 355–366, 2008.
- Junya Honda and Akimichi Takemura. Optimality of Thompson sampling for Gaussian bandits depends on priors. In *the Seventeenth International Conference on Artificial Intelligence and Statistics*, volume 33, pages 375–383, 2014.
- Kevin Jamieson, Matthew Malloy, Robert Nowak, and Sébastien Bubeck. lin^2 : An optimal exploration algorithm for multi-armed bandits. In *the 27th Conference on Learning Theory*, pages 423–439, 2014.
- Emilie Kaufmann, Nathaniel Korda, and Rémi Munos. Thompson sampling: An asymptotically optimal finite-time analysis. In *Algorithmic Learning Theory*, pages 199–213, 2012.
- Johannes Kirschner, Tor Lattimore, and Andreas Krause. Information directed sampling for linear partial monitoring. *arXiv preprint arXiv:2002.11182*, 2020.
- Robert Kleinberg and Tom Leighton. The value of knowing a demand curve: Bounds on regret for online posted-price auctions. In *the 44th Annual IEEE Symposium on Foundations of Computer Science*, pages 594–605, 2003.
- Junpei Komiyama, Junya Honda, and Hiroshi Nakagawa. Regret lower bound and optimal algorithm in finite stochastic partial monitoring. In *Advances in Neural Information Processing Systems* 28, pages 1792–1800, 2015.
- Nathaniel Korda, Emilie Kaufmann, and Remi Munos. Thompson sampling for 1-dimensional exponential family bandits. In *Advances in Neural Information Processing Systems*, pages 1448–1456, 2013.
- T. L. Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics* 6(1):4–22, 1985.
- Tor Lattimore and Csaba Szepesvári. An information-theoretic approach to minimax regret in partial monitoring. In *the 32nd Annual Conference on Learning Theory*, volume 99, pages 2111–2139, 2019.
- Xiuyuan Lu and Benjamin Van Roy. Ensemble sampling. *Advances in Neural Information Processing Systems* 30:3258–3266, 2017.
- My Phan, Yasin Abbasi-Yadkori, and Justin Domke. Thompson sampling and approximate inference. In *Advances in Neural Information Processing Systems*, pages 8804–8813, 2019.
- Antonio Piccolboni and Christian Schindelhauer. Discrete prediction games with arbitrary feedback and loss. In *COLT/EuroCOLT*, pages 208–223, 2001.
- Daniel J. Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, and Zheng Wen. A tutorial on Thompson sampling. *Foundations and Trends in Machine Learning* 1(1):1–96, 2018.
- Aldo Rustichini. Minimizing regret: The general case. *Games and Economic Behavior* 19(1):224–243, 1999.
- William R Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika* 25(3-4):285–294, 12 1933.
- Hastagir P Vanchinathan, Gábor Bartók, and Andreas Krause. Efficient partial monitoring with prior information. In *Advances in Neural Information Processing Systems*, pages 1691–1699, 2014.
- Julian Zimmert and Tor Lattimore. Connections between mirror descent, thompson sampling and the information ratio. In *Advances in Neural Information Processing Systems*, pages 11973–11982, 2019.
- Martin Zinkevich. Online convex programming and generalized in-finitesimal gradient descent. In *the Twentieth International Conference on Machine Learning*, pages 928–935. AAAI Press, 2003.

A Notation

Table 1 summarizes the symbols used in this paper.

Table 1: List of symbols used in this paper.

Symbol	Meaning
P_n	$(n - 1)$ -dimensional probability simplex
$\ \cdot\ $	Euclidian norm for vector and operator norm for matrix
$\ \cdot\ _p$	p -norm
$\ \cdot\ _{K_A}$	norm induced by positive semidefinite matrix K_A
$D_{KL}(p\ q)$	KL divergence from q to p
$B_r^n(p)$	n -dimensional Euclidian ball of radius r at point $p \in \mathbb{R}^n$
$N; M \in \mathbb{Z}^+$	the number of actions and outcomes
\mathcal{F}	set of feedback symbols
A	the number of feedback symbols
$p \in P_M$	opponent's strategy
T	time horizon
$L = (l_{ij}) \in \mathbb{R}^{N \times M}$	loss matrix
$H = (h_{ij}) \in \mathbb{R}^{N \times M}$	feedback matrix
$S_i \in \mathbb{R}^A, 1 \leq i \leq N$	signal matrix
$i(t)$	action taken at time t
$N_i(t)$	the number of times the action i is taken before time $t \in [T]$
$j(t)$	outcome taken by opponent at time t
$y(t)$	feedback observed at time t
$F_t(p)$	unnormalized posterior distribution in (1)
$f_t(p)$	probability density function corresponding to $F_t(p)$
$G_t(p)$	unnormalized proposal distribution for p in (2)
$g_t(p)$	probability density function corresponding to $G_t(p)$
$q_i^{(t)} \in P_M$	empirical feedback distribution of action i by time t
$q_{i:n} \in P_M$	empirical feedback distribution of action i after the action is taken n times
$C_i \in P_M$	cell of action i

B Posterior Distribution and Proposal Distribution in Section 3

In this appendix, we discuss representation of the posterior distribution and its relation with the proposal distribution.

Proposition 4. $F_t(p)$ in (1) is proportional to the posterior distribution of the opponent's strategy, and $F_t(p) = G_t(p)$ for all $p \in P_M$.

Proof. The posterior distribution of the opponent's strategy parameter $p = f(i(s); y(s))_{s=1}^t$ is rewritten as

$$\begin{aligned}
 p &= \frac{f(i(s); y(s))_{s=1}^t}{\int_{P_M} f(i(s); y(s))_{s=1}^t dp} \\
 &= \frac{\prod_{s=1}^t \sum_{j \in \mathcal{F}} p_j \mathbb{1}[j = i(s); y = y(s)]}{\prod_{s=1}^t \sum_{j \in \mathcal{F}} \sum_{y \in \mathcal{Y}} (S_{i,y} p)^{n_{i,y}}} \\
 &= \frac{\prod_{i=1}^N \exp(-n_i D_{KL}(q_i^{(t)} \| S_i p))}{\prod_{i=1}^N \exp(-n_i D_{KL}(q_i^{(t)} \| S_i p))}; \quad (5)
 \end{aligned}$$

where $S_{i,y}$ is the i -th row of the signal matrix S_i , and note that $q_i^{(t)}$ is the empirical feedback distribution of action i at time t , that is, $q_i^{(t)} = [n_{i,1} = n_i; \dots; n_{i,A} = n_i] \in P_A$ for $n_{i,y} = \sum_{s=1}^t \mathbb{1}[i(s) = i; y(s) = y]$ and $n_i = \sum_{y=1}^A n_{i,y}$.

Next, we show that $F_t(p) = G_t(p)$ holds for all $p \in P_M$. Using Pinsker's inequality, the unnormalized posterior distribution $F_t(p)$ can be bounded from above as

$$\begin{aligned}
 F_t(p) &= \prod_{i=1}^N \exp(-n_i D_{KL}(q_i^{(t)} \| S_i p)) \\
 &\leq \prod_{i=1}^N \exp\left(-\frac{1}{2} n_i k q_i^{(t)} S_i p k^2\right) \quad (\text{by Pinsker's inequality}) \\
 &= \prod_{i=1}^N \exp\left(-\frac{1}{2} \sum_{s=1}^N n_i k q_i^{(t)} S_i p k^2\right) \\
 &= \prod_{i=1}^N \exp\left(-\frac{1}{2} \sum_{s=1}^N n_i k q_i^{(t)} S_i p k^2\right) \quad \text{by } k q_i^{(t)} S_i p k^2 = k q_i^{(t)} S_i p k^2 \\
 &= G_t(p) : \tag{6}
 \end{aligned}$$

□

Remark. The unnormalized density $Q_t(p)$ is indeed Gaussian. Recalling that B_t and b_t are defined in (3) as

$$B_t = \sum_{i=1}^N n_i S_i^> S_i = \sum_{s=1}^N S_{i(s)}^> S_{i(s)} = B_{t-1} + S_{i(t)}^> S_{i(t)}; \quad b_t = \sum_{i=1}^N n_i S_i^> q_i^{(t)} = b_{t-1} + S_{i(t)}^> e_{y(t)}; \tag{7}$$

we have

$$\begin{aligned}
 \prod_{i=1}^N n_i k q_i^{(t)} S_i p k^2 &= \prod_{i=1}^N n_i (q_i^{(t)} S_i p)^> (q_i^{(t)} S_i p) \\
 &= p^> \sum_{i=1}^N n_i S_i^> S_i p; \quad 2 \prod_{i=1}^N n_i S_i^> q_i^{(t)} p + \prod_{i=1}^N n_i k q_i^{(t)} k^2 \\
 &\quad \left| \frac{i=1}{B_t} \{z\} \right| \quad \left| \frac{i=1}{b_t} \{z\} \right| \quad \left| \frac{i=1}{c_t} \{z\} \right| \\
 &= p^> B_t p \quad 2 b_t^> p + c_t \\
 &= (p - B_t^{-1} b_t)^> B_t (p - B_t^{-1} b_t) + c_t \quad b_t^> B_t^{-1} b_t : \tag{8}
 \end{aligned}$$

Therefore, we have

$$\exp\left(-\frac{1}{2} \sum_{i=1}^N n_i k q_i^{(t)} S_i p k^2\right) / \exp\left(-\frac{1}{2} (p - B_t^{-1} b_t)^> B_t (p - B_t^{-1} b_t)\right) : \tag{9}$$

C Proof of Proposition 2

We will see that the procedure sampling from $q_t(p)$ and Algorithm 3 are equivalent. First, we derive the Gaussian density $q_t(p)$ projected onto $p \in P^M$: $\prod_{i=1}^M p_i = 1$.

For simplicity, we omit the subscript and write, e.g., B instead of B_t . We define $p = [p^{(1)}; \dots; p^{(M)}]^> \in \mathbb{R}^{M-1} \times \mathbb{R}$. Let $h = B^{-1} b$, and define $h = [h^{(1)}; \dots; h^{(M)}]^> \in \mathbb{R}^{M-1} \times \mathbb{R}$. Let $B = \begin{bmatrix} C & d \\ d^> & f \end{bmatrix}$, where $C \in \mathbb{R}^{(M-1) \times (M-1)}$; $d \in \mathbb{R}^{M-1}$, and $f \in \mathbb{R}$. Also, let $b = [b^{(1)}; \dots; b^{(M)}]^> \in \mathbb{R}^{M-1} \times \mathbb{R}$.

Using the decomposition

$$(p - B^{-1} b)^> B (p - B^{-1} b) = \underbrace{p^> B p}_{(a)} - 2 \underbrace{h^> B p}_{(b)} + h^> B h; \tag{10}$$

we rewrite each term by restricting the domain so that it satisfies the condition $\prod_{i=1}^M p_i = 1$. Now the first term (a) is rewritten as

$$\begin{aligned} (a) &= p^{(\cdot)} \left[C p^{(\cdot)} + 2 p^{(\cdot)} d p_M + f p_M^2 \right] \\ &= p^{(\cdot)} \left[C p^{(\cdot)} + 2 p^{(\cdot)} d \prod_{i=1}^{M-1} p_i + f \prod_{i=1}^{M-1} p_i^2 \right] \end{aligned} \quad (11)$$

The term (a1) is rewritten as

$$\begin{aligned} (a1) &= p^{(\cdot)} d \prod_{i=1}^{M-1} p_i \\ &= p^{(\cdot)} d \prod_{i=1}^{M-1} d_{M-1} p_i \\ &= p^{(\cdot)} d \prod_{i=1}^{M-1} D = \frac{1}{2} d_{M-1} + 1_{M-1} d \end{aligned} \quad (12)$$

and the term (a2) is rewritten as

$$\begin{aligned} (a2) &= \prod_{i=1}^{M-1} p_i^2 \\ &= \prod_{i=1}^{M-1} \frac{1}{2} p_i + \prod_{i=1}^{M-1} p_i^2 \\ &= \frac{1}{2} \prod_{i=1}^{M-1} p_i + \prod_{i=1}^{M-1} p_i^2 \end{aligned} \quad (13)$$

Therefore,

$$(a) = p^{(\cdot)} \left[\frac{C + 2D + f \prod_{i=1}^{M-1} p_i}{B} p^{(\cdot)} + 2(f \prod_{i=1}^{M-1} p_i + d) p^{(\cdot)} + f \right] \quad (14)$$

With regard to the term (b), we have

$$\begin{aligned} (b) &= b^T p \\ &= b^{(M)} p_M + b^{(M-1)} p_{M-1} \\ &= (b^{(M)} + b^{(M-1)} p_{M-1}) p^{(\cdot)} + b^{(M)} \end{aligned} \quad (15)$$

Therefore,

$$\begin{aligned} (p - B^{-1}b)^T B (p - B^{-1}b) &= p^{(\cdot)} B p^{(\cdot)} - 2 \left(f \prod_{i=1}^{M-1} p_i + d + \frac{b^{(M)} + b^{(M-1)} p_{M-1}}{b} \right) p^{(\cdot)} + f \prod_{i=1}^{M-1} p_i^2 + h^T B h \\ &= (p^{(\cdot)} - B^{-1}b)^T B (p^{(\cdot)} - B^{-1}b) + f \prod_{i=1}^{M-1} p_i^2 + h^T B h \quad \text{by } h^T B h = b^T B^{-1}b \end{aligned} \quad (16)$$

From the above argument, the density $(p - B^{-1}b; B^{-1})$ is the Gaussian distribution of (p) on \mathbb{R}^M : $\prod_{i=1}^M p_i = 1$ g. Therefore, the $p = [p^{(\cdot)}; 1]$ for $p^{(\cdot)} \sim N(B^{-1}b; B^{-1})$ is supported over \mathbb{R}^M : $\prod_{i=1}^M p_i = 1$ g.

If the sample $p^{(\cdot)}$ from $N(B^{-1}b; B^{-1})$ is in P_{M-1} , then we can obtain the last element $p^{(M)}$ by $p^{(M)} = 1 - \prod_{i=1}^{M-1} p_i$. Otherwise, the probability that $p^{(\cdot)}$ is the first $M-1$ elements of the sample from (p) is zero, and hence $p^{(\cdot)}; p^{(M)}$ cannot be a sample from (p) . Therefore, sampling p_i from $g_i(p)$ and Algorithm 3 are equivalent.

D Relation between TSPM-Gaussian and BPM-TS

In this appendix, we discuss the relation between TSPM-Gaussian and BPM-TS ([Vanchinathan et al., 2014](#)).

Underlying Feedback Structure. Here, we discuss the underlying feedback structure behind TSPM-Gaussian and BPM-TS.

We first consider the underlying feedback structure behind BPM-TS. In the following, we see that the feedback structure

$$y(t) = S_{i(t)}p + S_{i(t)} ; \quad N(0; I_M) \quad (17)$$

induces the posterior distribution in BPM-TS. Under this feedback structure, we have $N(S_{i(t)}p; S_{i(t)}S_{i(t)}^>)$.

When we take the prior distribution (p) as $N(0; \frac{1}{2}I_M)$, the posterior distribution for the opponent's strategy parameter can be written as

$$\begin{aligned} & p \text{ f } i(s); y(s) g_{s=1}^t \\ & / (p) \prod_{s=1}^Y (y(s) \text{ j } i(s); p) \\ & = (p) \prod_{s=1}^Y P_y N(S_{i(s)}p; S_{i(s)}S_{i(s)}^>) \text{ f } y = y(s) g \\ & = \exp \left[-\frac{p^> p}{2 \cdot 0} \right] \prod_{s=1}^Y \exp \left[-\frac{1}{2} (y(s) - S_{i(s)}p)^> (S_{i(s)}S_{i(s)}^>)^{-1} (y(s) - S_{i(s)}p) \right] \\ & = \exp \left[-\frac{1}{2} p^> \left(\frac{1}{2} I_M + \sum_{s=1}^Y S_{i(s)}^> (S_{i(s)}S_{i(s)}^>)^{-1} S_{i(s)} \right) p \right] \\ & \quad - \frac{1}{2} \sum_{s=1}^Y y(s)^> (S_{i(s)}S_{i(s)}^>)^{-1} S_{i(s)} p + (\text{a term independent of } p) \\ & / \exp \left[-\frac{1}{2} (p^> B_t^{\text{BPM}} p - 2b_t^{\text{BPM}} p) \right] \\ & / \exp \left[-\frac{1}{2} (p - B_t^{\text{BPM}}^{-1} b_t^{\text{BPM}})^> B_t^{\text{BPM}} (p - B_t^{\text{BPM}}^{-1} b_t^{\text{BPM}}) \right]; \end{aligned} \quad (18)$$

where

$$B_t^{\text{BPM}} = \frac{1}{2} I_M + \sum_{s=1}^Y S_{i(s)}^> (S_{i(s)}S_{i(s)}^>)^{-1} S_{i(s)} = B_t^{\text{BPM}} + S_{i(t)}^> (S_{i(t)}S_{i(t)}^>)^{-1} S_{i(t)}; \quad (19)$$

$$b_t^{\text{BPM}} = \sum_{s=1}^Y S_{i(s)}^> (S_{i(s)}S_{i(s)}^>)^{-1} y(s) = b_t^{\text{BPM}} + S_{i(t)}^> (S_{i(t)}S_{i(t)}^>)^{-1} y(t); \quad (20)$$

Therefore, the posterior distribution $p \text{ j } f i(s); y(s) g_{s=1}^t$ is

$$\frac{1}{(2\pi)^M \text{ j } B_t^{\text{BPM}}^{-1}} \exp \left[-\frac{1}{2} (p - B_t^{\text{BPM}}^{-1} b_t^{\text{BPM}})^> B_t^{\text{BPM}} (p - B_t^{\text{BPM}}^{-1} b_t^{\text{BPM}}) \right]; \quad (21)$$

and this distribution indeed corresponds to the posterior distribution in BPM-TS ([Vanchinathan et al., 2014](#)) with $B_t^{\text{BPM}} = B_t^{-1}$.

Using the same argument, we can confirm that the feedback structure

$$y_t = S_i p + ; \quad N(0; I_M); \quad (22)$$

induces

$$g_t(p) := \frac{1}{(2\pi)^M |B_t|} \exp\left\{-\frac{1}{2} p^T B_t^{-1} p\right\}; \quad (23)$$

which corresponds to the posterior distribution for TSPM in linear partial monitoring.

Covariances in TSPM-Gaussian and BPM-TS. In the linear partial monitoring, TSPM assumes noise with covariance Σ_M , which is compatible with the fact that the discrete setting can be regarded as linear PM with Σ_M -sub-Gaussian noise. On the other hand, BPM-TS assumes covariance Σ_t and in general $\Sigma_t \succeq \Sigma_M$ holds. Therefore, BPM-TS assumes unnecessarily larger covariance, which makes learning slow down.

E Preliminaries for Regret Analysis

In this appendix, we give some technical lemmas, which are used for the derivation of the regret bound in Appendix F. Here, we write χ^2_k to denote $\chi^2_k(0)$. For $a, b \in \mathbb{R}$, let $a \wedge b$ be $\min\{a, b\}$ if $a, b > 0$ otherwise, and $a \vee b$ be $\max\{a, b\}$ if $a, b > 0$ otherwise. We use $\chi^2_k(a) := P_{X \sim \chi^2_k}$ to evaluate the behavior of the posterior samples, where χ^2_k is the chi-squared distribution with k degree of freedom.

E.1 Basic Lemmas

Fact 5 (Moment generating function of squared-Gaussian distribution) Let X be the random variable following the standard normal distribution. Then, the moment generating function of X^2 is $E \exp(-tX^2) = (1 - 2t)^{-1/2}$ for $t < 1/2$.

Lemma 6 (Chernoff bound for chi-squared random variable) Let X be the random variable following the chi-squared distribution with k degree of freedom. Then, for any $0 < \delta < 1/2$,

$$P\{X \geq (1 + \delta)k\} \leq e^{-\frac{\delta^2 k}{2}}; \quad (24)$$

Proof. By Markov's inequality, the LHS can be bounded as

$$\begin{aligned} P\{X \geq (1 + \delta)k\} &= P\left\{\sum_{i=1}^k X_i^2 \geq (1 + \delta)k\right\} \quad (X_1, \dots, X_k \text{ i.i.d. } N(0, 1)) \\ &= P\left\{\exp\left(\sum_{i=1}^k X_i^2\right) \geq \exp((1 + \delta)k)\right\} \\ &\leq e^{-a} E \exp\left(\sum_{i=1}^k X_i^2\right) \quad (\text{by Markov's ineq}) \\ &= e^{-a} (1 - 2)^{-\frac{k}{2}} \quad (\text{by Fact 5}); \end{aligned} \quad (25)$$

which completes the proof. \square

E.2 Property of Strong Local Observability

Recall that $\gamma_i = (L_i - L_1) \vee p > 0$ for $i \in [N]$, which is the difference of the expected loss of actions i and 1. For this define

$$\gamma := \frac{1}{2} \min_{i \in [N]} \frac{\gamma_i}{k} \wedge \min_{p \in \mathcal{C}_i} \frac{4}{3} k p \vee p; \quad (26)$$

which is used throughout the proof of this appendix and Appendix F. The following lemma provides the key property of the strong local observability condition.

Lemma 7. For any partial monitoring game with strong local observability and $\mathbf{a} \in \mathbb{R}^M$, any of the conditions 1-3 in the following is not satisfied:

1. $L_1 \vee p > L_k \vee p$ (Worse action looks better under p .)
2. $kS_1 p \leq S_1 p$

$$3. \|S_k p - S_k p\|_k \leq \dots$$

Proof. We prove by contradiction. Assume that there exists $p \in \mathbb{R}^M$ such that conditions 1-3 are simultaneously satisfied.

Now, by the conditions 2 and 3, we have

$$\begin{aligned} |S_1 p - S_1 p|_j &\leq \frac{1}{A}; \\ |S_k p - S_k p|_j &\leq \frac{1}{A}. \end{aligned} \quad (27)$$

Here, $|j|$ is the element-wise absolute value, and \leq means that the inequality holds for each element. Therefore,

$$\begin{pmatrix} S_1 \\ S_k \end{pmatrix} (p - p) \leq \frac{1}{2A}. \quad (28)$$

On the other hand, by the strong local observability condition, for any $i \in \{1, \dots, k\}$, there exists $z_{1;k} \in \mathbb{R}^{2A}$ such that

$$(L_1 - L_k)^T z_{1;k} = \begin{pmatrix} S_1 \\ S_k \end{pmatrix}^T. \quad (29)$$

Now, we have

$$\begin{aligned} z_{1;k}^T \begin{pmatrix} S_1 \\ S_k \end{pmatrix} &\leq (p - p)^T \begin{pmatrix} S_1 \\ S_k \end{pmatrix} \\ \|z_{1;k}\|_k &\leq \frac{1}{\sigma_k} \begin{pmatrix} S_1 \\ S_k \end{pmatrix}^T (p - p) \quad (\text{by Cauchy-Schwarz inequality}) \\ &\leq \frac{1}{\sigma_k} \sqrt{2A} \|z_{1;k}\|_k \quad (\text{by Eq. (28)}); \end{aligned} \quad (30)$$

and

$$\begin{aligned} z_{1;k}^T \begin{pmatrix} S_1 \\ S_k \end{pmatrix} &= (L_1 - L_k)^T z_{1;k} (p - p) \quad (\text{by Eq. (29)}) \\ &= (L_1 - L_k)^T z_{1;k} p + (L_k - L_1)^T z_{1;k} p \\ &\leq \frac{1}{\sigma_k} \begin{pmatrix} S_1 \\ S_k \end{pmatrix}^T p \quad (\text{by Condition 1 \& def. of } \sigma_k); \end{aligned} \quad (31)$$

Therefore, from (30) and (31), we have

$$\frac{1}{\sigma_k} \sqrt{2A} \|z_{1;k}\|_k \leq \frac{1}{\sigma_k} \begin{pmatrix} S_1 \\ S_k \end{pmatrix}^T p. \quad (32)$$

This inequality does not hold for all $k \in \{1, \dots, k\}$ for the predefined value of p , since we have

$$\frac{1}{\sqrt{2A}} \min_{k \in \{1, \dots, k\}} \frac{1}{\sigma_k} > \frac{1}{\sigma_k} \begin{pmatrix} S_1 \\ S_k \end{pmatrix}^T p. \quad (33)$$

Therefore, the proof is completed by contradiction. \square

Remark. The similar result holds when the optimal actions are replaced with actions $j \in \{1, \dots, k\}$ such that $\sigma_{j;k} := (L_j - L_k)^T p > 0$ by taking $z_{j;k}$ satisfying

$$\frac{1}{\sqrt{2A}} \min_{j \in \{1, \dots, k\}: \sigma_{j;k} > 0} \frac{1}{\sigma_{j;k}} > \frac{1}{\sigma_{j;k}} \begin{pmatrix} S_j \\ S_k \end{pmatrix}^T p. \quad (34)$$

From Lemma 7, we have the following corollary.

Corollary 8. For any $p \in \mathbb{R}^M$ satisfying $(L_i - L_j)^T p > 0$ and $\|S_1 p - S_1 p\|_k \leq \frac{1}{A}$, we have

$$\|S_1 p - S_1 p\|_k > \frac{1}{A}. \quad (35)$$

Proof. Note that $(L_i - L_j)^T p > 0$ for any $i \in \{1, \dots, k\}$. Therefore, the result directly follows from Lemma 7. \square

The next lemma is the property of Mahalanobis distance corresponding to

Lemma 9. Define $T_i := \{p \in \mathbb{R}^M : \|S_i p - S_i p^*\| > \epsilon\}$. Assume that $N_i(t) \geq n_i$, $\|S_i p^* - S_i p\| \leq \epsilon$. Then, for any $0 < \delta < 1/2$

$$\mathbb{P} \left(\inf_{p \in T_i} \|B_t^{1/2} (p - p^*)\|^2 \leq \frac{9}{16} \epsilon^2 n_i (1 - \delta)^{M-2} \right) \leq \delta \quad (36)$$

Proof. To bound the LHS of the above inequality, we bound $\|B_t^{1/2} (p - p^*)\|^2$ from below for $p \in T_i$. Using the triangle inequality and the assumptions, we have

$$\|S_i (p - p^*)\| \geq \|S_i p - S_i p^*\| - \|S_i p^* - S_i p\| \geq \epsilon - \epsilon = 0 \quad (37)$$

Therefore, we have

$$\begin{aligned} \|B_t^{1/2} (p - p^*)\|^2 &= \sum_{k \in [N]} N_k(t) \|S_k (p - p^*)\|^2 \quad (\text{by def. of } B_t) \\ &\geq \sum_{k \in [N]} n_k \|S_k (p - p^*)\|^2 \quad (N_i(t) \geq n_i) \\ &> \frac{9}{16} \epsilon^2 n_i \quad (\text{by Eq. (37)}) \end{aligned} \quad (38)$$

By the Chernoff bound for a chi-squared random variable in Lemma 6, we now have

$$\mathbb{P} \left(\sum_{k \in [N]} N_k(t) \|S_k (p - p^*)\|^2 \leq \frac{9}{16} \epsilon^2 n_i (1 - \delta)^{M-2} \right) \leq \delta \quad (39)$$

for any $\delta > 0$ and $0 < \delta < 1/2$. Hence, using the fact that $\|B_t^{1/2} (p - p^*)\|^2$ follows the chi-squared distribution with M degree of freedom, we have

$$\mathbb{P} \left(\inf_{p \in T_i} \|B_t^{1/2} (p - p^*)\|^2 \leq \frac{9}{16} \epsilon^2 n_i \right) \leq \delta \exp \left(-\frac{9}{16} \epsilon^2 n_i (1 - \delta)^{M-2} \right) \quad (40)$$

which completes the proof. \square

E.3 Statistics of Uninterested Actions

For any $k \in [K]$ and $n_k \in [T]$, define

$$Z_{n_k} := \sum_{k \in [K]} \sum_{n_k \in [T]} S_k p^*{}^2 \quad (41)$$

$$Z_{n_i} := \max_{k \in [K]} Z_{n_k} \quad (42)$$

In this section, we bound Z_{n_i} from above. Note that Z_{n_i} is independent of the randomness of Thompson sampling.

Lemma 10 (Upper bound for the expectation of Z_{n_i}).

$$\mathbb{E} Z_{n_i} \leq 4N \log T + \frac{A}{2} \log 2 + 1 \quad (43)$$

Proof. Recall that in linear partial monitoring, the feedback $y_t \in \mathbb{R}^A$ for action k is given as

$$y_t = S_k p^* + \xi_t, \quad \xi_t \sim \mathcal{N}(0; I_A) \quad (44)$$

at round $t \in [T]$. Therefore $y(t) = S_k p^* + \xi_t \sim \mathcal{N}(0; I_A)$. Since $\alpha_{k;n_k} = \frac{1}{n_k} \sum_{s \in [T]: i(s)=k} y(s)$ for any $n_k \in [T]$, we have

$$\alpha_{k;n_k} - S_k p^* = \frac{1}{n_k} \sum_{s \in [T]: i(s)=k} (y(s) - S_k p^*) \sim \mathcal{N}(0; I_A/n_k) \quad (45)$$

Therefore,

$$P_{\bar{n}_k(\mathbf{a}_k; n_k, S_k p)} \sim N(0; I_A); \quad (46)$$

and thus

$$n_k k \mathbf{a}_k; n_k, S_k p k^2 = k^2 P_{\bar{n}_k(\mathbf{a}_k; n_k, S_k p)} k^2 \stackrel{2}{A} : \quad (47)$$

Therefore, for any $\theta \leq 1=2$,

$$\begin{aligned} E \max_{n_k 2[T]} Z_{n_k} &= P \max_{n_k 2[T]} Z_{n_k} x dx \\ &= \int_0^1 [1 \wedge T P f Z_1 x g] dx \quad (\text{by the union bound}) \\ &= \int_0^1 h \int_0^1 1 \wedge T e^{-x(1-2)^{\frac{A}{2}}} dx \quad (\text{by } Z_1 \stackrel{2}{A} \text{ and Lemma 6}) \\ &= \int_0^1 dx + \int_0^1 T e^{-x(1-2)^{\frac{A}{2}}} dx \\ &= x + T \int_0^1 e^{-x(1-2)^{\frac{A}{2}}} dx \\ &= x + T(1-2)^{\frac{A}{2}} \frac{1-e^{-x}}{x} \Big|_0^1 \\ &= \frac{1}{2} \log T + \frac{A}{2} \log(1-2) + 1; \end{aligned} \quad (48)$$

where $x := \frac{1}{2} \log T + \frac{A}{2} \log(1-2)$. Therefore, taking $\theta = 1=4$, we have

$$\begin{aligned} E Z_{n_i} &= E \max_{k \in i} Z_{n_k} \\ &= E \max_{k \in i} Z_{n_k} \\ &= (N-1) \frac{1}{2} \log T + \frac{A}{2} \log(1-2) + 1 \\ &= 4N \log T + \frac{A}{2} \log 2 + 1; \end{aligned} \quad (49)$$

which completes the proof. \square

E.4 Mahalanobis Distance Process

Discussions in this section are essentially very similar to [Abbasi-Yadkori et al. \(2011, Lemma 11\)](#), but their results are not directly applicable and we give the full derivation for self-containedness. To maximize the applicability here we only assume sub-Gaussian noise rather than a Gaussian one.

Let t be zero-mean sub-Gaussian random variable, which satisfies

$$E e^{> t} = e^{\frac{k k^2}{2}} \quad (50)$$

for any $x \in \mathbb{R}^M$.

Lemma 11. For any vector $v \in \mathbb{R}^M$ and positive definite matrix $V \in \mathbb{R}^M \times \mathbb{R}^M$ such that $V \succeq I$,

$$E_t e^{\frac{k t + v k^2}{2}} = \frac{P \int \bar{v} \bar{j}}{j \sqrt{V^{-1} j}} e^{\frac{1}{2} v^T (V^{-1})^{-1} v}; \quad (51)$$

Proof. For any $x \in \mathbb{R}^M$

$$E_{N(0; V^{-1})} e^{> x} = e^{\frac{k k^2}{2}} : \quad (52)$$

Therefore, by letting $\alpha = t + v$ we see that

$$E_t e^{\frac{k^T t + vk^2}{2} \mathbf{1}} = E_{N(0;V^{-1})} e^{h^T (\alpha + v) \mathbf{i}} \quad (53)$$

As a result, by the definition of sub-Gaussian random variables, we have

$$\begin{aligned} E_t e^{\frac{k^T t + vk^2}{2} \mathbf{1}} &= E_{N(0;V^{-1})} E_t e^{h^T (\alpha + v) \mathbf{i}} \\ &= E_{N(0;V^{-1})} e^{h^T v \mathbf{i}} E_t e^{h^T \alpha \mathbf{i}} \\ &= E_{N(0;V^{-1})} e^{h^T v \mathbf{i}} E_{Z \sim N(0;I)} e^{k^T Z} \\ &= \frac{1}{(2\pi)^{d/2}} \int_{\mathbb{R}^d} \frac{1}{\sqrt{|V^{-1}|}} e^{h^T v \mathbf{i}} e^{k^T Z} e^{-\frac{1}{2} Z^T V^{-1} Z} dZ \\ &= \frac{1}{(2\pi)^{d/2}} \int_{\mathbb{R}^d} \frac{1}{\sqrt{|V^{-1}|}} e^{\frac{1}{2} (v^T (V^{-1})^{-1} v)} e^{\frac{1}{2} Z^T (V^{-1})^{-1} Z} dZ \\ &= \frac{1}{(2\pi)^{d/2}} \int_{\mathbb{R}^d} \frac{1}{\sqrt{|V^{-1}| |V^{-1}|}} e^{\frac{1}{2} (v^T (V^{-1})^{-1} v)} e^{\frac{1}{2} Z^T (V^{-1})^{-1} Z} dZ \\ &= \frac{1}{\sqrt{|V^{-1}|}} e^{\frac{1}{2} v^T (V^{-1})^{-1} v} \end{aligned} \quad (54)$$

□

Lemma 12.

$$E \exp \left\{ \frac{1}{2} k^T \rho_t - \frac{1}{2} k^T B_t k - \frac{1}{2} k^T \rho_{t-1} - \frac{1}{2} k^T B_{t-1} \rho_{t-1}; B_{t-1}; S_{i(t-1)} \right\} = \frac{\sqrt{|B_t|}}{\sqrt{|B_{t-1}|}} \quad (55)$$

Proof. Let $Z_t := \rho_t + \sum_{s=1}^t S_{i(s)}^> y(s)$, and we have

$$\begin{aligned} B_t &= I + \sum_{s=1}^t S_{i(s)}^> S_{i(s)}, \\ \rho_t &= \sum_{s=1}^t S_{i(s)}^> y(s) = B_t \rho + Z_t, \\ \rho_t &= B_t^{-1} \rho_t = \rho + B_t^{-1} Z_t. \end{aligned}$$

In the following we omit the conditioning on $(\rho_{t-1}; B_{t-1}; S_{i(t-1)})$ for notational simplicity.

Let us define $\tilde{C}_t := S_{i(t)} B_{t-1} S_{i(t)}^>$ and $\tilde{d}_t := S_{i(t)} B_{t-1}^{-1} Z_{t-1} = S_{i(t)} (\rho_{t-1} - \rho)$. Then, using the Sherman-Morrison-Woodbury formula we have

$$\begin{aligned}
& k\hat{p}_t - p k_{B_t}^2 - k\hat{p}_{t-1} - p k_{B_{t-1}}^2 \\
&= Z_t^> B_t^{-1} Z_t - Z_{t-1}^> B_{t-1}^{-1} Z_{t-1} \\
&= (Z_{t-1}^> + \hat{S}_{i(t)}) (B_{t-1}^{-1} - B_t^{-1} S_{i(t)}^>) (I + S_{i(t)} B_t^{-1} S_{i(t)}^>)^{-1} S_{i(t)} B_t^{-1} (Z_{t-1} + S_{i(t)}^>) - Z_{t-1}^> B_{t-1}^{-1} Z_{t-1} \\
&= (Z_{t-1}^> + \hat{S}_{i(t)}) B_t^{-1} (Z_{t-1} + S_{i(t)}^>) - Z_{t-1}^> B_{t-1}^{-1} Z_{t-1} \\
&\quad (Z_{t-1}^> + \hat{S}_{i(t)}) B_t^{-1} S_{i(t)}^> (I + S_{i(t)} B_t^{-1} S_{i(t)}^>)^{-1} S_{i(t)} B_t^{-1} (Z_{t-1} + S_{i(t)}^>) \\
&= \hat{S}_{i(t)} B_t^{-1} S_{i(t)}^> + 2 Z_{t-1}^> B_t^{-1} S_{i(t)}^> \\
&\quad (Z_{t-1}^> + \hat{S}_{i(t)}) B_t^{-1} S_{i(t)}^> (I + S_{i(t)} B_t^{-1} S_{i(t)}^>)^{-1} S_{i(t)} B_t^{-1} (Z_{t-1} + S_{i(t)}^>) \\
&= \hat{C}_t + 2 d_t^> - (d_t^> + \hat{C}_t) (I + C_t)^{-1} (d_t + C_t) \\
&= \hat{C}_t (I - (I + C_t)^{-1} C_t) + 2 d_t^> (I - (I + C_t)^{-1} C_t) - d_t^> (I + C_t)^{-1} d_t \\
&= \hat{C}_t (I + C_t)^{-1} + 2 d_t^> (I + C_t)^{-1} - d_t^> (I + C_t)^{-1} d_t \\
&= \hat{C}_t + C_t^{-1} d_t^2 \frac{2}{C_t(I+C_t)^{-1}} - d_t^> (I + C_t)^{-1} C_t^{-1} d_t - d_t^> (I + C_t)^{-1} d_t \\
&= \hat{C}_t + C_t^{-1} d_t^2 \frac{2}{C_t(I+C_t)^{-1}} - d_t^> (I + C_t)^{-1} (I + C_t)^{-1} d_t : \tag{56}
\end{aligned}$$

Therefore, Lemma 11 with $W := C_t(I + C_t)^{-1} = (I + C_t)C_t^{-1}$; $v := C_t^{-1}d_t$ yields

$$\begin{aligned}
& E \exp \frac{1}{2} k\hat{p}_t - p k_{B_t}^2 - k\hat{p}_{t-1} - p k_{B_{t-1}}^2 \\
&\quad q \frac{j(I + C_t)C_t^{-1}j}{j(I + C_t)C_t^{-1}j} e^{\frac{1}{2}d_t^> C_t^{-1}((I + C_t)C_t^{-1} - I)^{-1} C_t^{-1}d_t} e^{-\frac{1}{2}d_t^> (I + C_t)^{-1}(I + C_t)^{-1}d_t} \\
&\quad q \frac{j(I + C_t)C_t^{-1}j}{j(I + C_t)C_t^{-1}j} e^{\frac{1}{2}d_t^> C_t^{-1}(C_t^{-1})^{-1} C_t^{-1}d_t} e^{-\frac{1}{2}d_t^> (I + C_t)^{-1}(I + C_t)^{-1}d_t} \\
&= \frac{p}{s} \frac{j(I + C_t)j}{j(I + C_t)j} \\
&= \frac{jB_t j}{jB_{t-1} j}; \tag{57}
\end{aligned}$$

where see, e.g., [Abbasi-Yadkori et al. \(2011, Lemma 11\)](#) for the last equality. \square

E.5 Norms under Perturbations

In the following two lemmas, we give some analysis of norms under perturbations.

Lemma 13. Let A be a positive definite matrix. Let $\alpha \in \mathbb{R}^d$ and $\beta > 0$ be such that $\|\alpha\|_2 < \beta$. Then

$$\min_{x: \|x\|_2 \leq \beta} \max_{x^0: \|x^0\|_2 \leq \beta} (a + x + x^0)^T A (a + x + x^0) = \min_{x^0: \|x^0\|_2 \leq \beta} (a + x^0)^T A (a + x^0) : \tag{58}$$

Proof. By considering the Lagrangian multiplier we see that any stationary point of the function $(a + x^0)^T A (a + x^0)$ over $f(x; x^0) : \|x\|_2 \leq \beta; \|x^0\|_2 \leq \beta$ satisfies

$$\begin{aligned}
& A(a + x + x^0) - \lambda_1 x = 0; \\
& A(a + x + x^0) - \lambda_2 x^0 = 0; \\
& \lambda_1 \beta - \lambda_2 \beta = 2\beta; \\
& \lambda_1 x^0 - \lambda_2 x^0 = 2\beta; \tag{59}
\end{aligned}$$

and therefore $\lambda_1 x = \lambda_2 x^0$. Considering the last two conditions (59) we have $\lambda_2 = \lambda_1$, implying that

$$x^0 = (3A - 2I)\alpha \tag{60}$$

or

$$x^0 = (A - 2I)Aa \quad (61)$$

for λ_1 satisfying $x^0 > x^0 = \lambda_1^2$.

Note that it holds for any positive definite matrix B that

$$\frac{d^2}{d\lambda^2} a(B + I) \lambda^2 a = a(B + I) \lambda^4 a = (B + I) \lambda^2 a^2; \quad (62)$$

which is positive almost everywhere, meaning that $a(B + I) \lambda^2 a$ is strictly convex with respect to $\lambda \in \mathbb{R}$. Therefore, there exists at most two λ 's satisfying (60) and $x^0 > x^0 = \lambda^2$, and there exists at most two λ 's satisfying (61) and $x^0 > x^0 = \lambda^2$. In summary, there are at most four stationary points of $(a + x^0) > A(a + x^0)$ over $f(x; x^0) : \|x\| \leq 2; \|x^0\| \leq g$.

On the other hand, two optimization problems

$$\min_{x: \|x\| \leq 2} \min_{x^0: \|x^0\| \leq g} (a + x + x^0) > A(a + x + x^0) = \min_{x^0: \|x^0\| \leq g} \min_{x: \|x\| \leq 2} (a + x^0) > A(a + x^0) \quad (63)$$

and

$$\max_{x: \|x\| \leq 2} \max_{x^0: \|x^0\| \leq g} (a + x + x^0) > A(a + x + x^0) = \max_{x^0: \|x^0\| \leq g} \max_{x: \|x\| \leq 2} (a + x^0) > A(a + x^0) \quad (64)$$

can be easily solved by an elementary calculation and the optimal values are equal to those corresponding to (60).

Therefore, the optimal solutions of the two minimax problems

$$\max_{x: \|x\| \leq 2} \min_{x^0: \|x^0\| \leq g} (a + x + x^0) > A(a + x + x^0) \quad (65)$$

and

$$\min_{x: \|x\| \leq 2} \max_{x^0: \|x^0\| \leq g} (a + x + x^0) > A(a + x + x^0) \quad (66)$$

correspond to two points corresponding to (61).

We can see again from an elementary calculation that the optimal solutions for two optimization problems

$$\begin{aligned} \min_{x^0: \|x^0\| \leq g} (a + x^0) > A(a + x^0) \\ \max_{x^0: \|x^0\| \leq g} (a + x^0) > A(a + x^0) \end{aligned} \quad (67)$$

have the same necessary and sufficient conditions (61) and we complete the proof by noticing that (65) is less than (66). \square

Lemma 14. Let $A \in \mathbb{S}_1^n$ be a positive-definite matrix with minimum eigenvalue at least 0. Then, for any $\rho \in \mathbb{R}^d$ and $\epsilon > 0$ satisfying $\epsilon < \rho^0 - \rho$, $k=3$,

$$\rho^0 - \rho \leq k_A^2 \inf_{p: \|p\| \leq \rho} \sup_{p^0: \|p^0\| \leq \rho} k p^0 - \rho k_A^2 \frac{\rho}{n} \|S_1(\rho - p)\|; \quad (68)$$

Proof. Let $a = \rho - p$. By Lemma 13, we have

$$\begin{aligned} & \inf_{p: \|p\| \leq \rho} \sup_{p^0: \|p^0\| \leq \rho} k p^0 - \rho k_A^2 \\ &= \inf_{x: \|x\| \leq 2} \sup_{x^0: \|x^0\| \leq g} k a + x + x^0 k_A^2 \\ &= \inf_{x: \|x\| \leq 2} k a + x k_A^2; \end{aligned} \quad (69)$$

On the other hand, since $\sum_{k=2}^{\infty} C_k + B^d(0) < \infty$ we can take $\epsilon > 0$ such that $\sum_{k=2}^{\infty} C_k < \epsilon$. Hence,

$$\begin{aligned} z_{1;k}^{\geq} &= \sum_{S_1}^{\infty} (\rho - p) = (L_1 - L_k)^{\geq} (\rho - p) \\ &= (L_k - L_1)^{\geq} (\rho - p) + (L_1 - L_k)^{\geq} p + (L_k - L_1)^{\geq} p \\ &= (L_k - L_1)^{\geq} (\rho - p) + \epsilon \quad (\text{by } \sum_{k=2}^{\infty} C_k \text{ and def. of } \epsilon) \end{aligned} \quad (76)$$

From (75) and (76), we have

$$\sum_{k=2}^{\infty} (L_k - L_1)^{\geq} (\rho - p) \leq \frac{p}{2A} \sum_{k=2}^{\infty} C_k \quad (77)$$

Now, the left hand side of (77) is bounded from below as

$$\begin{aligned} \sum_{k=2}^{\infty} (L_k - L_1)^{\geq} (\rho - p) &\geq \sum_{k=2}^{\infty} (L_k - L_1) \frac{\rho - p}{k} \\ &= \sum_{k=2}^{\infty} (L_k - L_1) \frac{1}{k} \frac{\rho - p}{\max_i \frac{kL_1 - L_i}{kz_{1;j}}} \\ &= \sum_{k=2}^{\infty} (L_k - L_1) \frac{1}{k} \frac{\rho - p}{\max_i \frac{kL_1 - L_i}{kz_{1;j}}} \end{aligned} \quad (78)$$

On the other hand, using the definition of ρ , the right hand side of (77) is bounded from above as

$$\frac{p}{2A} \sum_{k=2}^{\infty} C_k < \sum_{k=2}^{\infty} C_k \quad (79)$$

Therefore, the proof is completed by contradiction. \square

E.6 Exit Time Analysis

We next consider the exit time. Let A_t be an event deterministic given F_t , and B_t be a random event such that if B_t occurred then A_t never occurs for $t = t + 1; t + 2; \dots$. Let $P_t; t = 1; 2; \dots; T$, be a stochastic process satisfying $P_t = P(B_t | F_t)$ a.s. and P_t^{-1} is a supermartingale with respect to the filtration induced by F_t .

Theorem 16. Let τ be the stopping time defined as

$$\tau = \min\{t \geq 1 : A_t \text{ occurs for some } t \in [1, \tau]\} \quad (80)$$

Then we almost surely have

$$E \sum_{t=1}^{\tau} 1[A_t] = \sum_{t=1}^{\tau} P_t = \tau \quad (81)$$

We prove this theorem based on the following lemma.

Lemma 17. Let $(Q_i)_{i=1}^{\infty} \in [0, 1]$ be an arbitrary stochastic process such that $(Q_i^{-1})_{i=1}^{\infty}$ is a supermartingale with respect to a filtration $(G_i)_{i=1}^{\infty}$. Then, for any $G_0 \subset G_1$,

$$E \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} (1 - Q_j) Q_i \leq E \sum_{j=1}^{\infty} Q_j \quad (82)$$

Proof. Let

$$\begin{aligned} N_k((Q_i; G_i)_{i=1}^{\infty}; G_0) &= E \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} (1 - Q_j) Q_i \\ \bar{N}_k((Q_i; G_i)_{i=1}^{\infty}; G_0) &= E \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} (1 - Q_j) Q_i \quad \text{where } Q_j = Q_k \text{ for } j > k. \end{aligned} \quad (83)$$

We show $\bar{N}_k((Q_i; G)_{i=1}^1; G_0) = E[Q_1^1 | G_0] = 1$ a.s. for any $(Q_i; G)_{i=1}^1, G_0 = G_1$ and $k \geq 2, N$ by induction. First, for $k = 1$ the statement holds since

$$\begin{aligned} \bar{N}_1((Q_i; G)_{i=1}^1; G_0) &= E \left[\prod_{i=1}^1 \sum_{j=1}^3 (1 - Q_i) G_0^5 \right] \\ &= E [Q_1^1 = 1 | G_0] \\ &= E [Q_1^1 = G_0 = 1] \end{aligned} \tag{84}$$

Next, assume that the statement holds for $(Q_i; G)_{i=1}^k, G_0 = G_1$ and $k \leq k_0$. Then, we almost surely have

$$\begin{aligned} \bar{N}_{k_0+1}((Q_i; G)_{i=1}^1; G_0) &= E \left[\prod_{i=1}^2 (1 - Q_i) \right] E \left[\prod_{j=2}^3 (1 - Q_j) G_1^5 G_0^5 \right] \\ &= E [(1 - Q_1)(1 + \bar{N}_{k_0}((Q_i; G)_{i=2}^1; G_1)) | G_0] \\ &= E [(1 - Q_1)E[Q_2^1 | G_1] | G_0] \quad (\text{assumption of the induction}) \\ &= E [Q_1^1 = G_0 = 1 - Q_1^1] \text{ is a supermartingale.} \end{aligned} \tag{85}$$

We obtain the lemma from

$$E \left[\prod_{i=1}^k \sum_{j=1}^3 (1 - Q_j) G_0^5 \right] = \bar{N}_k((Q_i; G)_{i=1}^1; G_0) = \bar{N}_k((Q_i; G)_{i=1}^1; G_0) \text{ a.s.} \tag{86}$$

□

Proof of Theorem 16. The statement is obvious for the case $t = T + 1$ and we consider the other case in the following.

Let τ_i be the time of the i -th occurrence of A_t . More formally, we define τ_i as the stopping time $\tau_1 = 1$ and

$$\tau_{i+1} = \begin{cases} \min_{t \geq \tau_i + 1} \{t \in [T] : P_{\tau_i}^T \mathbb{1}[A_{t^0}] = i + 1\} & \text{if } P_{\tau_i}^T \mathbb{1}[A_{\tau_i^0}] = i + 1; \\ \tau_i + 1 & \text{otherwise.} \end{cases} \tag{87}$$

Then $(P_i^0) = (P_{\tau_i}^0)$ is a stochastic process measurable by the filtration induced by $(F_i) = (F_{\tau_i})$. By Lemma 17 we obtain

$$\begin{aligned} E \left[\prod_{t=1}^n \sum_{j=1}^3 (1 - P_j^0) F_1^0 \right] &= E \left[\prod_{n=1}^n \sum_{t=1}^n \sum_{j=1}^3 (1 - P_j^0) F_1^0 \right] \\ &= 1 + E \left[\sum_{t=1}^{n-1} \sum_{j=1}^3 (1 - P_j^0) F_1^0 \right] \\ &= 1 + E \left[\sum_{i=1}^{n-2} \sum_{j=1}^3 (1 - P_j^0) F_1^0 \right] \\ &= 1 + E [P_1^0 = 1] \\ &= P_1^0 \end{aligned} \tag{88}$$

□

F Regret Analysis of TSPM Algorithm

In this appendix, we give the proof of Theorem 3. Note that the cells are defined for the decomposition of R^M , not P_M . In other words, the cell C_i is here defined as $C_i = \{p \in R^M : \text{action } i \text{ is optimal}\}$.

For the linear setting, the empirical feedback distribution $q_i^{(t)}$ and $q_{i,n}$ are defined as

$$q_i^{(t)} := \frac{1}{N_i(t)} \sum_{s \in [t-1]: i(s)=i} y(s); \quad (89)$$

$$q_{i,n} := \text{the value of } q_i^{(t)} \text{ after taking action } i \text{ for } n \text{ times.} \quad (90)$$

Recall that $\rho_t = B_t^{-1} b_t$, which is the mode of $q_t(p)$.

F.1 Regret Decomposition

Here, we break the regret into several terms. For any $n \in [N]$, we define events

$$A_i(t) := \{k S_i \rho_t \leq S_i p \leq k \frac{\rho}{4}\}; \quad (91)$$

$$A_i^c(t) := \{k S_i \rho_t \leq S_i p \leq k \frac{\rho}{8}\}; \quad (92)$$

We first decompose the regret as

$$\begin{aligned} \text{Reg}(T) &= \sum_{t=1}^T \sum_{i(t)} h_i \\ &= \sum_{t=1}^T \sum_{i(t)} h_i \mathbb{1}_{A_i(t)} + \max_{j \in [N]} \sum_{t=1}^T \sum_{j(t)} h_j \mathbb{1}_{A_j^c(t)} \\ &= \sum_{i \in [1]} \sum_{t=1}^T \sum_{i(t)=i} h_i \mathbb{1}_{A_i(t)} + \max_{j \in [2[N]]} \sum_{t=1}^T \sum_{j(t)=j} h_j \mathbb{1}_{A_j^c(t)} \\ &= \sum_{i \in [1]} \sum_{t=1}^T \underbrace{\sum_{i(t)=i} h_i \mathbb{1}_{A_i(t)}}_{(A)} + \underbrace{\sum_{i(t)=i} h_i \mathbb{1}_{A_i^c(t)}}_{(B)} + \max_{j \in [2[N]]} \sum_{t=1}^T \sum_{j(t)=j} h_j \mathbb{1}_{A_j^c(t)}; \end{aligned} \quad (93)$$

To decompose the last term, we define the following notation. We define for any $n \in [N]$

$$P_i(t) := P \{ \rho_t \geq 2 C_i \mid F_t, g \}; \quad (94)$$

We also define

$$C_{i,t} := C_i \setminus B_{\rho(t)}; \quad (95)$$

where B_{ρ} is defined in (72), and

$$i_t := \arg \max_{i \in [2[N]]} P \{ \rho_t \geq 2 C_{i,t} \mid F_t, g \}; \quad (96)$$

We denote p_t as an arbitrary point in $C_{i_t,t}$. Then, we define

$$A_i(t) := \{k S_i \rho_t \leq S_i p \leq k \frac{\rho}{8}\}; \quad (97)$$

Using these notations, the last term in (93) can be decomposed as

$$\begin{aligned} \sum_{t=1}^T \sum_{i(t)} h_i \mathbb{1}_{A_i^c(t)} &= \sum_{k=1}^N \sum_{i(t)} h_i \mathbb{1}_{\{p_t \in C_k; A_k^c(t)\}} \\ &= \sum_{k=1}^N \sum_{i(t)} h_i \mathbb{1}_{\{p_t \in C_k; A_k^c(t)\}} + \sum_{k=1}^N \sum_{i(t)} h_i \mathbb{1}_{\{p_t \in C_k; A_k(t); A_k^c(t)\}} \\ &= \underbrace{\sum_{k=1}^N \sum_{i(t)} h_i \mathbb{1}_{\{p_t \in C_k; A_k^c(t)\}}}_{(C)} + \underbrace{\sum_{k=1}^N \sum_{i(t)} h_i \mathbb{1}_{\{p_t \in C_k; A_k(t); A_k^c(t)\}}}_{(D)} + \underbrace{\sum_{k=2}^N \sum_{i(t)} h_i \mathbb{1}_{\{p_t \in C_k; A_k(t)\}}}_{(E)}; \end{aligned} \quad (98)$$

We will bound the expectation of each term in the following and complete the proof of Theorem 3 as

$$\begin{aligned}
 E[\text{Reg}(T)] &= \sum_{i \in \mathcal{I}} \left(O\left(\frac{1}{2} \log T\right) + O\left(\frac{N}{2} \log T\right) \right. \\
 &\quad \left. + \max_{j \in \mathcal{I}^{(N)}} \left(O\left(\frac{NM}{2} \log T\right) + O(1) + O(1) \right) \right) \\
 &= O\left(\max_{i \in \mathcal{I}^{(N)}} \left(\frac{N}{2} \log T; \frac{N^2 M \max_{i \in \mathcal{I}^{(N)}}}{2} \log T \right) \right) \\
 &= O\left(\frac{AN^2 M \max_{i \in \mathcal{I}^{(N)}}}{2} \log T \right); \tag{99}
 \end{aligned}$$

where the last transformation follows from the definition of (26).

F.2 Analysis for Case (A)

Lemma 18. For any $\epsilon > 0$,

$$E \sum_{t=1}^T \sum_{i \in \mathcal{I}} \left(\frac{64}{9} \log T + 2^{M-2} \right); \tag{100}$$

To prove Lemma 18, we prove the following lemma using Corollary 8 and Lemma 9.

Lemma 19. For any $0 < \epsilon < 1/2$,

$$P_{\mathbf{p}_t \in \mathcal{V}_i} \left(\sum_{j \in \mathcal{I}} A_j(t); N_i(t) > n_i \right) \leq \exp\left(-\frac{9}{16} n_i (1 - \epsilon)^{M-2}\right); \tag{101}$$

where $\mathcal{V}_i := \{\mathbf{p} \in \mathcal{C}_i : \|\mathbf{S}_i \mathbf{p} - \mathbf{S}_i \mathbf{p}_t\| \leq \epsilon\}$.

Proof. Since $\mathbf{p}_t \sim N(\mathbf{p}_t; \mathbf{B}_t^{-1})$ for $\mathbf{p}_t = \mathbf{B}_t^{-1} \mathbf{b}_t$, the squared Mahalanobis distance $\mathbf{B}_t^{-1} (\mathbf{p}_t - \mathbf{p}_t)^2$ follows the chi-squared distribution with M degree of freedom. Therefore, we have

$$P_{\mathbf{p}_t \in \mathcal{V}_i} \left(\sum_{j \in \mathcal{I}} A_j(t); N_i(t) > n_i \right) \leq \inf_{\mathbf{p} \in \mathcal{V}_i} \exp\left(-\frac{1}{2} (\mathbf{p} - \mathbf{p}_t)^2\right); \tag{102}$$

where $h(a) = P_{\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left(\sum_{i=1}^M X_i^2 \geq a \right)$. To use Lemma 9, we check the condition of Lemma 9 is indeed satisfied. First, it is obvious that the assumption $\sum_{i \in \mathcal{I}} n_i > 0$ and $\|\mathbf{S}_i \mathbf{p}_t - \mathbf{S}_i \mathbf{p}\| \leq \epsilon < 1/4$ are satisfied. Besides \mathcal{V}_i implies $\mathcal{T}_i = \{\mathbf{p} \in \mathcal{R}^M : \|\mathbf{S}_i \mathbf{p} - \mathbf{S}_i \mathbf{p}_t\| \leq \epsilon\}$ from Corollary 8. Thus, applying Lemma 9 concludes the proof. \square

Proof of Lemma 18 For any $n_i > 0$,

$$\begin{aligned}
 & \sum_{t=1}^T \sum_{i \in \mathcal{I}} \left(\frac{64}{9} \log T + 2^{M-2} \right) \\
 &= \sum_{t=1}^T \sum_{i \in \mathcal{I}} \left(\frac{64}{9} \log T + 2^{M-2} \right) \mathbb{1}_{\left(\sum_{j \in \mathcal{I}} A_j(t); N_i(t) > n_i \right)} + \sum_{t=1}^T \sum_{i \in \mathcal{I}} \left(\frac{64}{9} \log T + 2^{M-2} \right) \mathbb{1}_{\left(\sum_{j \in \mathcal{I}} A_j(t); N_i(t) \leq n_i \right)} \\
 &= \sum_{t=1}^T \sum_{i \in \mathcal{I}} \left(\frac{64}{9} \log T + 2^{M-2} \right) \mathbb{1}_{\left(\sum_{j \in \mathcal{I}} A_j(t); N_i(t) > n_i \right)}; \tag{103}
 \end{aligned}$$

The second term is bounded from above as

$$\begin{aligned}
 & \mathbb{E} \sum_{t=1}^T \mathbb{1}[i(t) = i; \mathcal{A}_1(t); A_i(t); N_i(t) > n_i] \\
 &= \mathbb{E} \sum_{t=1}^T \mathbb{P}[i(t) = i; \mathcal{A}_1(t); A_i(t); N_i(t) > n_i] \\
 &= \mathbb{E} \sum_{t=1}^T \mathbb{P}[i(t) = i; \mathcal{A}_1(t); A_i(t); N_i(t) > n_i \mid (i(t) = i \text{ implies } \mathcal{A}_1(t))] \\
 &= \mathbb{E} \sum_{t=1}^T \mathbb{P}[\mathcal{A}_1(t); N_i(t) > n_i]
 \end{aligned} \tag{104}$$

To obtain an upper bound for $\mathbb{P}[\mathcal{A}_1(t); N_i(t) > n_i]$, we use Lemma 19. By taking $\epsilon = \frac{16}{9} \frac{1}{2} \log T$ with $\delta = 1/4$, we have

$$\begin{aligned}
 & \mathbb{E} \sum_{t=1}^T \mathbb{1}[i(t) = i; \mathcal{A}_1(t); A_i(t)] \leq n_i + \sum_{t=1}^T \mathbb{P}[\mathcal{A}_1(t); N_i(t) > n_i] \\
 & \leq n_i + \sum_{t=1}^T \exp\left[-\frac{9}{16} \frac{2n_i}{(1-\delta)^{M=2}}\right] \quad (\text{by Lemma 19}) \\
 & = \frac{16}{9} \frac{1}{2} \log T + (1-\delta)^{M=2} \\
 & = \frac{64}{9} \log T + 2^{M=2};
 \end{aligned} \tag{105}$$

which completes the proof. \square

F.3 Analysis for Case (B)

Lemma 20. For any $i \in \{1, \dots, M\}$,

$$\mathbb{E} \sum_{t=1}^T \mathbb{1}[i(t) = i; A_i^c(t)] \leq \frac{256N \log T + \frac{A}{2} \log 2 + 1}{2} + \frac{16A^2}{2} \tag{106}$$

The regret in this case can intuitively be bounded because as the round proceeds the event makes $S_i(t)$ close to S_i^* , which implies that the expected number of times the event $A_i^c(t)$ occurs is not large.

Before going to the analysis of Lemma 20, we prove useful inequalities between $S_i(t)$, $k_i^{(t)}$, and S_i^* .

Lemma 21. Assume $N_i(t) > 0$. Then,

$$k_i^{(t)} \leq S_i^* k^2 \frac{Z_{n_i}}{N_i(t)} + k_i^{(t)} \leq S_i^* k^2 \tag{107}$$

Proof. Recall that \hat{p}_t is the maximizer of $g_t(p)$, and we have

$$\hat{p}_t = \arg \max_{p \in \mathcal{R}^M} g_t(p) = \arg \max_{p \in \mathcal{R}^M} \sum_{i=1}^M \exp\left[-\frac{1}{2} N_i(t) k_i^{(t)}\right] S_i p k^2 = \arg \min_{p \in \mathcal{R}^M} \sum_{i=1}^M N_i(t) k_i^{(t)} S_i p k^2 \tag{108}$$

Using this and the definition of Z_{ni} , we have

$$\begin{aligned} N_i(t)kq_i^{(t)} &= S_i p_t k^2 \sum_{k=1}^X N_k(t)kq_k^{(t)} = S_k p_t k^2 \\ &= S_i p k^2 \sum_{k=1}^{k_2^{[N]}} N_k(t)kq_k^{(t)} = S_k p k^2 \\ &= S_i p k^2 \sum_{k=1}^{k_2^{[N]}} N_k(t)kq_k^{(t)} = S_i p k^2 : \end{aligned} \quad (109)$$

Dividing by $N_i(t)$ on the both sides completes the proof. \square

Lemma 22. Assume that $A_i^c(t)$ and $N_i(t) > 0$ hold. Then,

$$kq_i^{(t)} = S_i p k > \frac{1}{2} \frac{S_i Z_{ni}}{N_i(t)} : \quad (110)$$

Proof. By the triangle inequality,

$$\begin{aligned} kq_i^{(t)} &= S_i p k - k S_i p_t \frac{S_i p k k q_i^{(t)}}{S_i p_t k} \\ &> \frac{1}{4} \frac{S_i Z_{ni}}{N_i(t) + kq_i^{(t)}} = S_i p k^2 \quad (\text{by } A_i^c(t) \text{ and Lemma 21}) \\ &> \frac{1}{4} \frac{S_i Z_{ni}}{N_i(t)} - k q_i^{(t)} = S_i p k \quad \text{by } \frac{1}{x+y} > \frac{1}{x} - \frac{1}{y} \text{ for } x, y > 0 ; \end{aligned} \quad (111)$$

which is equivalent to (110). \square

Proof of Lemma 20 We first bound the expectation conditioned Z_{ni} , and then take the expectation for Z_{ni} . Now,

$$\begin{aligned} &E \sum_{t=1}^T \mathbb{1}[i(t) = i; A_i^c(t)] Z_{ni} \\ &= E \sum_{t=1}^T \mathbb{1}[i(t) = i; A_i^c(t); N_i(t) > \frac{64Z_{ni}}{2}] Z_{ni} \\ &\quad + E \sum_{t=1}^T \mathbb{1}[i(t) = i; A_i^c(t); N_i(t) \leq \frac{64Z_{ni}}{2}] Z_{ni} \\ &= \frac{64Z_{ni}}{2} + E \sum_{t=1}^T \mathbb{1}[i(t) = i; A_i^c(t); N_i(t) > \frac{64Z_{ni}}{2}] Z_{ni} \quad (i(t) = i \text{ for all } t \in [T]) : \end{aligned} \quad (112)$$

The first term becomes $256N \log T + \frac{A}{2} \log 2 + 1 = 2$ by taking expectation over Z_{ni} using Lemma 10. Then, we bound the second term. From Lemma 21 and $N_i(t) > \frac{64Z_{ni}}{2}$ im-

ply $kq^{(t)}$ $S_i p k > = 16$. Therefore,

$$\begin{aligned}
& \mathbb{E} \prod_{t=1}^T \mathbb{1}_{i(t) = i; A_i^c(t); N_i(t) > \frac{64Z_{ni}}{2}} \quad \# \\
& \mathbb{E} \prod_{t=1}^T \mathbb{1}_{i(t) = i; kq^{(t)} S_i p k > \frac{1}{16}} \quad \# \\
& \mathbb{E} \prod_{t=1}^T \mathbb{1}_{i(t) = i; \left[\prod_{y \in [A]} j(q^{(t)})_y (S_i)_y p j > \frac{1}{16} \right]} \quad \# \\
& \mathbb{E} \prod_{y=1}^n \prod_{t=1}^T \mathbb{1}_{i(t) = i; j(q^{(t)})_y (S_i)_y p j > \frac{1}{16}} \quad \# \\
& \mathbb{E} \prod_{y=1}^n \prod_{n_i=1}^{n_i} \prod_{t=1}^T \mathbb{1}_{i(t) = i; N_i(t) = n_i; j(q^{(t)})_y (S_i)_y p j > \frac{1}{16}} \quad \#\# \\
& = \mathbb{E} \prod_{y=1}^n \prod_{n_i=1}^{n_i} \prod_{t=1}^T \mathbb{1}_{i(t) = i; N_i(t) = n_i; j(q^{(t)})_y (S_i)_y p j > \frac{1}{16}} \\
& \quad \text{(The event } i(t) = i; N_i(t) = n_i \text{ occurs at most once for each } i.) \\
& \prod_{y=1}^n \prod_{n_i=1}^{n_i} P \left[\prod_{t=1}^T j(q^{(t)})_y (S_i)_y p j > \frac{1}{4} \right] \\
& \prod_{y=1}^n \prod_{n_i=1}^{n_i} 2 \exp \left(-2n_i \frac{1}{4} \right) \quad \text{(by Hoeffding's inequality)} \\
& 2A \prod_{n_i=1}^n \exp \left(-\frac{n_i^2}{8A} \right) \\
& = 2A \frac{1}{\exp \left(\frac{1}{8A} \right)} \\
& 2A \frac{1}{\exp \left(\frac{1}{8A} \right)} \quad \text{(by } e^x \geq 1 + x \text{)} \\
& = \frac{16A^2}{2} : \tag{113}
\end{aligned}$$

By summing up the above argument, the proof is completed. \square

F.4 Analysis for Case (C)

Before going to the analysis of cases (C), (D), and (E), we recall some notations. Recall that

$$P_i(t) = P_{\{p_t \in C_i\}} \mathbb{1}_{F_t g}; \tag{114}$$

$C_{i,t} = C_i \setminus B_{\rho(p_t)}$, $i_t = \arg \max_{i \in [N]} P_{\{p_t \in C_{i,t}\}} \mathbb{1}_{F_t g}$, and p_t is an arbitrary point in $C_{i_t,t}$. Also recall that

$$A_i(t) = \sum_{p_t \in C_i} S_i p k \frac{1}{8} : \tag{115}$$

Lemma 23. For any $i \in [N]$,

$$\mathbb{E} \prod_{t=1}^T \mathbb{1}_{p_t \in C_i; A_i^c(t)} \leq \frac{N}{p_0} \frac{2^M \log T}{2} + e^{-kp} \frac{1}{k^2=2} + \frac{1}{T} + \frac{L}{M} \frac{1}{1 - e^{-\frac{1}{2^5}}} : \tag{116}$$

Before proving the above lemma, we give two lemmas.

Lemma 24.

$$P_{\rho_t} \leq C_{i,t} + F_t \quad P_{\rho_t} \leq C_{i,t} + F_t \quad p_0 = N; \quad (117)$$

where $p_0 := 1 - h(\epsilon^2)$.

Proof. First, we prove

$$P_{\rho_t} \leq \sum_{i \in [N]} C_{i,t} + F_t \leq 1 - h(\epsilon^2); \quad (118)$$

This follows from

$$\begin{aligned} P_{\rho_t} &\leq \sum_{i \in [N]} C_{i,t} + F_t = P_{\rho_t} \leq B_{\rho_t} + F_t \\ &\leq h \inf_{p: \rho_t \leq p} k B_t^{1-2} (p - \rho_t) k^2 \\ &\leq h k p - \rho_t k^2 \\ &\leq h(\epsilon^2); \end{aligned} \quad (119)$$

Using the definition of ρ_t completes the proof. \square

Lemma 25. For any $i \in [N]$, the event $A_i^c(t)$ implies $k S_i \rho_t - S_i p k \leq 16$.

Proof. Using the triangle inequality, we have

$$\begin{aligned} k S_i \rho_t - S_i p k &\leq k S_i \rho_t - S_i p k + k S_i \rho_t - S_i p k \\ &\leq 8 k S_i k \rho_t - p k \\ &\leq 8 k S_i k \frac{1}{16 \max k S_i k} \\ &\leq 8 \cdot 16 = 128; \end{aligned} \quad (120)$$

Proof of Lemma 23 For any n_0 , which is specified later, we have

$$\begin{aligned} &E \sum_{t=1}^T \sum_{i \in [N]} 1_{\rho_t \leq C_i; A_i^c(t)} \\ &= E \sum_{t=1}^T \sum_{i \in [N]} 1_{\rho_t \leq C_i; A_i^c(t); N_i(t) < n_0} + E \sum_{t=1}^T \sum_{i \in [N]} 1_{\rho_t \leq C_i; A_i^c(t); N_i(t) \geq n_0} \end{aligned} \quad (121)$$

The first term can be bounded by $(p_0 = N)^{-1} n_0$ from Lemma 24. The rigorous proof can be obtained by the almost same argument as the following analysis of the second term using Theorem 16.

Then, we will bound the second term. Specifically, we will prove that $\frac{M \log T}{(\epsilon = 16)^2}$,

$$E \sum_{t=1}^T \sum_{i \in [N]} 1_{\rho_t \leq C_i; A_i^c(t); N_i(t) \geq n_0} = O(1); \quad (122)$$

First we have

$$\begin{aligned} &E \sum_{t=1}^T \sum_{i \in [N]} 1_{\rho_t \leq C_i; A_i^c(t); N_i(t) \geq n_0} \\ &= \sum_{m=n_0}^{\infty} E \sum_{t=1}^T \sum_{i \in [N]} 1_{\rho_t \leq C_i; A_i^c(t); N_i(t) = m}; \end{aligned} \quad (123)$$

Let

$$= \min_{t: p_t \geq C_i; A_i^c(t); N_i(t) = m} (T+1) \quad (124)$$

be the first time such that $p_t \geq C_i; A_i^c(t)$ and $N_i(t) = m$ occur. Letting $A_t := p_t \geq C_i; A_i^c(t); N_i(t) = m$, $B_t := f_i(t) = ig$ and $P_t := p_0 = N$ in Theorem 16, we have

$$E \sum_{t=1}^T \frac{X_t}{p_0} \mathbb{1}_{p_t \geq C_i; A_i^c(t); N_i(t) = m} \leq \frac{N}{p_0} P_f \quad Tg: \quad (125)$$

Here T implies that

$$\begin{aligned} k\hat{p} - p k_B &= (\hat{p} - p) \sum_{j \in [N]} N_j(S_j^> S_j) A(\hat{p} - p) \\ &= m \sum_{j \in [N]} S_j^> S_j (\hat{p} - p) \\ &= m k_{S_i} (\hat{p} - p)^2 = m (=16)^2; \end{aligned} \quad (126)$$

where the last inequality follows from Lemma 25. Therefore we have

$$E \exp(k\hat{p} - p k_B = 2) \leq E \mathbb{1}_{[T]} \exp(k\hat{p} - p k_B = 2) \exp(m (=16)^2 = 2) P_f \quad Tg: \quad (127)$$

Note that $\sum_{j \in [N]} B_j \leq (1 + TL = M)^M$ for $L = \max_i \text{trace}(S_i^> S_i) = \max_i k_{S_i} k_F$ by Lemma 10 of [Abbasi-Yadkori et al. \(2011\)](#), where k_F is the Frobenius norm. Therefore we have

$$\begin{aligned} E[\exp(k\hat{p} - p k_B = 2)] &\leq E \sum_{j \in [N]} \frac{\exp(k\hat{p} - p k_B = 2)}{B_j} \\ &\leq (1 + TL = M)^{M=2} E \frac{\exp(k\hat{p} - p k_B = 2)}{B_j} \\ &\leq (1 + TL = M)^{M=2} E \frac{\exp(k\hat{p}_0 - p k_{B_0} = 2)}{B_{0j}} \end{aligned} \quad (128)$$

$$= 1 + \frac{TL}{M} e^{kp k^2 = 2}; \quad (129)$$

where (128) holds since $\frac{\exp(k\hat{p} - p k_B = 2)}{B_j}$ is a supermartingale from Lemma 12. Combining (125), (127), and (129), we obtain

$$\begin{aligned} E \sum_{t=1}^T \frac{X_t}{p_0} \mathbb{1}_{p_t \geq C_i; A_i^c(t); N_i(t) = m} &\leq \frac{N}{p_0} \left(1 + \frac{TL}{M} e^{kp k^2 = 2} \right)^{M=2} E \sum_{m=n_0}^T e^{m (=16)^2 = 2} \\ &\leq \frac{N}{p_0} \left(1 + \frac{TL}{M} e^{kp k^2 = 2} \right)^{M=2} \frac{e^{n_0^2 = 2}}{1 - e^{(-16)^2 = 2}}; \end{aligned} \quad (130)$$

By choosing $n_0 = \frac{M \log T}{(-16)^2}$ we obtain the lemma. \square

F.5 Analysis for Case (D)

Lemma 26. For any $i \in [N]$,

$$E \sum_{t=1}^h \frac{X_t}{p_0} \mathbb{1}_{p_t \geq C_i; A_i(t); A_i^c(t)} \leq \frac{48M + 2N}{9} \frac{1}{p_0}; \quad (131)$$

Remark. To prove the regret upper bound, it is enough to prove Lemma 26 only for $i=1$. However, for the sake of generality, we prove the lemma for any $i \in [N]$.

Before proving Lemma 26, we give two following lemmas.

Lemma 27. For any $i \in [N]$, the event $A_i(t)$ implies $A_i^c(t)$.

Proof. Using the triangle inequality, we have

$$\begin{aligned} & \|S_i(t) - S_i^c(t)\|_k \leq \|S_i(t) - S_i^c(t)\|_k + \|S_i^c(t) - S_i^c(t)\|_k \\ & = 8 + \|S_i(t)\|_k \frac{\|S_i^c(t) - S_i^c(t)\|_k}{16 \max_k \|S_i(t)\|_k} \leq 4; \end{aligned} \quad (132)$$

which completes the proof. \square

Now, Lemma 26 can be intuitively proven because from Lemma 27, $A_i(t)$ implies $A_i^c(t)$, and the events $A_i(t)$ and $A_i^c(t)$ does not simultaneously occur many times.

Let t_1, \dots, t_m be the time of the r th m times that the event $p_t \in C_i; A_i(t); N_i(t) = n_i$ occurred (not $p_t \in C_i; A_i(t); N_i(t) = n_i$). In other words, we define

- t_1 : the first time that $p_t \in C_i; A_i(t)$ and $N_i(t) = n_i$ occurred
- t_2 : the second time that $p_t \in C_i; A_i(t)$ and $N_i(t) = n_i$ occurred
- ...

Now we prove the following lemma using Lemma 9.

Lemma 28. For any $0 < \epsilon \leq 1/2$,

$$P \left[\bigcap_{i=1}^n A_i^c(t); \bigcap_{k=1}^m t_k \leq \epsilon \right] \leq \exp \left[-\frac{9}{16} \sum_{i=1}^n n_i (1 - \epsilon)^{M=2} \right]; \quad (133)$$

Proof. Recall that $\mathbb{T}_i = \{p \in \mathbb{R}^M : \|S_i(p) - S_i^c(p)\|_k > \epsilon\}$. We follow a similar argument as the analysis for Lemma 19. Since $p_t \sim N(B_t^{-1}b; B_t^{-1})$, the squared Mahalanobis distance $k_{B_t}^{1=2}(p - p_t)_k^2$ follows the chi-squared distribution with M degree of freedom. Hence, $\mathbb{P}(a) = P_{X \sim \frac{1}{M} \chi^2_M}$ ag, we have

$$P \left[\bigcap_{i=1}^n A_i^c(t); \bigcap_{k=1}^m t_k \leq \epsilon \right] \leq \inf_{p \in \mathbb{T}_i} P_{B_t}^{1=2}(p - p_t)_k^2; \quad (134)$$

Then, Eq. (133) directly follows from Lemma 9. \square

Proof of Lemma 26 From Lemma 27, the event $A_i(t)$ implies $A_i^c(t)$. Hence, it is enough to derive the upper bound for

$$E \left[\sum_{t=1}^n \mathbb{1}_{\{p_t \in C_i; A_i(t); A_i^c(t)\}} \right] \quad (135)$$

instead of the bound for

$$E \left[\sum_{t=1}^n \mathbb{1}_{\{p_t \in C_i; A_i(t); A_i^c(t)\}} \right]; \quad (136)$$

Proof. We evaluate each term in the summation using Theorem 16 with

$$\begin{aligned} A_t &= f_{p_t} 2 C_1; k S_i(p_t - p) k = 8; N_1(t) = n g; \\ B_t &= f_{p_t} 2 C_1 g; \end{aligned} \tag{140}$$

for $n \geq [T]$. Recall that

$$g_t(p) = \frac{1}{(2^M j B_t^{-1} j)} \exp \left(\frac{1}{2} k p - p k_{B_t}^2 \right) \tag{141}$$

is the probability density function of $F_t = f_{B_t}; b, g$. Using the definition in (80), it holds for any $n \geq [T]$ that

$$\begin{aligned} P_{f_{B_t} j F_t} g &= \int_{\mathbb{R}^2} f_{p_t} 2 C_1 j F_t g \\ &= \int_{\mathbb{R}^{2 C_1}} g(p) dp \\ &= \int_{p: k p - p k \geq 3} g(p) dp \\ &= \sup_{p: k p - p k \geq 2} \int_{p^0: k p^0 - p k} g(p^0) dp^0 \\ &= \sup_{p: k p - p k \geq 2} \inf_{p^0: k p^0 - p k} g(p^0) \text{Vol}(f_{p^0}: k p^0 - p k) \\ &= \frac{(p - \bar{p})^M}{(M=2+1)} \sup_{p: k p - p k \geq 2} \inf_{p^0: k p^0 - p k} g(p^0) \\ &= \frac{(p - \bar{p})^M j B_t j}{(M=2+1)} \exp \left(\frac{1}{2} \inf_{p: k p - p k \geq 2} \sup_{p^0: k p^0 - p k} k p^0 - p k_{B_t}^2 \right) \\ &= \frac{(p - \bar{p})^M j B_t j}{(M=2+1)} \exp \left(\frac{k p - p k_{B_t}^2}{2} - \frac{p - \bar{p}}{n} \right); \end{aligned} \tag{142}$$

where (142) follows since $p : k p - p k \geq 3 g f_{p^0}: k p^0 - p k g$ for any p_0 such that $k p_0 - p k \geq 2$, and the last inequality follows from Theorem 15. To apply Theorem 15, we used Lemma 27.

Now we define a stochastic process corresponds to (143) as

$$P_t = \frac{(p - \bar{p})^M j B_t j}{(M=2+1)} \exp \left(\frac{k p_t - p k_{B_t}^2}{2} - \frac{p - \bar{p}}{n} \right); \tag{144}$$

Then, by Lemma 12,

$$\begin{aligned} E[P_{t+1}^1 j F_t] &= \frac{(M=2+1)}{(p - \bar{p})^M} e^{-\frac{p - \bar{p}}{n}} E \left[\frac{1}{j B_{t+1} j} \exp \left(\frac{k p_t - p k_{B_{t+1}}^2}{2} \right) F_t; S_{i(t)} \right] \\ &= \frac{(M=2+1)}{(p - \bar{p})^M} e^{-\frac{p - \bar{p}}{n}} E \left[\frac{1}{j B_t j} \exp \left(\frac{k p_t - p k_{B_t}^2}{2} \right) F_t \right] \\ &= P_t^1; \end{aligned} \tag{145}$$

which means that P_t^1 is a supermartingale. Therefore we can apply Theorem 16 and obtain

$$\begin{aligned} E \left[\sum_{t=1}^X 1_{[A_t \geq C_1; k S_i(p_t - p) k = 8; N_1(t) = n]} \right] &= E \left[\sum_{t=1}^X P_t^1 \right] \\ &= E \left[P_1^1 \right] \\ &= \frac{(M=2+1) e^{-\frac{2 k p - k^2 = 2}{p - \bar{p}}}}{(p - \bar{p})^M} e^{-\frac{p - \bar{p}}{n}}; \end{aligned} \tag{146}$$

Finally we have

$$\begin{aligned}
 & \int_0^1 \int_0^1 \dots \int_0^1 \prod_{t=1}^n [p_t \geq c_t; k S_t(p_t - p)] \, dx_1 \dots dx_n \\
 &= \int_0^1 \int_0^1 \dots \int_0^1 \prod_{t=1}^n [p_t \geq c_t; k S_t(p_t - p)] \, dx_1 \dots dx_n \\
 &= \frac{(M+2+1)e^{-2kp} p^{k^2-2}}{(p-2)^M} \int_0^1 \int_0^1 \dots \int_0^1 \prod_{t=1}^n [p_t \geq c_t; k S_t(p_t - p)] \, dx_1 \dots dx_n \\
 &= \frac{(M+2+1)e^{-2kp} p^{k^2-2}}{(p-2)^M} \frac{2}{(p-2)^2} \int_0^1 \int_0^1 \dots \int_0^1 \prod_{t=1}^n [p_t \geq c_t; k S_t(p_t - p)] \, dx_1 \dots dx_n \\
 &= \frac{2^{M+2+3} (M+2+1)e^{-2kp} p^{k^2-2}}{2^{M+2} (M+2)^{M+1}}; \tag{147}
 \end{aligned}$$

which completes the proof. \square

G Property of Dynamic Pricing Games

In this appendix, we will see a property of dp-easy games.

Proposition 30. Consider any dp-easy games with $c > 1$. Then, any two actions in the game are neighbors.

Remark. In section 5, we considered dp-easy games with $c > 1$ to prove Proposition 30.

Proof. Take any two different actions $j, k \in [N]$ such that $j < k$. From the definition of the loss matrix in dp-easy games, we have $e_j \in C_j$ and $e_k \in C_k$.

First, we will find $\alpha \in [0, 1]$ such that

$$e_j + (1 - \alpha)e_k \in C_j \setminus C_k; \tag{148}$$

From the definition of the loss matrix, the element of $L(e_j + (1 - \alpha)e_k) \in P_M$ is

$$\begin{aligned}
 & \frac{\alpha}{c} < i < j < k < N \\
 & : \frac{c + (1 - \alpha)(i - j + 1)}{c} < i < j + 1 < i < k < N \\
 & : \frac{c + (1 - \alpha)(k - i - N)}{c} < i < j + 1 < i < k < N
 \end{aligned} \tag{149}$$

It is easy to see that the indices which give the minimum value (149) is j or k . Thus, to achieve the condition (148), the following should be satisfied,

$$j = \frac{c + (1 - \alpha)(k - i - N)}{c}; \tag{150}$$

which is equivalent to

$$\alpha = \frac{k - j}{c + k} (=: \beta); \tag{151}$$

Note that we have $\beta < 1$ for any $c > 1$.

Next, we introduce the following definitions.

$$p^{(j;k)} := \sum_n e_j + (1 - \beta)e_k \in C_j \setminus C_k; \tag{152}$$

$$\text{Ball}^{(j;k)} := \{p \in P_M : \|p - p^{(j;k)}\| \leq \beta\}; \tag{153}$$

$$L^{(x)} := L(p^{(j;k)} + x) \in \mathbb{R}^N; \tag{154}$$

To prove the proposition, it is enough to prove the following: there exists ϵ_0 , $\text{Ball}^{(j;k)} \subset C_j \cap C_k$.

To prove this, it is enough to prove that, there exists ϵ_0 ,

$$\min_{x \in \mathbb{R}^M} \min_{i \in [N]} \min_{j \in [N]} \min_{k \in [N]} (L^{(x)})_i - (L^{(x)})_j - (L^{(x)})_i - (L^{(x)})_k > 0: \quad (155)$$

We will prove (155) in the following. Take any $x \in \mathbb{R}^M$ and

$$:= \min_{i:1 \leq i \leq j} \frac{1}{2} \frac{j-i}{kL_j - L_i k} \wedge \min_{i:j < i < k} \frac{1}{2} \frac{(1-\epsilon)(k-i)}{kL_i - L_k k} \wedge \min_{i:k < i \leq N} \frac{1}{2} \frac{c+j}{kL_j - L_i k}: \quad (156)$$

Note that the used here is different from the one used in the proof of the regret upper bounds.

Case (A): When $1 \leq i < j$, using Cauchy Schwarz inequality, we have

$$\begin{aligned} (L^{(x)})_i - (L^{(x)})_j - (L^{(x)})_i - (L^{(x)})_k &= (L^{(x)})_i - (L^{(x)})_j \\ &= (j-i) \frac{1}{2} \frac{j-i}{kL_j - L_i k} \\ &= (j-i) \frac{1}{2} \frac{(j-i)}{kL_j - L_i k} \\ &> 0: \end{aligned} \quad (157)$$

The arguments for cases (B) and (C) follow in the similar manner as case (A).

Case (B): When $j < i < k$, we have

$$\begin{aligned} (L^{(x)})_i - (L^{(x)})_j - (L^{(x)})_i - (L^{(x)})_k &= (L^{(x)})_i - (L^{(x)})_k \\ &= (k-i) \frac{1}{2} \frac{(1-\epsilon)(k-i)}{kL_i - L_k k} \\ &= (1-\epsilon) \frac{1}{2} (k-i) \\ &> 0: \end{aligned} \quad (158)$$

Case (C) When $k < i \leq N$, we have

$$\begin{aligned} (L^{(x)})_i - (L^{(x)})_j - (L^{(x)})_i - (L^{(x)})_k &= (L^{(x)})_i - (L^{(x)})_j \\ &= (c+j) \frac{1}{2} \frac{c+j}{kL_j - L_i k} \\ &= (c+j) \frac{1}{2} \\ &> 0: \end{aligned} \quad (159)$$

Summing up the argument for cases (A) to (C), the proof is completed. \square

H Details and Additional Results of Experiments

Here we give the specific values of the opponent's strategy used in Section 5 and show the extended experimental results for performance comparison. Table 2 summarizes the values of opponent's strategy used in this appendix and Section 5. Figure 3 shows the empirical comparison of the proposed algorithms against the benchmark methods, and Figure 4 shows the number of the rejected times. We can see the same tendency as Section 5, that is, TSPM performs the best and the number of rejections does not increase with the time step.

Table 2: The values of the opponent's strategy.

# of outcomes M	opponent's strategy p
2	[0.7; 0.3]
3	[0.5; 0.3; 0.2]
4	[0.3; 0.3; 0.3; 0.1]
5	[0.2; 0.3; 0.3; 0.1; 0.1]
6	[0.2; 0.2; 0.3; 0.1; 0.1; 0.1]
7	[0.2; 0.2; 0.3; 0.1; 0.1; 0.05; 0.05]

(a) dp-easy, $N = M = 2$

(b) dp-easy, $N = M = 3$

(c) dp-easy, $N = M = 4$

(d) dp-easy, $N = M = 5$

(e) dp-easy, $N = M = 6$

(f) dp-easy, $N = M = 7$

(g) dp-hard, $N = M = 2$

(h) dp-hard, $N = M = 3$

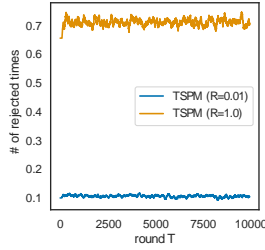
(i) dp-hard, $N = M = 4$

(j) dp-hard, $N = M = 5$

(k) dp-hard, $N = M = 6$

(l) dp-hard, $N = M = 7$

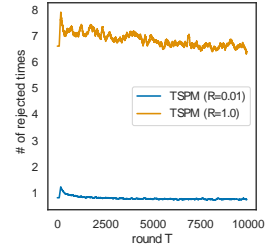
Figure 3: Regret-round plots of the algorithms. The solid lines indicate the average over 100 independent trials. The thin fillings are the standard error.



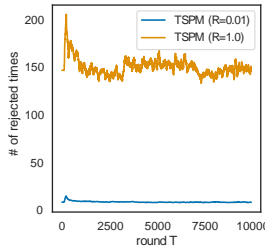
(a) dp-easy, $N = M = 2$



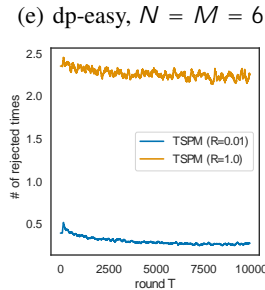
(b) dp-easy, $N = M = 3$



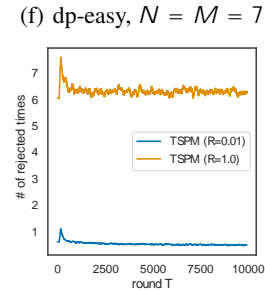
(c) dp-easy, $N = M = 4$



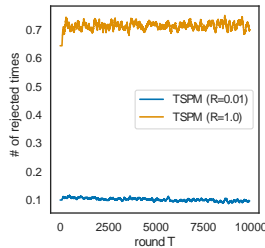
(d) dp-easy, $N = M = 5$



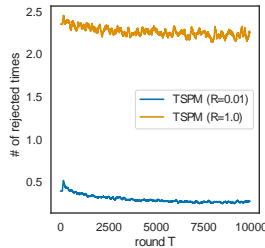
(e) dp-easy, $N = M = 6$



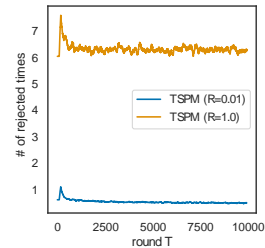
(f) dp-easy, $N = M = 7$



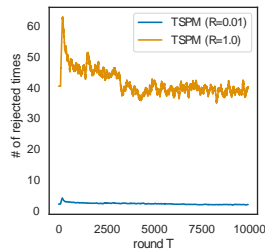
(g) dp-hard, $N = M = 2$



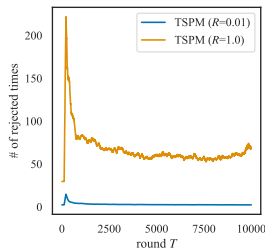
(h) dp-hard, $N = M = 3$



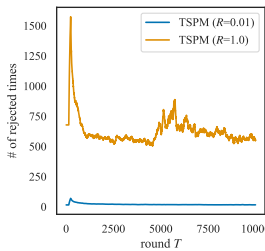
(i) dp-hard, $N = M = 4$



(j) dp-hard, $N = M = 5$



(k) dp-hard, $N = M = 6$



(l) dp-hard, $N = M = 7$

Figure 4: The number of rejected times by the accept-reject sampling. The solid lines indicate the average over 100 independent trials after taking moving average with window size 100.