
Efficient Partial Monitoring with Prior Information

Hastagiri P Vanchinathan

Dept. of Computer Science
ETH Zürich, Switzerland
hastagiri@inf.ethz.ch

Gábor Bartók

Dept. of Computer Science
ETH Zürich, Switzerland
bartok@inf.ethz.ch

Andreas Krause

Dept. of Computer Science
ETH Zürich, Switzerland
krausea@ethz.ch

Abstract

Partial monitoring is a general model for online learning with limited feedback: a learner chooses actions in a sequential manner while an opponent chooses outcomes. In every round, the learner suffers some loss and receives some feedback based on the action and the outcome. The goal of the learner is to minimize her cumulative loss. Applications range from dynamic pricing to label-efficient prediction to dueling bandits. In this paper, we assume that we are given some prior information about the distribution based on which the opponent generates the outcomes. We propose BPM, a family of new efficient algorithms whose core is to track the outcome distribution with an ellipsoid centered around the estimated distribution. We show that our algorithm provably enjoys near-optimal regret rate for locally observable partial-monitoring problems against stochastic opponents. As demonstrated with experiments on synthetic as well as real-world data, the algorithm outperforms previous approaches, even for very uninformed priors, with an order of magnitude smaller regret and lower running time.

1 Introduction

We consider *Partial Monitoring*, a repeated game where in every time step a learner chooses an action while, simultaneously, an opponent chooses an outcome. Then the player receives a loss based on the action and outcome chosen. The learner also receives some feedback based on which she can make better decisions in subsequent time steps. The goal of the learner is to minimize her cumulative loss over some time horizon.

The performance of the learner is measured by the *regret*, the excess cumulative loss of the learner compared to that of the best fixed constant action. If the regret scales linearly with the time horizon, it means that the learner does not approach the performance of the best action, that is, the learner fails to learn the problem. On the other hand, sublinear regret indicates that the disadvantage of the learner compared to the best fixed strategy fades with time.

Games in which the learner receives the outcome as feedback after every time step are called *online learning with full information*. This special case of partial monitoring has been addressed by (among others) Vovk [1] and Littlestone and Warmuth [2], who designed the randomized algorithm Exponentially Weighted Averages (EWA) as a learner strategy. This algorithm achieves $\Theta(\sqrt{T \log N})$ expected regret against any opponent, where N is the number of actions and T is the time horizon. This regret growth rate is also proven to be optimal.

Another well-studied special case is the so-called *multi-armed bandit problem*. In this feedback model, the learner gets to observe the loss she suffered in every time step. That is, the learner does not receive any information about losses of actions she did not choose. Asymptotically optimal results were obtained by Audibert and Bubeck [3], who designed the Implicitly Normalized Forecaster (INF) that achieves the minimax optimal $\Theta(\sqrt{TN})$ regret growth rate.¹

¹The algorithm Exp3 due to Auer et al. [4] achieves the same rate up to a logarithmic factor.

However, not all online learning problems have one of the above feedback structures. An important example for a problem that does not fit in either full-information or bandit problems is *dynamic pricing*. Consider the problem of a vendor wanting to sell his products to customers for the best possible price. When a customer comes in, she (secretly) decides on a maximum price she is willing to buy his product for, while the vendor has to set a price without knowing the customer’s preferences. The loss of the vendor is some preset constant if the customer did not buy the product, and an “opportunity loss”, when the product was sold cheaper than the customer’s maximum. The feedback, on the other hand, is merely an indicator whether the transaction happened or not.

Dynamic pricing is just one of the practical applications of partial monitoring. *Label efficient prediction*, in its simplest form, has three actions: the first two actions are guesses of a binary outcome but provide no information, while the third action provides information about the outcome for some unit loss as the price. This can be thought of an abstract form of *spam filtering*: the first two actions correspond to putting an email to the inbox and the spam folder, the third action corresponds to asking the user if the email is spam or not. Another problem that can be cast as partial monitoring is that of *dueling bandits* [5, 6] in which the learner chooses a pair of actions in every time step, the loss she suffers is the average loss of the two actions, and the feedback is which action was “better”.

In online learning, we distinguish different models of how the opponent generates the outcomes. In the mildest version called *stochastic* or *stationary memoryless*, the opponent chooses an outcome distribution before the game starts and then selects outcomes in an iid random manner drawn from the chosen distribution. The *oblivious adversarial* opponent chooses the outcomes arbitrarily, but without observing the actions of the learner. This selection method is equivalent to choosing an outcome sequence ahead of time. Finally, the *non-oblivious* or *adaptive adversarial* opponent chooses outcomes arbitrarily with the possibility of looking at past actions of the learner. In this work, we focus on strategies against stochastic opponents.

Related work Partial monitoring was first addressed in the seminal paper of Piccolboni and Schindelhauer [7], who designed and analyzed the algorithm FeedExp3. The algorithm’s main idea is to maintain an unbiased estimate for the loss of each action in every time step, and then use these estimates to run the full-information algorithm (EWA). Piccolboni and Schindelhauer [7] proved an $O(T^{3/4})$ upper bound on the regret (not taking into account the number of actions) for games for which learning is at all possible. This bound was later improved by Cesa-Bianchi et al. [8] to $O(T^{2/3})$, who also constructed an example of a problem for which this bound is optimal.

From the above bounds it can be seen that not all partial-monitoring problems have the same level of difficulty: while bandit problems enjoy an $O(\sqrt{T})$ regret rate, some partial-monitoring problems have $\Omega(T^{2/3})$ regret. To this end, Bartók et al. [9] showed that partial-monitoring problems with finitely many actions and outcomes can be classified into four groups: *trivial* with zero regret, *easy* with $\tilde{\Theta}(\sqrt{T})$ regret, *hard* with $\Theta(T^{2/3})$ regret, and *hopeless* with linear regret. The distinguishing feature between easy and hard problems is the *local observability condition*, an algebraic condition on the feedback structure that can be efficiently verified for any problem. Bartók et al. [9] showed the above classification against stochastic opponents with the help of algorithm BALATON. This algorithm keeps track of estimates of the loss difference of “neighboring” action pairs and eliminates actions that are highly likely to be suboptimal.

Since then, several algorithms have been proposed that achieve the $\tilde{O}(\sqrt{T})$ regret bound for easy games [10, 11]. All these algorithms rely on the core idea of estimating the expected loss difference between pairs of actions.

Our contributions In this paper, we introduce BPM (Bayes-update Partial Monitoring), a new family of algorithms against iid stochastic opponents that rely on a novel way of the usage of past observations. Our algorithms maintain a confidence ellipsoid in the space of outcome distributions, and update the ellipsoid based on observations following a Bayes-like update. Our approach enjoys better empirical performance and lower computational overhead; another crucial advantage is that we can incorporate prior information about the outcome distribution by means of an initial confidence ellipsoid. We prove near-optimal minimax expected regret bounds for our algorithm, and demonstrate its effectiveness on several partial monitoring problems on synthetic and real data.

2 Problem setup

Partial monitoring is a repeated game where in every round, a learner chooses an action while the opponent chooses an outcome from some finite action and outcome sets. Then, the learner observes a feedback signal (from some given set of symbols) and suffers some loss, both of which are deterministic functions of the action and outcome chosen. In this paper we assume that the opponent chooses the outcomes in an iid stochastic manner. The goal of the learner is to minimize her cumulative loss.

The following definitions and concepts are mostly taken from Bartók et al. [9]. An instance of partial monitoring is defined by the *loss matrix* $L \in \mathbb{R}^{N \times M}$ and the *feedback table* $H \in \Sigma^{N \times M}$, where N and M are the cardinality of the action set and the outcome set, respectively, while Σ is some alphabet of symbols. That is, if learner chooses action i while the outcome is j , the loss suffered by the learner is $L[i, j]$, and the feedback received is $H[i, j]$.

For an action $1 \leq i \leq N$, let ℓ_i denote the column vector given by the i^{th} row of L . Let Δ_M denote the M -dimensional probability simplex. It is easy to see that for any $p \in \Delta_M$, if we assume that the opponent uses p to draw the outcomes (that is, p is the *opponent strategy*), the expected loss of action i can be expressed as $\ell_i^\top p$.

We measure the performance of an algorithm with its *expected regret*, defined as the expected difference of the cumulative loss of the algorithm and that of the best fixed action in hindsight:

$$R_T = \max_{1 \leq i \leq N} \sum_{t=1}^T (\ell_{I_t} - \ell_i)^\top p,$$

where T is some time horizon, I_t ($t = 1, \dots, T$) is the action chosen in time step t , and p is the outcome distribution the opponent uses.

In this paper, we also assume we have some prior knowledge about the outcome distribution in the form of a *confidence ellipsoid*: we are given a distribution $p_0 \in \Delta_M$ and a symmetric positive semidefinite covariance matrix $\Sigma_0 \in \mathbb{R}^{M \times M}$ such that the true outcome distribution p^* satisfies

$$\|p_0 - p^*\|_{\Sigma_0^{-1}} = \sqrt{(p_0 - p^*)^\top \Sigma_0^{-1} (p_0 - p^*)} \leq 1.$$

We use the term “confidence ellipsoid” even though our condition is not probabilistic; we do not assume that p^* is drawn from a Gaussian distribution before the game starts. On the other hand, the way we track p^* is derived by Bayes updates with a Gaussian conjugate prior, hence the name. We would also like to note that having the above prior knowledge is without loss of generality. For “large enough” Σ_0 , the whole probability simplex is contained in the confidence ellipsoid and thus partial monitoring without any prior information reduces to our setting.

The following definition reveals how we use the loss matrix to recover the structure of a game.

Definition 1 (Cell decomposition, Bartók et al. [9, Definition 2]). *For any action $1 \leq i \leq N$, let \mathcal{C}_i denote the set of opponent strategies for which action i is optimal:*

$$\mathcal{C}_i = \{p \in \Delta_M : \forall 1 \leq j \leq N, (\ell_i - \ell_j)^\top p \leq 0\}.$$

We call the set \mathcal{C}_i the optimality cell of action i . Furthermore, we call the set of optimality cells $\{\mathcal{C}_1, \dots, \mathcal{C}_N\}$ the cell decomposition of the game.

Every cell \mathcal{C}_i is a convex closed polytope, as it is defined by a linear inequality system. Normally, a cell has dimension $M - 1$, which is the same as the dimensionality of the probability simplex. It might happen however, that a cell is of lower dimensionality. Another possible degeneracy is when two actions share the same cell. In this paper, for ease of presentation, we assume that these degeneracies do not appear. For an illustration of cell decomposition, see Figure 1(a).

Now that we know the regions of optimality, we can define when two actions are *neighbors*. Intuitively, two actions are neighbors if their optimality cells are neighbors in the strong sense that they not only meet in “one corner”.

Definition 2 (Neighbors, Bartók et al. [9, page 4]). *Two actions i and j are neighbors, if the intersection of their optimality cells $\mathcal{C}_i \cap \mathcal{C}_j$ is an $M - 2$ -dimensional convex polytope.*

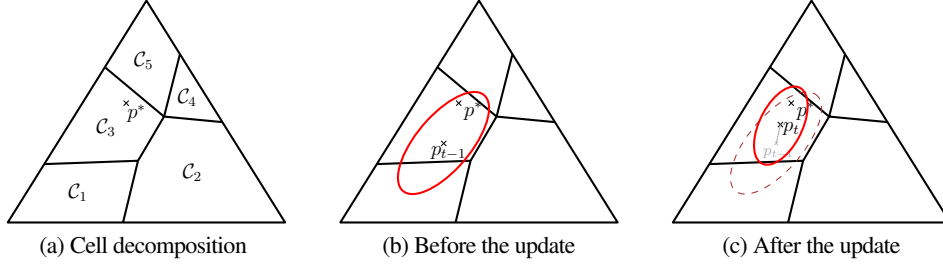


Figure 1: (a) An example for a cell decomposition with $M = 3$ outcomes. Under the true outcome distribution p^* , action 3 is optimal. Cells \mathcal{C}_1 and \mathcal{C}_3 are neighbors, but \mathcal{C}_2 and \mathcal{C}_5 are not. (b) The current estimate p_{t-1} is far away from the true distribution, the confidence ellipsoid is large. (c) After updating, p_t is closer to the truth, the confidence ellipsoid shrinks.

To optimize performance, the learner’s primary goal is to find out which cell the opponent strategy lies in. Then, the learner can choose the action associated with that cell to play optimally. Since the feedback the learner receives is limited, this task of finding the optimal cell may be challenging.

The next definition enables us to utilize the feedback table H .

Definition 3 (Signal matrix, Bartók et al. [9, Definition 1]). *Let $\{\alpha_1, \alpha_2, \dots, \alpha_{\sigma_i}\} \subseteq \Sigma$ be the set of symbols appearing in row i of the feedback table H . We define the signal matrix $S_i \in \{0, 1\}^{\sigma_i \times M}$ of action i as*

$$S_i[k, j] = \mathbb{I}(H[i, j] = \alpha_k).$$

In words, S_i is the indicator table of observing symbols $\alpha_1, \dots, \alpha_{\sigma_i}$ under outcomes $1, \dots, M$ given that the action chosen is i . For an example, consider the case when the i^{th} row of H is $(a \ b \ a \ c)$. Then,

$$S_i = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

A very useful property of the signal matrix is that if we represent outcomes with M -dimensional unit vectors, then S_i can be used as a linear transformation to arrive at the unit-vector representation of the observation.

The following condition is key in distinguishing easy and hard games:

Definition 4 (Local observability, Bartók et al. [9, Definition 3]). *Let actions i and j be neighbors. These actions are said to be locally observable if $\ell_i - \ell_j \in \text{Im } S_i^\top \oplus \text{Im } S_j^\top$. Furthermore, a game is locally observable if all of its neighboring action pairs are locally observable.*

Bartók et al. [9] showed that finite stochastic partial-monitoring problems that admit local observability have $\tilde{\Theta}(\sqrt{T})$ minimax expected regret. In the following, we present our new algorithm family that achieves the same regret rate for locally observable games against stochastic opponents.

3 BPM: New algorithms for Partial Monitoring based on Bayes updates

The algorithms we propose can be decomposed into two main building blocks: the first one keeps track of a belief about the true outcome distribution and provides us with a set of *feasible* actions in every round. The second one is responsible for selecting the action to play from this action set. Pseudocode for the algorithm family is shown in Algorithm 1.

3.1 Update Rule

The method of updating the belief about the true outcome distribution (p^*) is based on the idea that we pretend that the outcomes are generated from a Gaussian distribution with covariance $\Sigma = I_M$ and unknown mean. We also pretend we have a Gaussian prior for tracking the mean. The parameters of this prior are denoted by p_0 (mean) and Σ_0 (covariance). In every time step, we perform a Gaussian Bayes-update using the observation received.

Algorithm 1 BPM

input: L, H, p_0, Σ_0 **initialization:** Calculate signal matrices S_i **for** $t = 1$ **to** T **do** Use selection rule (cf., Sec. 3.2) to choose an action I_t Observe feedback Y_t Update posterior: $\Sigma_t^{-1} = \Sigma_{t-1}^{-1} + P_{I_t}$ and $p_t = \Sigma_t(\Sigma_{t-1}^{-1}p_{t-1} + S_{I_t}^\top(S_{I_t}S_{I_t}^\top)^{-1}Y_t)$;**end for**

Full-information case As a gentle start, we explain how the update rule would look like if we had full information about the outcome in each time step. The update in this case is identical with the standard Gaussian one-step update:

$$\begin{aligned} \Sigma_t &= \Sigma_{t-1} - \Sigma_{t-1}(\Sigma_{t-1} + I)^{-1}\Sigma_{t-1} & \text{or equiv.} & \quad \Sigma_t^{-1} = \Sigma_{t-1}^{-1} + I, \\ \mu_t &= \Sigma_t(\Sigma_{t-1}^{-1}\mu_{t-1} + X_t) & \text{or equiv.} & \quad \mu_t = \mu_{t-1} + \Sigma_t(X_t - \mu_{t-1}). \end{aligned}$$

Here we use subindex $t-1$ for the prior parameters and t for the posterior parameters in time step t , and denote by X_t the outcome (observed in this case), encoded by an M -dimensional unit vector.

General case Moving away from the full-information case, we face the problem of not observing the outcome, only some symbol that is governed by the signal matrix of the action we chose and the outcome itself. If we denote, as above, the outcome at time step t by an M -dimensional unit vector X_t , then the observation symbol can be thought of as a unit vector given by $Y_t = S_i X_t$, provided the chosen action is i . It follows that what we observe is a linear transformation of the sample from the outcome distribution.

Following the Bayes update rule and assuming we chose action i at time step t , we derive the corresponding Gaussian posterior given that the likelihood of the observation is $\pi(Y|p) \sim N(S_i p, S_i S_i^\top)$. After some algebraic manipulations we get that the posterior distribution is Gaussian with covariance $\Sigma_t = (\Sigma_{t-1}^{-1} + P_i)^{-1}$ and mean $p_t = \Sigma_t(\Sigma_{t-1}^{-1}p_{t-1} + P_i X_t)$, where $P_i = S_i^\top(S_i S_i^\top)^{-1}S_i$ is the orthogonal projection to the image space of S_i^\top . Note that even though X_t is not observed, the update can be performed, since $P_i X_i = S_i^\top(S_i S_i^\top)^{-1}S_i X_t = S_i^\top(S_i S_i^\top)^{-1}Y_t$.

A significant advantage of this method of tracking the outcome distribution as opposed to keeping track of loss difference estimates (as done in previous works), is that feedback from one action can provide information about losses across all the actions. We believe that this property has a major role in the empirical performance improvement over existing methods.

An important part in analyzing our algorithm is to show that, despite the fact that the outcome distribution is not Gaussian, the update tracks the true outcome distribution well. For an illustration of tracking the true outcome distribution with the above update, see Figures 1(b) and 1(c).

3.2 Selection rules

For selecting actions given the posterior parameters, we propose two versions for the selection rule:

1. Draw a random sample p from the distribution $N(p_{t-1}, \Sigma_{t-1})$, project the sample to the probability simplex, then choose the action that minimizes the loss for outcome distribution p . This rule is a close relative of Thompson-sampling. We call this version of the algorithm BPM-TS.
2. Use p_{t-1} and Σ_{t-1} to build a confidence ellipsoid for p^* , enumerate all actions whose cells intersect with this ellipsoid, then choose the action that was chosen the fewest times so far (called BPM-LEAST).

Our experiments demonstrate the performance of both versions. We analyze version BPM-LEAST.

4 Analysis

We now analyze BPM-LEAST that uses the Gaussian updates, and considers a set of feasible actions based on the criterion that an action is feasible if its optimality cell intersects with the ellipsoid

$$\left\{ p: \|p - p_t\|_{\Sigma_t^{-1}} \leq 1 + \sqrt{\frac{1}{2} N \log MT} \right\}.$$

From these feasible actions, it picks the one that has been chosen the fewest times up to time step t . For this version of the algorithm, the following regret bound holds.

Theorem 1. *Given a locally observable partial-monitoring problem (L, H) with prior information p_0, Σ_0 , the algorithm BPM-LEAST achieves expected regret*

$$R_T \leq C \sqrt{TN \log(MT)},$$

where C is some problem-dependent constant.

The above constant C depends on two main factors, both of them related to the feedback structure. The first one is the sum of the smallest eigenvalues of $S_i S_i^\top$ for every action i . The second is related to the local observability condition. As the condition says, for every neighboring action pairs i and j , $\ell_i - \ell_j \in \text{Im} S_i^\top \oplus \text{Im} S_j^\top$. This means that there exist v_{ij} and v_{ji} vectors such that $\ell_i - \ell_j = S_i^\top v_{ij} - S_j^\top v_{ji}$. The constant depends on the maximum 2-norm of these v_{ij} vectors.

The proof of the theorem is deferred to the supplementary material. In a nutshell, the proof is divided into two main parts. First we need to show that the update rule—even though the underlying distribution is not Gaussian—serves as a good tool for tracking the true outcome distribution. After some algebraic manipulations, the problem reduces to a finding a high probability upper bound for norms of weighted sums of noise vectors. To this end, we used the martingale version of the *matrix Hoeffding* inequality [12, Theorem 1.3].

Then, we need to show that the confidence ellipsoid shrinks fast enough that if we only choose actions whose cell intersect with the ellipsoid, we do not suffer a large regret. In the core of proving this, we arrive at a term where we need to upper bound $\|\ell_i - \ell_j\|_{\Sigma_t}$, for some neighboring action pairs (i, j) , and we show that due to local observability and the speed at which the posterior covariance shrinks, this term can be upper bounded by roughly $1/\sqrt{t}$.

5 Experiments

First, we run extensive evaluations of BPM on various synthetic datasets and compare the performance against CBP [10] and FeedExp3 [7]. The datasets used in the simulated experiments are identical to the ones used by Bartók et al. [10] and thus allow us to benchmark against the current state of the art. We also provide results of BPM on a dataset that was collected by Singla and Krause [13] from real interactions with many users on the Amazon Mechanical Turk (AMT) [14] crowdsourcing platform. We present the details of the datasets used and the summarize our results and findings in this section.

5.1 Implementation Details

In order to implement BPM, we made the following implementation choices:

1. To use BPM-LEAST (see Section 3.2), we need to recover the current feasible actions. We do so by sampling multiple (10000) times from concentric Gaussian ellipsoids centred at the current mean (p_t) and collect feasible actions based on which cells the samples lie in. We resort to sampling for ease of implementation because otherwise we deal with the problem of finding the intersection between an ellipsoid and a simplex in M -dimensional space.
2. To implement BPM-TS, we draw p from the distribution $N(p_{t-1}, \Sigma_{t-1})$. We then project it back to the simplex to obtain a probability distribution on the outcome space.

Primarily due to sampling, both our algorithms are computationally more efficient than the existing approaches. In particular, BPM-TS is ideally suited for real world tasks as it is several orders of magnitude faster than existing algorithms during all our experiments. In each iteration, BPM-TS only needs to draw one sample from a multivariate gaussian and does not need any cell decompositions or expensive computations to obtain high dimensional intersections.

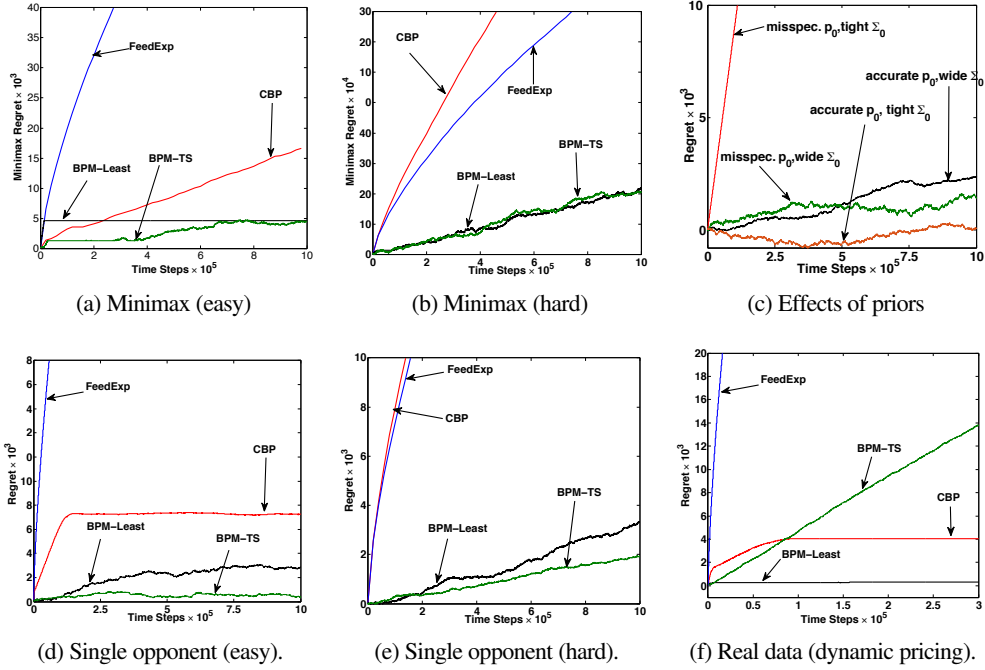


Figure 2: (a,b,d,e) Comparing BPM on the locally non-observable game ((a,d) benign opponent; (c,e) hard opponent). Hereby, (a,b) show the pointwise maximal regret over 15 scenarios, and (d,e) show the regret against a single opponent strategy. (c) shows the effect of a misspecified prior. (f) is the performance of the algorithms on the real dynamic pricing dataset.

5.2 Simulated dynamic pricing games

Dynamic pricing is a classic example of partial monitoring (see the introduction), and we compare the performance of the algorithms on the small but not locally observable game described in Bartók et al. [10]. The loss matrix and feedback tables for the dynamic pricing game are given by:

$$L = \begin{pmatrix} 0 & 1 & \dots & N-1 \\ c & 0 & \dots & N-2 \\ \vdots & \ddots & \ddots & \vdots \\ c & \dots & c & 0 \end{pmatrix}; \quad H = \begin{pmatrix} y & y & \dots & y \\ n & y & \dots & y \\ \vdots & \ddots & \ddots & \vdots \\ n & \dots & n & y \end{pmatrix}.$$

Recall that the game models a repeated interaction of a seller with buyers in a market. Each buyer can either buy the product at the price (signal “ y ”) or deny the offer (signal “ n ”). The corresponding loss to the seller is either a known constant c (representing opportunity cost) or the difference between offered price and the outcome of the customer’s latent valuation of the product (willingness-to-pay). A similar game models procurement processes as well. Note that this game does not satisfy local observability. While our theoretical results require this condition, in practice, if the opponent does not intentionally select harsh regions on the simplex, BPM remains applicable. Under this setting, expected *individual* regret is a reasonable measure of performance. That is, we measure the expected regret for fixed opponent strategies. We also consider the *minimax* expected regret, which measures worst-case performance (pointwise maximum) against multiple opponent strategies.

Benign opponent While the dynamic pricing game is not locally observable in general, certain opponent strategies are easier to compete with than others. Specifically, if the stochastic opponent chooses an outcome distribution that is away from the intersection of the cells that do not have local observability, the learning happens in “non-dangerous” or benign regions. We present results under this setting for simulated dynamic pricing with $N = M = 5$. The results shown in Figures 2(a) and 2(d) illustrate the benefits of both variants of BPM over previous approaches. We achieve an order of magnitude reduction in the regret suffered w.r.t. both the minimax and the individual regret.

Harsh opponent For the same problem, with opponent chooses close to the boundary of the cells of two non-locally observable actions, the problem becomes harder. Still, BPM dramatically outperforms the baselines and suffers very little regret as shown in Figures 2(b) and 2(e).

Effect of the prior We study the effects of a misspecified prior in Figure 2(c). As long as the initial confidence interval specified by the prior covariance is large enough to contain the opponent’s distribution, an incorrectly specified prior mean does not have an adverse effect on the performance of BPM. As expected, if the prior confidence ellipse used by BPM does not contain the opponent’s outcome distribution, however, the regret grows linear in time. Further, if the prior is very informative (accurately specified prior mean and tight confidence ellipse), very little regret is suffered.

5.3 Results on real data

Dataset description We simulate a procurement game based on real data. Parameter estimation was done by posting a Human Intelligence Task (HIT) on the Amazon Mechanical Turk (AMT) platform. Motivated by an application in viral marketing, users were asked about the price they would accept for (hypothetically) letting us post promotional material to their friends on a social networking site. The survey also collected features like age, geographic region, number of friends in the social network, activity levels (year of joining, time spent per day etc.). Note that since the HIT was just a survey and the questions were about a hypothetical scenario, participants had no incentives to misreport their responses. Complete responses were collected from approx. 800 participants. See [13] for more details.

The procurement game We simulate a procurement auction by playing back these responses offline. The game is very similar in structure to dynamic pricing, with the optimal action being the best fixed price that maximized the marketer’s value or equivalently, minimized the loss. We sampled iid from the survey data and perturbed the samples slightly to simulate a stream of 300000 potential users. At each iteration, we simulate a user with a private valuation generated as a function of her attributes. We discretized the offer prices and the private valuations to be one of 11 values and set the opportunity cost of losing a user due to low pricing to be 0.5. Thus we recover a partial monitoring game with 11 actions and 11 outcomes with a 0/1 feedback matrix.

Results We present the results of our evaluation on this dataset in Figure 2(f). Notice that although the game is not locally observable, the outcome distribution does not seem to be in a difficult region of the cell decomposition as the adaptive algorithms (CBP and both versions of BPM) perform well. We note that the total regret suffered by BPM-LEAST is a *factor of 10 lower* than the regret achieved by CBP on this dataset. The plots are averaged over 30 runs of the competing algorithms on the stream. To the best of our knowledge, this is the first time partial monitoring has been evaluated on a real world problem of this size.

6 Conclusions and future work

We introduced a new family of algorithms for locally observable partial-monitoring problems against stochastic opponents. We also enriched the model of partial monitoring with the possibility of incorporating prior information about the outcome distribution in the form of a confidence ellipsoid. The new insight of our approach is that instead of tracking loss differences, we explicitly track the true outcome distribution. This approach not only eases computational overhead but also helps achieve low regret by being able to transfer information between actions. In particular, BPM-TS runs orders of magnitude faster than any existing algorithms, opening the path for the model of partial monitoring to be applied on realistic settings involving large numbers of actions and outcomes.

Future work includes extending our method for adversarial opponents. Bartók [11] already uses the idea of tracking the true outcome distribution with the help of a *confidence parallelotope*, which is rather close to our approach, but has the same shortcomings as other algorithms that track loss differences: it can not transfer information between actions. Extending our results to problems with large action and outcome spaces is also an important direction: if we have some prior information about the similarities between outcomes and/or actions, we have a chance for a reasonable regret.

Acknowledgments This research was supported in part by SNSF grant 200021_137971, ERC StG 307036 and a Microsoft Research Faculty Fellowship.

References

- [1] V. G. Vovk. Aggregating strategies. In *COLT*, pages 371–386, 1990.
- [2] Nick Littlestone and Manfred K. Warmuth. The weighted majority algorithm. *Inf. Comput.*, 108(2):212–261, 1994.
- [3] Jean-Yves Audibert and Sébastien Bubeck. Minimax policies for adversarial and stochastic bandits. In *COLT*, 2009.
- [4] Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM J. Comput.*, 32(1):48–77, 2002.
- [5] Yisong Yue, Josef Broder, Robert Kleinberg, and Thorsten Joachims. The K-armed dueling bandits problem. *Journal of Computer and System Sciences*, 78(5):1538–1556, 2012.
- [6] Nir Ailon, Thorsten Joachims, and Zohar Karnin. Reducing dueling bandits to cardinal bandits. *arXiv preprint arXiv:1405.3396*, 2014.
- [7] Antonio Piccolboni and Christian Schindelhauer. Discrete prediction games with arbitrary feedback and loss. In *COLT/EuroCOLT*, pages 208–223, 2001.
- [8] Nicolò Cesa-Bianchi, Gábor Lugosi, and Gilles Stoltz. Regret minimization under partial monitoring. *Math. Oper. Res.*, 31(3):562–580, 2006.
- [9] Gábor Bartók, Dávid Pál, and Csaba Szepesvári. Minimax regret of finite partial-monitoring games in stochastic environments. *Journal of Machine Learning Research - Proceedings Track (COLT)*, 19:133–154, 2011.
- [10] Gábor Bartók, Navid Zolghadr, and Csaba Szepesvári. An adaptive algorithm for finite stochastic partial monitoring. In *Proceedings of the 29th International Conference on Machine Learning, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012*, 2012.
- [11] Gábor Bartók. A near-optimal algorithm for finite partial-monitoring games against adversarial opponents. In *COLT 2013 - The 26th Annual Conference on Learning Theory, June 12-14, 2013, Princeton University, NJ, USA*, pages 696–710, 2013.
- [12] Joel A Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of Computational Mathematics*, 12(4):389–434, 2012.
- [13] Adish Singla and Andreas Krause. Truthful incentives in crowdsourcing tasks using regret minimization mechanisms. In *International World Wide Web Conference (WWW)*, 2013.
- [14] Amazon Mechanical Turk platform. URL <https://www.mturk.com>.

A Proof of Theorem 1

A.1 Validity of the update

We assume that p^* , the true opponent strategy, is within some distance from our initial prior p_0 , measured in Σ_0^{-1} -distance:

$$\|p_0 - p^*\|_{\Sigma_0^{-1}} \leq 1.$$

First we observe that the update can be rewritten in a cumulative form, to see how the parameters change from the initial prior (p_0, Σ_0) :

$$\begin{aligned}\Sigma_t^{-1} &= \Sigma_0^{-1} + \sum_{s=1}^t P_{I_s} \\ \Sigma_t^{-1} p_t &= \Sigma_0^{-1} p_0 + \sum_{s=1}^t P_{I_s} X_s.\end{aligned}$$

Now let us investigate the Σ_t^{-1} -distance of p_t from p^* !

$$\|p_t - p^*\|_{\Sigma_t^{-1}} = \left\| \Sigma_t \Sigma_0^{-1} p_0 + \Sigma_t \sum_{s=1}^t P_{I_s} X_s - p^* \right\|_{\Sigma_t^{-1}}$$

Now we decompose the samples X_s to mean and noise with the new notation $X_s = p^* + \epsilon_s$, yielding

$$\begin{aligned}\|p_t - p^*\|_{\Sigma_t^{-1}} &= \left\| \Sigma_t \Sigma_0^{-1} p_0 + \Sigma_t \sum_{s=1}^t P_{I_s} (p^* + \epsilon_s) - p^* \right\|_{\Sigma_t^{-1}} \\ &= \left\| \Sigma_t \Sigma_0^{-1} p_0 + \underbrace{\Sigma_t \left(\sum_{s=1}^t P_{I_s} - \Sigma_t^{-1} \right)}_{-\Sigma_0^{-1}} p^* + \Sigma_t \sum_{s=1}^t P_{I_s} \epsilon_s \right\|_{\Sigma_t^{-1}} \\ &\leq \left\| \Sigma_t \Sigma_0^{-1} (p_0 - p^*) \right\|_{\Sigma_t^{-1}} + \left\| \Sigma_t \sum_{s=1}^t P_{I_s} \epsilon_s \right\|_{\Sigma_t^{-1}}.\end{aligned}$$

We deal with the two resulting terms separately.

$$\begin{aligned}\left\| \Sigma_t \Sigma_0^{-1} (p_0 - p^*) \right\|_{\Sigma_t^{-1}}^2 &= (p_0 - p^*)^\top \underbrace{\Sigma_0^{-1} \Sigma_t \Sigma_0^{-1}}_{(I - \Sigma_0^{-1} \sum_{s=1}^t P_{I_s})^{-1}} (p_0 - p^*) \\ &\leq \|p_0 - p^*\|_{\Sigma_0^{-1}}^2 \leq 1.\end{aligned}$$

The second term is harder. Basically this is the term where we “pay the price” for not having started with a Gaussian distribution. We need to show that

$$\left\| \Sigma_t \sum_{s=1}^t P_{I_s} \epsilon_s \right\|_{\Sigma_t^{-1}} = \left\| \sum_{s=1}^t \sqrt{\Sigma_t} P_{I_s} \epsilon_s \right\|$$

is bounded with high probability. For any given action sequence, the above expression is a sum of independent random matrices. Now we recite a concentration inequality we need:

Theorem 2 (Matrix Hoeffding Theorem [12, Theorem 1.3]). *Consider a finite sequence $\{X_k\}$ of independent, random, self-adjoint matrices with dimension d , and let $\{A_k\}$ be a sequence of fixed self adjoint matrices. Assume that each random matrix satisfies*

$$\mathbb{E}X_k = 0 \quad \text{and} \quad X_k^2 \preceq A_k^2 \quad \text{almost surely.}$$

Then, for all $t \geq 0$,

$$P\left(\left\| \sum_k X_k \right\|_2 \geq t\right) \leq d \exp(-t^2/8\sigma^2) \quad \text{where} \quad \sigma^2 = \left\| \sum_k A_k^2 \right\|_2.$$

The above theorem can be extended to rectangular matrices, using the “dilation trick”²: for rectangular matrices $B_k \in \mathbb{R}^{d_1 \times d_2}$, we use the theorem with

$$X_k = \begin{pmatrix} 0 & B_k \\ B_k^\top & 0 \end{pmatrix} \in \mathbb{R}^{d_1+d_2}.$$

In our case, $X_s = \sqrt{\Sigma_t} P_{I_s} \epsilon_s$. Also note that here we need the martingale version of the inequality, which also holds, according to Section 7 of Tropp [12]. After algebraic manipulations, we arrive at

$$P\left(\left\|\sum_s P_s \epsilon_s\right\|_{\Sigma_t^{-1}} \geq \sqrt{\frac{1}{2} N \log \frac{M+1}{\delta}}\right) \leq \delta.$$

Putting together the terms we get that

$$\|p_t - p^*\|_{\Sigma_t^{-1}} \leq 1 + \sqrt{\frac{1}{2} N \log \frac{M+1}{\delta}}$$

with probability at least $1 - \delta$.

A.2 Regret

Now we turn our attention to calculating the regret of the algorithm that chooses the action that is chosen the fewest times so far among the actions whose optimality cells intersect with the current confidence ellipsoid. To accommodate the error for the outcome distribution not being Gaussian, we use the ellipsoid defined as

$$\left\{p : \|p - p_t\|_{\Sigma_t^{-1}} \leq 1 + \sqrt{\frac{1}{2} N \log \frac{M+1}{\delta}}\right\}.$$

The regret in a turn results from choosing a suboptimal action. Let us assume wlog that the optimal action is action 1, the true opponent strategy is p^* , and the chosen action is action k . Then, the instantaneous regret is

$$r_t = (\ell_k - \ell_1)^\top p^*.$$

Now if we pick a point p in the intersection of the cell of action k and the confidence ellipse, we can connect p^* and p with a line segment. That segment goes through the cells of, say, $1 = i_0, i_1, \dots, i_d = k$. Then we can write

$$\begin{aligned} (\ell_k - \ell_1)^\top p^* &= \sum_{j=1}^d (\ell_{i_j} - \ell_{i_{j-1}})^\top p^* \\ &= \sum_{j=1}^d (\ell_{i_j} - \ell_{i_{j-1}})^\top (p^j - p^*), \end{aligned}$$

where we denote by p^j the point where our line segment intersects the boundary of cells i_{j-1} and i_j . The above equation is true because for every j , $(\ell_{i_j} - \ell_{i_{j-1}})^\top p_j = 0$. Now we upper bound, for every j , the term

$$(\ell_{i_j} - \ell_{i_{j-1}})^\top (p^j - p^*) \leq \|\ell_{i_j} - \ell_{i_{j-1}}\|_{\Sigma_t} \|p^j - p^*\|_{\Sigma_t^{-1}},$$

with the help of Hölder’s inequality. We know from the previous section that $\|p^j - p^*\|_{\Sigma_t^{-1}}$ can be upper bounded with high probability. It remains to upper bound the first term.

With the help of the local observability condition, we have

$$\ell_{i_j} - \ell_{i_{j-1}} = S_{i_j}^\top v_{i_j, i_{j-1}} - S_{i_{j-1}}^\top v_{i_{j-1}, i_j},$$

for some $v_{i_{j-1}, i_j}, v_{i_j, i_{j-1}}$, and thus the problem reduces to upper bounding $\|S_i^\top v_{i, i'}\|_{\Sigma_t}$ for all $1 \leq i, i' \leq N$:

$$\begin{aligned} \|S_i^\top v_{i, i'}\|_{\Sigma_t}^2 &= \left\| \sqrt{S_i \Sigma_t S_i^\top} v_{i, i'} \right\|_2^2 \\ &\leq \|S_i \Sigma_t S_i^\top\|_2 \|v_{i, i'}\|_2^2 \\ &\leq \|S_i \Sigma_t S_i^\top\|_2 V_{\max}^2, \end{aligned}$$

²See remark 3.11 in Tropp [12].

where $V_{\max} = \max_{1 \leq i, i' \leq N} \|v_{i, i'}\|_2$.

$$\begin{aligned}
\|S_i \Sigma_t S_i^\top\|_2 &= \left\| S_i \left(\Sigma_0^{-1} + \sum_{s=1}^t P_{I_s} \right)^{-1} S_i^\top \right\|_2 \\
&\leq \left\| S_i (\Sigma_0^{-1} + n_i P_i)^{-1} S_i^\top \right\|_2 \\
&= \left\| (S_i^+)^+ (\Sigma_0^{-1} + n_i P_i)^{-1} (S_i^{\top+})^+ \right\|_2 \\
&\leq \left\| \left(n_i S_i^{\top+} S_i^\top (S_i S_i^\top)^{-1} S_i S_i^+ \right)^+ \right\|_2 \\
&\leq \frac{c_i}{n_i}
\end{aligned}$$

for some constant c_i .

Putting everything together we have that the instantaneous regret at time step t is

$$r_t \leq 2V_{\max} K_i \sqrt{\frac{C_1}{n_{\min}}} \left(C + \sqrt{\frac{1}{2} N \log \frac{M+1}{\delta}} \right).$$

Since our algorithm picks the action that is chosen the fewest number of times, it ensures that $n_{\min} \geq t/N$. Summing up the instantaneous regret for every turn we get the desired result

$$R_T \leq C_2 \sqrt{TN \log MT / \delta} \quad \text{w.p.} \geq 1 - \delta.$$

Setting δ to $1/\sqrt{T}$, we get the desired result.