

# Exact Passive-Aggressive Algorithms for Multiclass Classification Using Bandit Feedbacks

**Maanik Arora**

MAANIK.ARORA@RESEARCH.IIIT.AC.IN

*Machine Learning Lab, Kohli Center on Intelligent Systems, IIIT Hyderabad, India*

**Naresh Manwani**

NARESH.MANWANI@IIIT.AC.IN

*Machine Learning Lab, Kohli Center on Intelligent Systems, IIIT Hyderabad, India*

**Editors:** Sinno Jialin Pan and Masashi Sugiyama

## Abstract

In many real-life classification problems, we may not get exact class labels for training samples. One such example is bandit feedback in multiclass classification. In this setting, we only get to know whether our predicted label is correct or not. Due to which, we are left in uncertainty about the actual class label when we predict the wrong class. This paper proposes exact passive-aggressive online algorithms for multiclass classification under bandit feedback (EPABF). The proposed approach uses an exploration-exploitation strategy to guess the class label in every trial. To update the weights, we solve a quadratic optimization problem under multiple class separability constraints and find the exact solution. We do this by finding active constraints using the KKT conditions of the optimization problem. These constraints form a support set that determines the classes for which the weight vector needs to be updated. We propose three different variants of the weight update rule, which vary based on the aggressiveness to correct the mistake. These are called EPABF, EPABF-I, and EPABF-II. We also provide mistake bounds for the proposed EPABF, EPABF-I, and EPABF-II. Experiments demonstrated that our proposed algorithms perform better than other bandit feedback-based approaches and comparably to the full information approaches.

## 1. Introduction

Online learning of multiclass classifiers is an important and well-studied problem in machine learning. The goal is to classify examples into a set of classes. Digit recognition (Ma and Zhang, 2015), text classification (McCallum, 1999), recommendation systems (Li et al., 2016), optical character recognition (OCR), face recognition, tagging locations in vacation photos etc. are some of the well known applications of multiclass classifier.

Standard multiclass setting assumes that the exact (true) class labels of the training examples are given. This is called the full information case. However, in many real-life problems, getting true labels of the examples is a challenging task. Instead, we may get partial information about the class label. One such case is where we only know whether the predicted class label is the same as the true label. This is called *bandit feedback* setting (Sham M. Kakade, 2008). For example, when a user types a query in the recommendation system, she is presented with different links based on the user profile. A click on one of the link represents positive feedback, and no click corresponds to a negative response. On the other hand, the algorithm does not know what would have happened when other links been

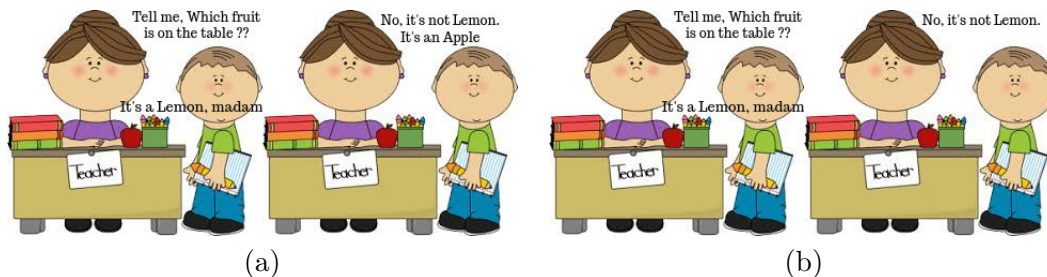


Figure 1: Example illustrating multiclass classification with (a) full information case and (b) bandit feedback case.

presented. Figure 1 presents the difference between the full information case and bandit information case using an example.

Consider the learning with bandit feedback as follows. The learner receives an instance  $\mathbf{x}^t$  at trial  $t$ . The learner has to predict its label out of a set of  $K$  labels. Based on the current parameters, it predicts the label  $\hat{y}^t$ , where  $\hat{y}^t \in [K]$ . Then it receives a partial feedback  $\mathbb{1}_{\{\hat{y}^t \neq y^t\}}$ . Here  $y^t$  is the true label of  $\mathbf{x}^t$ . So if the prediction is correct, that is  $\hat{y}^t = y^t$ , the feedback is 1 and learner exactly knows that the true class is  $\hat{y}^t$ . On the other hand, if  $\hat{y}^t \neq y^t$ , the bandit feedback is 0 and the learner gets to know that  $\hat{y}^t$  is not the true class. But, the learner still does not know the true class in this case. The performance of the learner is measured using the cumulative sum of mistakes in  $T$  trials, which is  $\sum_{t=1}^T \mathbb{1}_{\{\hat{y}^t \neq y^t\}}$ . Since the learning algorithm is unaware of the true label, it makes multiclass classification in the bandit setting harder than multiclass classification problems with a full information setting. In [Sham M. Kakade \(2008\)](#), authors extend the multiclass Perceptron ([Crammer and Singer, 2003](#)) to learn classifier using bandit feedback.

Inspired by passive-aggressive updates ([Crammer et al., 2006](#)), [Zhong and Dauce \(2015\)](#) propose bandit passive-aggressive algorithms (BPA) for multiclass classification using bandit feedbacks. Updates proposed in BPA are much more aggressive compared to Banditron ([Sham M. Kakade, 2008](#)). BPA shows better performance compared to Banditron ([Zhong and Dauce, 2015](#)). BPA finds the updates by solving a quadratic optimization problem in each trial. However, for multiclass classifier learning, it ignores the structure of multiple classes and reduces the problem to binary classification.

This paper proposes exact passive-aggressive algorithms under bandit feedback (EPABF) for online multiclass classification in the bandit setting. Our approach takes care of the issues in the method proposed in [Zhong and Dauce \(2015\)](#) by considering the separability conditions among different classes. At every trial, KKT optimality conditions are used to find the active constraints. These active constraints correspond to the support classes whose parameters will be updated. A similar approach was used for handling full information case ([Matsushima et al., 2010](#)). We observe that the proposed approaches perform better than the other bandit algorithms. Our key contributions in this paper are as follows.

1. We derive update rules for EPABF, EPABF-I, and EPABF-II under bandit feedback setting. EPABF algorithms update the parameters at trial  $t$  by minimizing convex

quadratic optimization problems. We find the exact solution to these optimization problems.

2. We propose a support class set finding algorithms for all the three variants and show their correctness.
3. We provide mistake bounds for all three variants EPABF, EPABF-I, and EPABF-II.
4. We perform simulations on different datasets and show the proposed algorithms' effectiveness compared to other bandit feedback algorithms and full information algorithms.

## 2. Related Work

### Learning Under Full Information Setting

In the full information case, the learner knows the exact labels of all the examples in the training set. Popular batch learning algorithms for multiclass learning are discussed in (Crammer and Singer, 2002; Hsu and Lin, 2002; Ou and Murphey, 2007). In the online setting, we get instances sequentially, so the algorithm has to make predictions even when learning continuously. When it makes the wrong prediction, then it updates the parameters of the classifier else; it does not change the parameters. Perceptron is one of the earliest online algorithms for learning binary classifiers (Rosenblatt, 1960). Extensions of Perceptron algorithm for multiclass is proposed in Crammer and Singer (2003); Fink et al. (2006). Cesa-Bianchi et al. (2005) proposed a second-order Perceptron algorithm, which uses a second-order derivative in updating the parameters. Passive-aggressive learning is another framework in which the algorithm makes more aggressive updates to achieve the smallest loss on current misclassified example (Crammer et al., 2006) by solving a constrained quadratic optimization problem in each trial. However, for multiclass classifier learning, it ignores the structure of multiple classes and reduces the binary classification problem. Matsushima et al. (2010) proposed an exact passive-aggressive approach for the multiclass classification problem, which addresses the above issue. It first determines the set of classes called support classes when an input is received. Then the weight vectors are updated corresponding to all the support classes.

### Learning Under Bandit Feedback

In many practical situations, the learner does not know the exact class label of the training examples. Instead, it may just observe some partial feedback. One such feedback is whether the predicted label is the same as the actual label. This is called bandit feedback. Several classification algorithms exist that address the bandit feedback setting. Banditron (Sham M. Kakade, 2008) extends the Perceptron algorithm to deal with the bandit feedback. Banditron uses an exploitation-exploration scheme proposed in Auer et al. (2003). When it updates, it replaces the gradient of the loss function with an unbiased estimator. When the data is linearly separable, the expected number of mistakes made by Banditron is  $O(\sqrt{T})$  in  $T$  rounds. In the general case, the expected number of mistakes of Banditron is  $O(T^{2/3})$ . Another bandit algorithm, named NEWTRON (Elad Hazan, 2011), is based

on the online Newton method. It uses strongly convex loss objective function (adding regularization term with the loss function) and Follow-The-Regularized-Leader (FTRL) strategy to achieve  $O(\log T)$  regret bound in the best case and  $O(T^{2/3})$  regret bound in the worst case. Second-order Perceptron is also extended in bandit feedback setting by [Crammer and Gentile \(2011\)](#). It uses upper-confidence bounds (UCB) ([Auer et al., 2002](#)) based approach to make trade-off between exploration and exploitation and achieves regret bound of  $O(\sqrt{T} \log(T))$ .

Recently, [Beygelzimer et al. \(2019\)](#) proposed an algorithm based on Kernel Perceptron. Authors proposed efficient online multiclass linear classification algorithms with bandit feedback when the data is linearly separable by a margin. They showed that their algorithm achieves a near-optimal bound of  $O(K^{-1})$  under strong linear separability condition ([Beygelzimer et al., 2019](#)). Among the second-order algorithms, ([Beygelzimer et al., 2017](#)) proposed a second-order algorithm with  $O(\sqrt{T})$  regret for the bandit online multiclass problem.

[Zhong and Dauce \(2015\)](#) extended the passive-aggressive online multiclass learning approach proposed in ([Crammer et al., 2006](#)) in the bandit feedback setting. They solved a constraint optimization problem in each trial to update the parameters which inherit similar issues as in the method proposed in ([Crammer et al., 2006](#)) (i.e., they solve a relaxed version of the quadratic optimization problem). In this paper, we use the ideas presented in [Matsushima et al. \(2010\)](#) to address the issues in BPA ([Zhong and Dauce, 2015](#)). We propose an exact passive-aggressive approach for multiclass classifiers under bandit feedback.

### 3. Passive-Aggressive Online Learning Under Full Information

The passive-aggressive online learning of multiclass classifier is proposed in [Crammer et al. \(2006\)](#). The learning is performed in a sequence of trials. At trial  $t$ , the learner is presented an example  $x_t \in \mathbb{R}^d$  and is required to predict the label out of a set of  $K$  labels. Let  $[K] = \{1, \dots, K\}$  denote the set of labels. Let  $w_1^t, \dots, w_K^t$  be the parameters of the multiclass classifier at the beginning of trial  $t$ . The label predicted by the algorithm at round  $t$  is given by,  $\hat{y}_t = \arg \max_{j \in [K]} (w_j^t \cdot x_t)$ . Then the algorithm observes the correct label  $y^t$ . The loss incurred by the algorithm is as follows.

$$l_{\text{hinge}}^t = \max[0; 1 - w_{y^t}^t \cdot x_t + \max_{j \in \hat{y}^t} w_j^t \cdot x_t] \tag{1}$$

where  $l_{\text{hinge}}$  is hinge loss for multiclass classification and is a convex surrogate of 0-1 loss [Crammer et al. \(2006\)](#). The passive-aggressive approach ([Crammer et al., 2006](#)) proposes three different ways to update the parameters as follows.

1. PA: In this approach, parameters  $w_1^{t+1}, \dots, w_K^{t+1}$  are found such that they are closest to  $w_1^t, \dots, w_K^t$  and the loss incurred on  $x^t$  becomes zero.

$$w_1^{t+1}, \dots, w_K^{t+1} = \arg \min_{w_1, \dots, w_K} \frac{1}{2} \sum_{j=1}^K \|w_j - w_j^t\|^2$$

$$\text{s.t: } \max[0; 1 - w_{y^t}^t \cdot x_t + \max_{j \in \hat{y}^t} w_j \cdot x_t] = 0$$

2. PA-I: PA makes very aggressive updates as it tries to find the parameters which incur zero loss on  $x^t$ . To relax that, PA-I allows some error. Thus, PA-I makes a trade-off between the closeness to  $w_1^t, \dots, w_K^t$  and aggressiveness to minimize the loss on example  $x^t$ .  $C$  is the hyper-parameter which controls the aggressiveness.

$$w_1^{t+1}, \dots, w_K^{t+1} = \arg \min_{w_1, \dots, w_K} \frac{1}{2} \sum_{j=1}^K k w_j - w_j^t k^2 + C$$

$$\text{s.t: } 0; w_{y^t} - x^t \max_{j \in y^t} w_j - x^t \leq 1$$

3. PA-II: PA-II also allows some error, but it minimizes the square of the loss in contrast to PA-I. Thus, it finds new parameters as follows.

$$w_1^{t+1}, \dots, w_K^{t+1} = \arg \min_{w_1, \dots, w_K} \frac{1}{2} \sum_{j=1}^K k w_j - w_j^t k^2 + C^2$$

$$\text{s.t: } w_{y^t} - x^t \max_{j \in y^t} w_j - x^t \leq 1$$

Exact solutions of these three approaches are proposed in [Matsushima et al. \(2010\)](#). The update equations proposed in [Matsushima et al. \(2010\)](#) are as follows.

$$w_v^{t+1} = w_v^t - \frac{1}{v} (1 - I_{f_{v=y^t}g}) x^t + I_{f_{v=y^t}g} \sum_{v \in y^t} \frac{1}{v} x^t \quad v = 1 \dots K$$

where

$$\frac{1}{v} = \begin{cases} \frac{1}{kx^t k^2} I_v^t - \frac{1}{jS^t_{j+1}} \sum_{j \in 2S^t} I_j^t & ; \text{ PA} \\ \frac{1}{kx^t k^2} I_v^t \max \left\{ \frac{j2S^t I_j^t}{jS^t_j} + \frac{Ckx^t k^2}{jS^t_j} ; \frac{j2S^t I_j^t}{jS^t_{j+1}} \right\} & ; \text{ PA-I} \\ \frac{I_v^t}{kx^t k^2} \frac{kx^t k^2 \frac{1}{2C}}{(jS^t_{j+1}) kx^t k^2 + \frac{jS^t_j}{2C}} \sum_{j \in 2S^t} I_j^t & ; \text{ PA-II} \end{cases}$$

In the above,  $S^t$  denotes the support set which contains the indices of active constraints ([Matsushima et al. \(2010\)](#)). Exact passive-aggressive updates proposed in [Matsushima et al. \(2010\)](#) are shown to perform better than the updates proposed [Crammer et al. \(2006\)](#). However, passive-aggressive approach in [Matsushima et al. \(2010\)](#); [Crammer et al. \(2006\)](#) is developed under the full information case. That is, the algorithm knows the exact label  $y^t$  of  $x^t$  at every trial  $t$ . In this paper, we will propose the exact passive-aggressive approach under bandit feedback.

#### 4. Proposed Approach: Exact Passive Aggressive Learning Under Bandit Feedback

Here again, learning happens in online fashion. At every trial  $t$ , an example  $x^t \in \mathbb{R}^d$  is presented to the algorithm. Let  $w_1^t, \dots, w_K^t$  be the parameters of the multiclass classifier at the beginning of the trial  $t$ , where  $K$  is the total number of classes. Then the algorithm

computes the label  $\hat{y}^t = \arg \max_{j \in [K]} w_j^t \cdot x^t$ . Before asking for the bandit feedback, we form a probability distribution over different classes as follows.

$$P^t(r) = (1 - \epsilon) I_{f_{r=y^t g}} + \frac{\epsilon}{K}$$

Here  $\epsilon \in (0, 0.5)$  is exploration parameter and  $I_{f_{r=y^t g}}$  is the Indicator function. We randomly sample a label  $\varphi^t$  from  $P^t(r)$ . We predict the sampled label  $\varphi^t$ . Then, we ask for the bandit feedback  $I_{f_{\varphi^t=y^t g}}$ , which means if  $\varphi^t = y^t$ , the feedback is 1, else it is 0. Thus, if  $\varphi^t \neq y^t$ , the algorithm does not know which of the remaining  $K - 1$  labels is the true label. The objective of the learning algorithm is to minimize the number of prediction mistakes, i.e.  $\sum_{t=1}^T I_{f_{\varphi^t \neq y^t g}}$ . However, we can't minimize  $I_{f_{\varphi^t \neq y^t g}}$  as it is not continuous. Thus, we propose a new loss functions as follows.

We know that the  $l_{\text{hinge}}^t$  (Eq. (1)) is a convex surrogate of the 0-1 loss. However, one can easily see that

$$l_{\text{hinge}}^t = \max(0; 1 - a_t w_{\varphi^t}^t \cdot x^t + \max_{r \in y^t} w_r^t \cdot x^t) = \max_{r=1}^K \max(0; 1 - a_t w_{\varphi^t}^t \cdot x^t + w_r^t \cdot x^t) = \sum_{r=1}^K l_r^t$$

Thus,  $\sum_{r=1}^K l_r^t$  also forms a convex surrogate for 0-1 loss.

We use  $\hat{\varphi}_r^t = \max(0; 1 - a_t w_{\varphi^t}^t \cdot x^t + w_r^t \cdot x^t)$  as an estimator for  $l_r^t$  in presence of bandit feedback, where  $a_t = I_{f_{\varphi^t=y^t g}} = P(\varphi^t)$ . The total loss incurred at time  $t$  is given as  $\sum_{r=1}^K \hat{\varphi}_r^t$ . Note that each of  $\hat{\varphi}_r^t$  is a random variable as it depends on  $\varphi^t$ . We observe the following.

$$E[\hat{\varphi}_r^t] = E[\max(0; 1 - a_t w_{\varphi^t}^t \cdot x^t + w_r^t \cdot x^t)] = \max(0; E[1 - w_{\varphi^t}^t \cdot x^t a_t + w_r^t \cdot x^t])$$

$$\max(0; 1 - w_{y^t}^t \cdot x^t + w_r^t \cdot x^t) = l_r^t$$

$E[\sum_{r=1}^K \hat{\varphi}_r^t]$  upper bounds  $\sum_{r \in y^t} l_r^t$ . Thus, minimizing  $\sum_{r=1}^K \hat{\varphi}_r^t$  is an appropriate objective at every trial. We now derive the update equations under exact passive aggressive framework using bandit feedback. As discussed earlier, there are three variants of the passive-aggressive approach, namely, PA, PA-I, and PA-II. We name them as EPABF (Exact Passive Aggressive approach under Bandit Feedback), EPABF-I and EPABF-II respectively.

#### 4.1. EPABF Updates

In EPABF approach, after receiving the bandit feedback  $I_{f_{\varphi^t=y^t g}}$ , the new parameters of the classifier are found as follows. If  $\sum_{r=1}^K \hat{\varphi}_r^t = 0$ , then we do not update the parameters. On the other hand, if the  $\sum_{r=1}^K \hat{\varphi}_r^t > 0$ , then we find the new parameters by solving the following optimization problem.

$$w_1^{t+1} \dots w_K^{t+1} = \arg \min_{w_1, \dots, w_K} \frac{1}{2} \sum_{v=1}^K \|w_v\|^2 \quad \text{s.t.} \quad \sum_{r=1}^K \hat{\varphi}_r^t = 0$$

$$= \arg \min_{w_1, \dots, w_K} \frac{1}{2} \sum_{v=1}^K \|w_v\|^2 \quad \text{s.t.} \quad a_t w_{\varphi^t}^t \cdot x^t - w_r \cdot x^t \leq 1; \forall r \in [K]$$

This is a quadratic optimization problem with  $K$  linear constraints (one constraint corresponding to each class).  $\mathcal{L}S^t$  denotes the support class set which contain indices of active constraint. Final update equations for EPABF are as follows.

$$w_r^{t+1} = w_r^t + a_t \sum_{i \in \mathcal{L}S^t} x_i^t; \quad r = 1 \dots K$$

where

$$w_r^t = \sum_{i \in \mathcal{L}S^t} \frac{1}{kx^t k^2} \theta_r + \frac{a_t \theta_{y^t}}{1 + jS^t j a_t^2} \frac{a_t^2 \sum_{i \in \mathcal{L}S^t} \theta_i}{1 + jS^t j a_t^2} \quad (2)$$

The derivation of the step sizes is given in the supplementary le. Thus, we can determine the step sizes of  $\eta_r^t$  and obtain a complete update rule if the support class set  $\mathcal{L}S^t$  is known. The overall approach of learning classifier using EPABF is described in Algorithm 1.

---

#### Algorithm 1 Exact Passive-Aggressive Approach Under Bandit Feedback (EPABF)

---

```

1: Input:  $\mathcal{Q}$  (0;0.5), Training Set  $\mathcal{Q}$ 
2: Initialize:  $w_1^0, \dots, w_K^0$ 
3: for  $t = 1, \dots, T$  do
4:   Receive  $x^t \in \mathcal{Q}$ 
5:   Set  $y^t = \arg \max_{r \in [K]} (w_r^t \cdot x^t)$ 
6:   for  $r = 1, \dots, K$  do
7:     Define  $P^t(r) = (1 - \eta_r^t) |_{f_{r=y^t} g} + \bar{\kappa}$ 
8:   end for
9:   Sample  $\mathcal{P}^t \sim P^t$ . Receive feedback  $|_{f_{y^t} g}$ .
10:  for  $r = 1, \dots, K$  do
11:    Compute  $\theta_r^t = [1 - a_t w_{\mathcal{P}^t}^t \cdot x^t + w_r^t \cdot x^t]_+$ 
12:  end for
13:   $S^t = \text{SCA}(\theta_1, \dots, \theta_K; a_t)$ 
14:  if  $\mathcal{P}^t \not\subseteq S^t$  then
15:    for  $r \in S^t$  do
16:      
$$\eta_r^t = \frac{\theta_r}{kx^t k^2} + \frac{kx^t k^2 a_t \theta_{y^t}}{1 + jS^t j a_t^2} \frac{a_t^2 \sum_{i \in \mathcal{L}S^t} \theta_i}{1 + jS^t j a_t^2}$$

17:    end for
18:  else if  $\mathcal{P}^t \not\subseteq S^t$  then
19:    for  $r \in S^t$  do
20:      
$$\eta_r^t = \frac{1}{kx^t k^2} \theta_r \frac{a_t^2 \sum_{i \in \mathcal{L}S^t} \theta_i}{1 + jS^t j a_t^2}$$

21:    end for
22:  end if
23:  Set  $\eta_r^t = 0; \forall r \notin S^t$ 
24:  for  $r = 1 \dots K$  do
25:     $w_r^{t+1} = w_r^t + a_t |_{f_{r=y^t} g} \sum_{i \in \mathcal{L}S^t} \eta_i^t x_i^t$ 
26:  end for
27: end for
28: Output:  $w_1^{T+1}, \dots, w_K^{T+1}$ 

```

---



---

#### Algorithm 2 Support Class Set Finding Algorithm for EPABF (SCA)

---

```

1: Input:  $\theta_1, \dots, \theta_K; a_t$ 
2: Initialize:  $S^t := \emptyset$ 
3: Sort  $\theta_1, \dots, \theta_K$  in descending order. Let  $(1) \dots (K)$  be the sorted order.
4: for  $j = 1 \dots K$  do
5:   if  $\mathcal{P}^t \not\subseteq S^t$  then
6:     if  $\frac{a_t^2 \sum_{k \in S^t} \theta_{(k)}}{1 + jS^t j a_t^2} < \theta_{(j)}$  then
7:        $S^t = S^t \cup \{j\}$ 
8:     end if
9:   else if  $\mathcal{P}^t \subseteq S^t$  then
10:    if  $\frac{a_t^2 \sum_{k \in S^t} \theta_{(k)}}{1 + jS^t j a_t^2} < \theta_{(j)} + \frac{a_t \theta_{y^t}}{1 + jS^t j a_t^2}$  then
11:       $S^t = S^t \setminus \{j\}$ 
12:    end if
13:  end if
14: end for
15: Output:  $S^t$ 

```

---

Determining Support Class Set  $S^t$  The EPABF algorithm assumes requires the support class set  $S^t$  as an input. Following theorem states the necessary and sufficient condition for any class belonging to the support class set.

Theorem 1 Assume that at trial  $t$ ,  $S^t \subseteq \mathcal{C}$ . Let  $(k)$  be the  $k$ -th class in the sorted order of  $\beta_1^t, \dots, \beta_K^t$ . The necessary and sufficient condition for  $(k)$  to be in set  $S^t$  is as follows.

$$\beta_{(k)}^t > \frac{\sum_{j=1}^k a_t^2 \beta_{(j)}^t}{1 + (k-1)a_t^2} \frac{a_t \beta_{(k)}^t}{1 + (k-1)a_t^2} \quad \beta_{(k)}^t \in S^t$$

$$\beta_{(k)}^t < \frac{\sum_{j=1}^k a_t^2 \beta_{(j)}^t}{1 + (k-1)a_t^2} \quad \beta_{(k)}^t \notin S^t$$

The proof of this theorem is given in supplementary file. Theorem 1 allows us to algorithmically determine the support classes  $S^t$ . We simply sort classes according to  $\beta_r^t$  values in descending order, and the top  $J$  classes in the sorted list are the support classes (where  $J$  is the largest  $k$  that satisfies the Theorem 1 statement). Support class set finding algorithm (SCA) is described in Algorithm 2.

#### 4.2. EPABF-I Updates

EPABF-I finds new parameters which are closest to the previous parameters and achieve minimum loss on the current example. This is done by solving the following problem.

$$w_1^{t+1} \dots w_K^{t+1} = \arg \min_{w_1, \dots, w_K} \frac{1}{2} \sum_{v=1}^K k w_v - w_v^t k^2 + C \sum_{v=1}^K w_v$$

$$s.t: a_t w_{\beta^t} x^t = w_r x^t = 1 \quad r; \quad r = 0; \quad r \in [K] \quad (3)$$

Here  $C > 0$  is a parameter to the trade-off between the two parts of the objective function. The larger the value of  $C$ , the focus is more on minimizing the loss. The resulting parameter update equations are as follows.

$$w_r^{t+1} = w_r^t + \frac{a_t}{i \in S^t} x^t + a_t \mathbb{1}_{f_r = \beta^t} \sum_{i \in S^t} x^t; \quad r = 1 \dots K$$

where

$$i = \begin{cases} \min_{j \in S^t} C; \frac{1}{k x^t k^2} \beta_r + \frac{a_t \beta_{\beta^t}}{a_t^2 j S^t j + 1} \frac{a_t^2 \sum_{v \in S^t} \beta_v}{a_t^2 j S^t j + 1} & \beta^t \in S^t \\ \min_{j \notin S^t} C; \frac{1}{k x^t k^2} \beta_r - \frac{a_t^2 \sum_{v \in S^t} \beta_v}{a_t^2 j S^t j + 1} & \beta^t \notin S^t \end{cases}$$

A complete derivation of the update equation and  $i$  is given in the supplementary file. To determine the support set  $S^t$ , EPABF-I uses a slightly different approach compared to Algorithm 2. EPABF-I uses an iterative process to find the values of  $i$  until all the values converge. Support class algorithm (SCA-I) for EPABF-I is given in Algorithm 3.



### 4.3. EPABF-II Updates

EPABF-II algorithm also makes a trade-o between the closeness to the current weight vector and aggressiveness to minimize the loss on the current instance. But in EPABF-II we minimize square of the loss while considering the trade-o .

$$w_1^{t+1} \dots w_K^{t+1} = \arg \min_{w_1 \dots w_K} \frac{1}{2} \sum_{v=1}^K k w_v w_v^t k^2 + C \sum_{v=1}^K \frac{1}{2} w_v^2$$

$$s.t: a_t w_{\varphi^t} x^t w_r x^t = 1 \quad r; r \in [K]$$

---

#### Algorithm 3 Support Class Set Finding Algorithm for EPABF-I (SCA-I)

---

```

1: Input:  $\theta_1, \dots, \theta_K; a_t$ 
2: Initialize:  $S^t := \emptyset$ 
3: Sort  $\theta_1, \dots, \theta_K$  in descending order. Let  $(1) \dots (K)$  be the sorted order.
4: while  $\theta_1, \dots, \theta_K$  do not converge do
5:   for  $r = 1 \dots K$  do
6:     if  $\varphi^t \notin S^t$  then
7:        $w_{(r)}^t = \min \left( C; \frac{1}{kx^t k^2} \theta_{(r)} + \frac{a_t \theta_{\varphi^t}}{a_t^2 j S^t j + 1} \frac{a_t^2 P_{v \in S^t} \theta_v}{a_t^2 j S^t j + 1} \right)$ 
8:     else if  $\varphi^t \in S^t$  then
9:        $w_{(r)}^t = \min \left( C; \frac{1}{kx^t k^2} \theta_{(r)} + \frac{a_t^2 P_{v \in S^t} \theta_v}{a_t^2 j S^t j + 1} \right)$ 
10:    end if
11:    if  $w_{(r)}^t > 0$  then
12:       $S^t = S^t \cup \{(r)\}$ 
13:    else
14:       $S^t = S^t \setminus \{(r)\}$ 
15:    end if
16:  end for
17: end while
18: Output:  $S^t$ 

```

---

#### Algorithm 4 Support Class Set Finding Algorithm for EPABF-II (SCA-II)

---

```

1: Input:  $\theta_1, \dots, \theta_K; a_t$ 
2: Initialize:  $S^t := \emptyset$ 
3: Sort  $\theta_1, \dots, \theta_K$  in descending order. Let  $(1) \dots (K)$  be the sorted order.
4: for  $j = 1 \dots K$  do
5:   if  $\varphi^t \notin S^t$  then
6:     if  $\frac{a_t^2 P_{k \in S^t} \theta_{(k)}}{a_t^2 j S^t j + 1} + \frac{1}{2C x^t k^2} \theta_{(j)} < \theta_{(j)}$ 
7:        $S^t = S^t \cup \{(j)\}$ 
8:     end if
9:   else if  $\varphi^t \in S^t$  then
10:    if  $\frac{a_t^2 P_{k \in S^t} \theta_{(k)}}{a_t^2 j S^t j + 1} + \frac{a_t \theta_{\varphi^t}}{a_t + \frac{1}{2C x^t k^2}} < \theta_{(j)}$ 
11:      then
12:         $S^t = S^t \cup \{(j)\}$ 
13:      end if
14:    end for
15: Output:  $S^t$ 

```

Parameter  $C > 0$  is used to make trade o between the two parts of the objective function. For larger value of  $C$ , more emphasis is given to minimize the square of the loss. The resulting parameter update equations are as follows.

$$w_r^{t+1} = w_r^t + \frac{1}{\sum_{i \in S^t} x_i^t} (x_r^t + a_t I_{f_r = \varphi^t} g_i^t); \quad r = 1 \dots K$$

where

$$\rho_r^t = \frac{\frac{1}{kx^t k^2} + \frac{1}{2Ckx^t k^2}}{\frac{1}{kx^t k^2} + \frac{1}{2Ckx^t k^2}} \left( \frac{B_r \rho_r^t}{a_t^2 j^{2j} + 1} + \frac{a_t \rho_r^t}{a_t + \frac{1}{2Ckx^t k^2}} \right) + \frac{a_t^2 \prod_{j=2}^{S^t} \rho_j^t}{a_t + \frac{1}{2Ckx^t k^2}} \frac{C}{A} \rho_r^t \geq S^t$$

$$\rho_r^t = \frac{1}{kx^t k^2} + \frac{1}{2Ckx^t k^2} \left( \frac{B_r \rho_r^t}{a_t^2 j^{2j} + 1} + \frac{a_t \rho_r^t}{a_t + \frac{1}{2Ckx^t k^2}} \right) + \frac{a_t^2 \prod_{j=2}^{S^t} \rho_j^t}{a_t^2 j^{2j} + 1 + \frac{1}{2Ckx^t k^2}} \frac{C}{A} \rho_r^t \geq S^t$$

The derivations of the step size is given in supplementary le. Using this step size, we get following method to determine the support class.

Determining Support Class Set  $S^t$  Following theorem states the necessary and sufficient condition for any class belonging to the support class set.

Theorem 2 Assume that at trial  $t$ ,  $S^t \in \mathcal{S}$ . Let  $(k)$  be the  $k$ -th class in the sorted order of  $\rho_1; \dots; \rho_K$ . The necessary and sufficient condition for  $(k)$  belongs to  $S^t$  is as follows.

$$\rho_{(k)}^t > \frac{\frac{a_t^2 \prod_{j=1}^{k-1} \rho_{(j)}^t}{a_t^2 (k-1) + 1} + \frac{a_t \rho_{(k)}^t}{a_t + \frac{1}{2Ckx^t k^2}}}{\frac{a_t^2 \prod_{j=1}^{k-1} \rho_{(j)}^t}{a_t^2 (k-1) + 1 + \frac{1}{2Ckx^t k^2}}} \rho_{(k)}^t \geq S^t$$

The proof of this theorem can be found in the supplementary le.

### 5. Mistake Bound Analysis

Here, we find the mistake bounds for the proposed algorithms EPABF, EPABF-I and EPABF-II. Note that,  $E[\sum_{v=1}^K \rho_v^t] \leq \sum_{v=1}^K I_v^t$  and  $E[\sum_{v=1}^K I_v^t] \leq \sum_{v=1}^K I_{v, \text{hinge}}^t \leq \sum_{v=1}^K I_{v, 0}^t$ . Thus,  $\sum_{t=1}^T E[\sum_{v=1}^K \rho_v^t]$  upper bounds the number of mis-classifications over a sequence of  $T$  examples. Similarly,  $E[\sum_{v=1}^K (I_v^t)^2] \leq \sum_{v=1}^K (E[\rho_v^t])^2 + \sum_{v=1}^K (I_v^t)^2 \leq \sum_{v=1}^K I_{v, 0}^t$ . Thus,  $\sum_{t=1}^T E[\sum_{v=1}^K (I_v^t)^2]$  also upper bounds the number of mis-classifications over a sequence of  $T$  examples.

Theorem 3 (Mistake Bound of EPABF) Let  $x^1; \dots; x^T$  be the sequence of examples presented to EPABF. Let  $u_1; \dots; u_K$  be the parameters of an arbitrary linear classifier. Let  $\rho_r^t = [1 - a_t u_{\rho_r^t} x^t + u_r x^t]_+$ ;  $r \in [K]$ . Then, EPABF algorithm satisfies following mistake bound.

$$\sum_{t=1}^T E[\sum_{v=1}^K (\rho_v^t)^2] \leq R \sum_{v=1}^V \frac{1}{ku_v k^2 + 4} \sum_{t=1}^T E[\sum_{v=1}^K (\rho_v^t)^2]$$

Where  $R = \max_{t \in [T]} kx^t k_2$  and  $\sum_{v=1}^V = \frac{K^3}{2} + 1$ .

The proof of the theorem can be found in supplementary le. The regret bound increases with  $K$  (number of classes). In the ideal case ( $l_v^t = 0; \forall v \in [K]; \delta t \in [T]$ ), the regret bound reduces to  $R \leq \frac{1}{2} \sum_{v=1}^K ku_v k^2$ , which is similar to the regret achieved in the full information case (Matsushima et al., 2010).

**Theorem 4 (Mistake Bound of EPABF-I)** Let  $x^1 :::: x^T$  be the sequence of examples presented to EPABF-I. Let  $u_1; ::::; u_K$  be the parameters of an arbitrary linear classifier. Let  $\hat{e}_r^t = [1 - a_t u_{\hat{y}^t} \cdot x^t + u_r \cdot x^t]_+$ ;  $r \in [K]$ . Let  $K$  be the total number of classes and  $R = \max_{t \in [T]} \|x^t\|_2$ . Then, EPABF-I achieves the following mistake bound.

$$\frac{1}{R^2} \sum_{t=1}^T \sum_{v=1}^K \langle \hat{e}_v^t, x^t \rangle^2 + \frac{1}{2} \sum_{v=1}^K \frac{u_v^2}{ku_v k^2} \leq TKR^2 + \frac{TK^4 R^2}{2} + T$$

The proof of the theorem can be found in supplementary le. The general regret bound is of order  $O(\sqrt{T})$  which is at par with the majority of the bandit feedback algorithms. The regret bound increases with increase in number of classes. In the ideal case ( $l_v^t = 0; \forall v \in [K]; \delta t \in [T]$ ), the bound reduces to  $\frac{1}{2} \sum_{v=1}^K ku_v k^2$ .

**Theorem 5 (Mistake Bound of EPABF-II)** Let  $x^1 :::: x^T$  be the sequence of examples presented to EPA. Let  $u_1; ::::; u_K$  be the parameters of an arbitrary linear classifier. Let  $\hat{e}_r^t = [1 - a_t u_{\hat{y}^t} \cdot x^t + u_r \cdot x^t]_+$ ;  $r \in [K]$  and  $R = \max_{t \in [T]} \|x^t\|_2$ . Then, EPABF-II algorithm satisfies following mistake bound.

$$\sum_{t=1}^T \sum_{v=1}^K \langle \hat{e}_v^t, x^t \rangle^2 \leq M \sum_{v=1}^K ku_v k^2 + 2CR^2 \sum_{t=1}^T \sum_{v=1}^K \langle \hat{e}_v^t, x^t \rangle^2$$

where  $M = \frac{R^2 + \frac{1}{2C}}{2K + \frac{1}{C}}$

The proof of the theorem can be found in supplementary le. In the ideal case when  $l_v^t = 0; \forall v \in [K]; \delta t \in [T]$ , the bound reduces to  $M \sum_{v=1}^K ku_v k^2$ .

## 6. Experiments

In this section, we provide experimental results to show the effectiveness of the proposed approach on synthetic and UCI datasets (Dua and Gra (2017)).

**Datasets Used** We report experimental results on two synthetic datasets and 5 UCI datasets (Dua and Gra (2017)). Synthetic datasets are generated as follows.

1. SynSep: It is a 10-class, 100-dimensional data set of 1,00,000 examples generated as follows. We randomly generate vectors from  $[ -1; 1 ]^{100}$ . We use  $u_1; \dots; u_{10} \in \mathbb{R}^{100}$  as the discriminant vectors for 10 classes. All the entries of  $u_i$  are 1 except for indices  $i + 10j$ ;  $j = 0 \dots 9$ . Entries at indices  $i + 10j$ ; ( $j = 0 \dots 9$ ) are -1. Class label for a given example  $x$  is assigned as  $y_x = \arg \max_{i \in \{1 \dots 10\}} (u_i \cdot x)$ . This data is linearly separable.
2. SynNonSep: This dataset is constructed in the same way as SynSep except that we introduce 5% label noise. This makes the data set linearly non-separable.

The UCI datasets (Dua and Gra, 2017) used are Ecoli, Abalone, Satimage, Iris and USPS. Abalone dataset is not balanced as there is asymmetry in its class distribution, so we divided the Rings attribute into 4 intervals as 1-7, 8-14, 15-21, 22-29. Thus, we reduce it to 4 class classification problem.

**Benchmark Algorithms:** We present experimental comparisons of the proposed algorithms (EPABF, EPABF-I, and EPABF-II) with Banditron (Sham M. Kakade (2008)), Bandit Passive-Aggressive (BPA) (Zhong and Dauce (2015)). We also compare the proposed approaches with PA, PA-I, and PA-II algorithms (Matsushima et al. (2010)), which are full information based algorithms. This means PA, PA-I, and PA-II use exact class labels for training. We do this to see how much performance degradation happens due to bandit feedback.

Experimental results are shown in Figures 4. For every dataset, we have ran each algorithm for 100,000 iterations. We have plotted the error curves by averaging the error rates (ratio of incorrectly classified and number of rounds) over 100 different runs. The final plots for each dataset have the average instantaneous error rate on the Y-axis and the number of trials on the X-axis.

(a) (b) (c)

Figure 2: Converged Error Rates versus  $\epsilon$  and  $C$  for Iris Dataset. Note that  $\epsilon$  and  $C$  are on the log scale.

**Choosing Optimal  $\gamma$  and  $C$  for the Proposed Algorithms:** EPABF algorithm takes  $\gamma$  is a user-defined parameter. Similarly, EPABF-I and EPABF-II require  $\gamma$  and  $C$  values as input. To choose the best values of these hyper-parameters, we use the following approach. We explain it for Iris dataset. Figure 2 shows the trend of converged error rate with  $\gamma$  and  $C$  for EPABF algorithms for Iris dataset. We have chosen the best value of  $\gamma$  and  $C$  for which the minimum error rate is achieved. We found the optimal values of  $\gamma$  and  $C$  for all the datasets using the same method. The optimal values of  $C$  and  $\gamma$  for different datasets are shown below in the Tables 1, 2 and 3. The final results shown are based on taking these optimal hyper-parameters values. We use the same approach for selecting the hyper-parameters for the benchmark algorithms also.

Table 1: EPABF		Table 2: EPABF-I			Table 3: EPABF-II		
Dataset		Dataset	$C$		Dataset	$C$	
SynSep	0.0001	SynSep	0.1	0.0003	SynSep	0.1	0.0001
SynNonSep	0.001	SynNonSep	0.1	0.001	SynNonSep	0.1	0.0001
Ecoli	0.04	Ecoli	0.01	0.05	Ecoli	0.01	0.04
Abalone	0.008	Abalone	0.01	0.02	Abalone	0.01	0.03
Satimage	0.065	Satimage	0.1	0.06	Satimage	0.1	0.06
Iris	0.035	Iris	0.03	0.025	Iris	0.03	0.04
USPS	0.06	USPS	0.1	0.1	USPS	0.1	0.1

**Comparison of EPABF algorithms with Other Bandit Algorithms:** Figure 3 presents comparison results of the proposed EPABF and its variants with other bandit feedback based algorithms (i.e., BPA and Banditron). We observe that the proposed EPABF variants outperform Banditron by a significant margin in terms of converged error rates for all the datasets. Compared to BPA, the proposed algorithms perform significantly better except for the Ecoli, USPS, and Iris datasets. For Ecoli, USPS, and Iris datasets, proposed approaches still perform better than BPA though marginally. For Iris and Satimage, EPABF variants take more time to converge compared to BPA and Banditron but achieve a lower error rate.

We see that from SynSep to SynNonSep, the performances of BPA and Banditron drop by huge amount compared to EPABF variants.

**Comparison with Full Information Case:** Figure 4 presents comparison results of the proposed algorithms with full information algorithms, namely, PA, PA-I, and PA-II. For Ecoli and USPS datasets, the proposed EPABF algorithms perform comparably to the full information counterparts. For Iris, Satimage, and USPS datasets, the full information algorithms perform better than the proposed approach.

On the other hand, for SynSep, SynNonSep, and Abalone datasets, the proposed bandit algorithms perform better than the full information based algorithms. A similar kind of behavior was observed in Zhong and Dauce (2015). This might happen because of the following reason. In the proposed algorithms, we are using an exploration-exploitation scheme. Unlike the full information algorithms, we do not always choose the best-predicted label, which is still selecting the best labels. Thus, over a broad set of examples, the actual

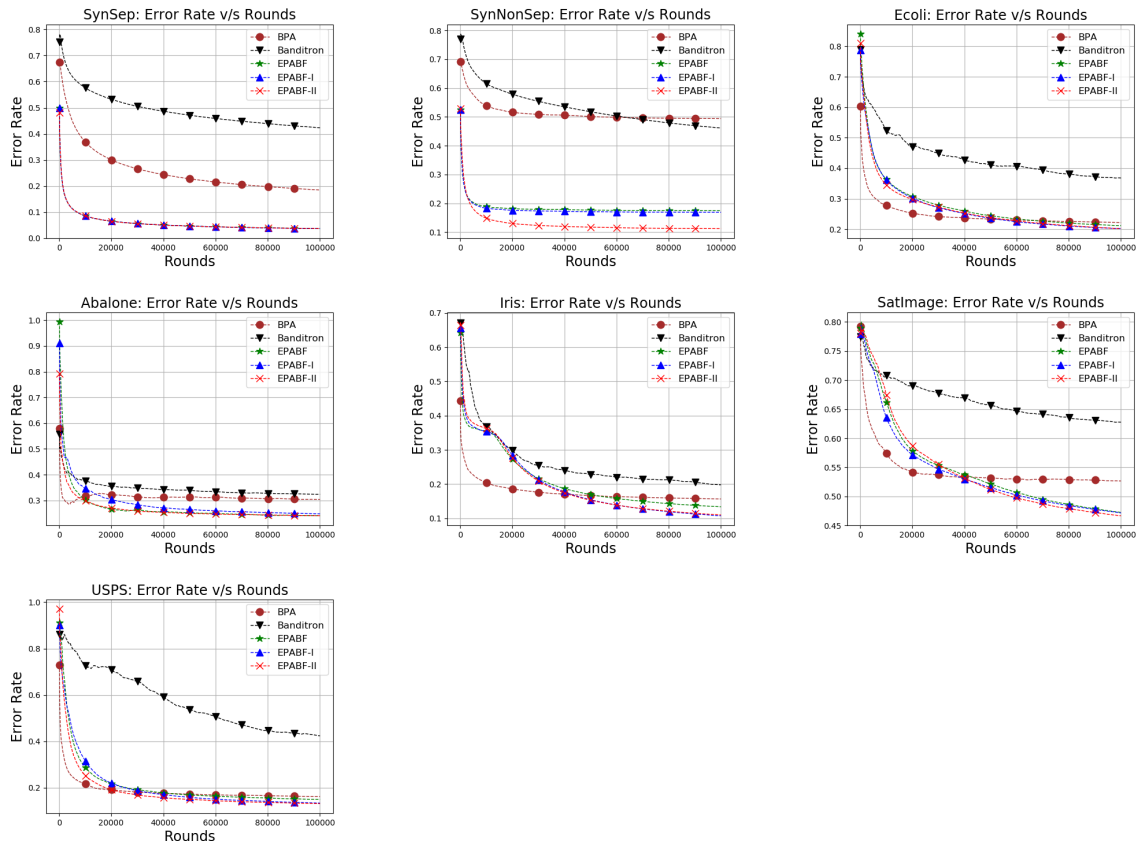


Figure 3: Performance comparison of EPABF, EPABF-I and EPABF-II with Benchmark Bandit Algorithms (BPA and Banditron).

labels match the labels predicted by EPABFs. This can lead to better performances of EPABF algorithms compared to full information based algorithms.

## 7. Conclusion

This paper proposed three variants of passive-aggressive online algorithms for multiclass classification using bandit feedback, namely EPABF, EPABF-I, and EPABF-II. We used exploration-exploitation to guess the best label of an example. To update the parameters, we solve a constrained optimization problem in each trial. The constraints are used to capture the separability conditions of the multiclass problem. Finally, we update the parameter of those classes whose corresponding constraints are active. These classes are called support classes. We also provided regret bounds for all three variants. We also offered an experimental comparison of all the algorithms on various datasets. The results showed that our EPABF algorithms perform better than existing bandit feedback based algorithms and perform comparably to the full feedback based algorithms on most of the given datasets.

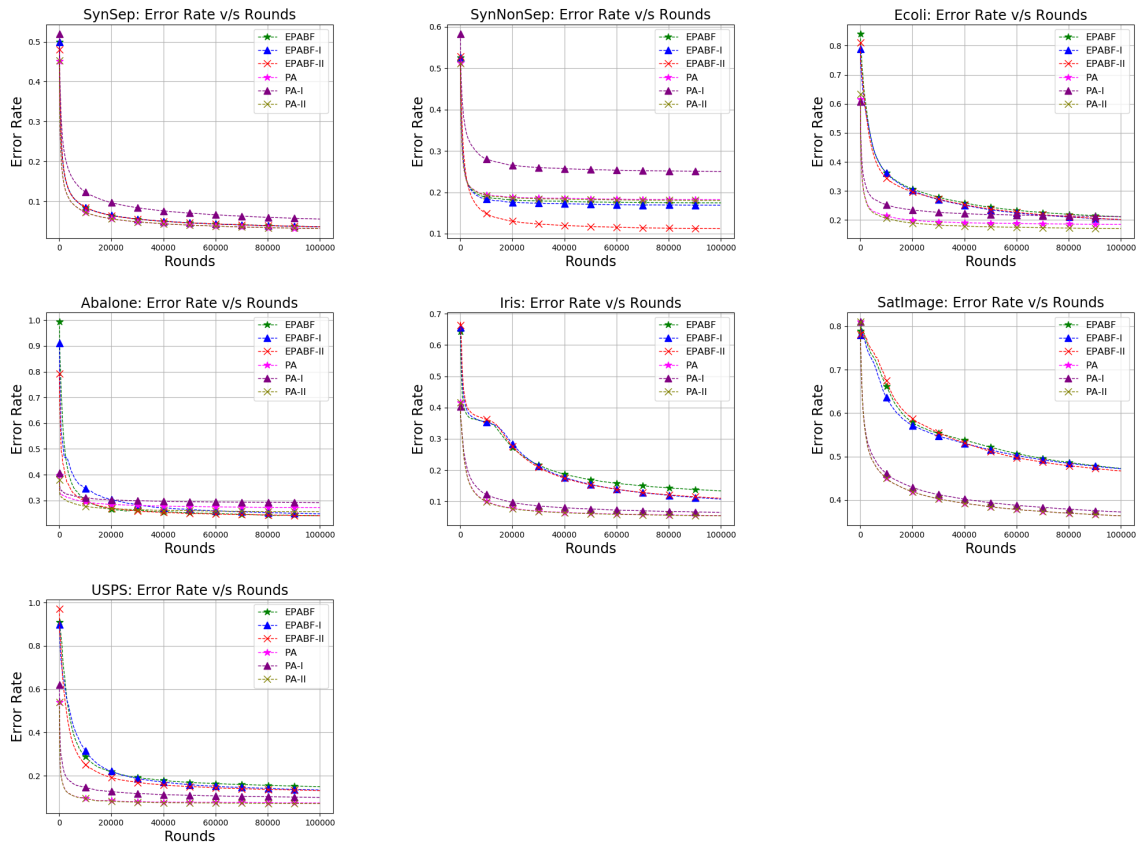


Figure 4: Performance comparison of EPABF, EPABF-I and EPABF-II with full information algorithms (PA, PA-I and PA-II).

References

Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Mach. Learn.*, 47(2-3):235–256, May 2002.

Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM J. Comput.*, 32(1):48–77, January 2003.

Alina Beygelzimer, Francesco Orabona, and Chicheng Zhang. Efficient online bandit multiclass learning with  $\tilde{O}(\sqrt{T})$  regret. 02 2017.

Alina Beygelzimer, David Pal, Balazs Szorenyi, Devanathan Thiruvengatachari, Chen-Yu Wei, and Chicheng Zhang, editors. *Bandit Multiclass Linear Classification: Efficient Algorithms for the Separable Case*, ICML, 2019.

Nicolò Cesa-Bianchi, Alex Conconi, and Claudio Gentile. A second-order perceptron algorithm. *SIAM Journal on Computing*, 34:640–668, January 2005.

