

# Online Learning: A Comprehensive Survey

**Steven C. H. Hoi**

CHHOI@SMU.EDU.SG

*School of Information Systems, Singapore Management University, Singapore*

**Doyen Sahoo**

DOYENS@SMU.EDU.SG

*School of Information Systems, Singapore Management University, Singapore*

**Jing Lu**

LVJING12@JD.COM

*JD.com*

**Peilin Zhao**

MASONZHAO@TENCENT.COM

*Tencent AI Lab*

**Editor:** XYZ

## Abstract

*Online learning* represents a family of machine learning methods, where a learner attempts to tackle some predictive (or any type of decision-making) task by learning from a sequence of data instances one by one at each time. The goal of online learning is to maximize the accuracy/correctness for the sequence of predictions/decisions made by the online learner given the knowledge of correct answers to previous prediction/learning tasks and possibly additional information. This is in contrast to traditional *batch* or *offline machine learning* methods that are often designed to learn a model from the entire training data set at once. Online learning has become a promising technique for learning from continuous streams of data in many real-world applications. This survey aims to provide a comprehensive survey of the online machine learning literature through a systematic review of basic ideas and key principles and a proper categorization of different algorithms and techniques. Generally speaking, according to the types of learning tasks and the forms of feedback information, the existing online learning works can be classified into three major categories: (i) *online supervised learning* where full feedback information is always available, (ii) *online learning with limited feedback*, and (iii) *online unsupervised learning* where no feedback is available. Due to space limitation, the survey will be mainly focused on the first category, but also briefly cover some basics of the other two categories. Finally, we also discuss some open issues and attempt to shed light on potential future research directions in this field.

**Keywords:** Online learning, Online convex optimization, Sequential decision making

## 1. Introduction

Machine learning plays a crucial role in modern data analytics and artificial intelligence (AI) applications. Traditional machine learning paradigms often work in a batch learning or offline learning fashion (especially for supervised learning), where a model is trained by some learning algorithm from an entire training data set at once, and then the model is deployed for inference without (or seldom) performing any update afterwards. Such learning methods suffer from expensive re-training cost when dealing with new training data, and thus are poorly scalable for real-world applications. In the era of big data, traditional batch learning paradigms become more and more restricted, especially when live data grows and

evolves rapidly. Making machine learning scalable and practical especially for learning from continuous data streams has become an open grand challenge in machine learning and AI.

Unlike traditional machine learning, *online learning* is a subfield of machine learning and includes an important family of learning techniques which are devised to learn models incrementally from data in a sequential manner. Online learning overcomes the drawbacks of traditional batch learning in that the model can be updated instantly and efficiently by an online learner when new training data arrives. Besides, online learning algorithms are often easy to understand, simple to implement, and often founded on solid theory with rigorous regret bounds. Along with urgent need of making machine learning practical for real big data analytics, online learning has attracted increasing interest in recent years.

This survey aims to give a comprehensive survey of *online learning* literature. Online learning<sup>1</sup> has been extensively studied across different fields, ranging from machine learning, data mining, statistics, optimization and applied math, to artificial intelligence and data science. This survey aims to distill the core ideas of online learning methodologies and applications in literature. This survey is written mainly for machine learning audiences, and assumes readers with basic knowledge in machine learning. While trying our best to make the survey as comprehensive as possible, it is very difficult to cover every detail since online learning research has been evolving rapidly in recent years. We apologize in advance for any missing papers or inaccuracies in description, and encourage readers to provide feedback, comments or suggestions. Finally, as a supplemental document to this survey, readers may check our updated version online at: <http://libol.stevenhoi.org/survey>.

## 1.1 What is Online Learning?

Traditional machine learning paradigm often runs in a batch learning fashion, e.g., a supervised learning task, where a collection of training data is given in advance to train a model by following some learning algorithm. Such a paradigm requires the entire training data set to be made available prior to the learning task, and the training process is often done in an offline environment due to the expensive training cost. Traditional batch learning methods suffer from some critical drawbacks: (i) low efficiency in both time and space costs; and (ii) poor scalability for large-scale applications because the model often has to be re-trained from scratch for new training data.

In contrast to batch learning algorithms, online learning is a method of machine learning for data arriving in a sequential order, where a learner aims to learn and update the best predictor for future data at every step. Online learning is able to overcome the drawbacks of batch learning in that the predictive model can be updated instantly for any new data instances. Thus, online learning algorithms are far more efficient and scalable for large-scale machine learning tasks in real-world data analytics applications where data are not only large in size, but also arriving at a high velocity.

---

1. The term of “online learning” in this survey is *not* related to “e-learning” in the online education field.

## 1.2 Tasks and Applications

Similar to traditional (batch) machine learning methods, online learning techniques can be applied to solve a variety of tasks in a wide range of real-world application domains. Examples of online learning tasks include the following:

*Supervised learning tasks:* Online learning algorithms can be derived for supervised learning tasks. One of the most common tasks is classification, aiming to predict the categories for a new data instance belongs to, on the basis of observing past training data instances whose category labels are given. For example, a commonly studied task in online learning is online binary classification (e.g., spam email filtering) which only involves two categories (“spam” vs “benign” emails); other types of supervised classification tasks include multi-class classification, multi-label classification, and multiple-instance classification, etc.

In addition to classification tasks, another common supervised learning task is regression analysis, which refers to the learning process for estimating the relationships among variables (typically between a dependent variable and one or more independent variables). Online learning techniques are naturally applied for regression analysis tasks, e.g., time series analysis in financial markets where data instances naturally arrive in a sequential way. Besides, another application for online learning with financial time series data is online portfolio selection where an online learner aims to find a good (e.g., profitable and low-risk) strategy for making a sequence of decisions for portfolio selection.

*Bandit learning tasks:* Bandit online learning algorithms, also known as Multi-armed bandits (MAB), have been extensively used for many online recommender systems, such as online advertising for internet monetization, product recommendation in e-commerce, movie recommendation for entertainment, and other personalized recommendation, etc.

*Unsupervised learning tasks:* Online learning algorithms can be applied for unsupervised learning tasks. Examples include clustering or cluster analysis — a process of grouping objects such that objects in the same group (“cluster”) are more similar to each other than to objects in other clusters. Online clustering aims to perform incremental cluster analysis on a sequence of instances, which is common for mining data streams.

*Other learning tasks:* Online learning can also be used for other kinds of machine learning tasks, such as learning for recommender systems, learning to rank, or reinforcement learning. For example, collaborative filtering with online learning can be applied to enhance the performance of recommender systems by learning to improve collaborative filtering tasks sequentially from continuous streams of ratings/feedback information from users.

Last but not least, we note that online learning techniques are often used in two major scenarios. One is to improve efficiency and scalability of existing machine learning methodologies for batch machine learning tasks where a full collection of training data must be made available before the learning tasks. For example, Support Vector Machines (SVM) is a well-known supervised learning method for batch classification tasks, in which classical SVM algorithms (e.g., QP or SMO solvers (Platt et al., 1999)) suffer from poor scalability for very large-scale applications. In literature, various online learning algorithms have been explored for training SVM in an online (or stochastic) learning manner (Poggio, 2001; Shalev-Shwartz et al., 2011), making it more efficient and scalable than conventional batch SVMs. The other scenario is to apply online learning algorithms to directly tackle online streaming data analytics tasks where data instances naturally arrive in a sequential manner

and the target concepts may be drifting or evolving over time. Examples include time series regression, such as stock price prediction, where data arrives periodically and the learner has to make decisions immediately before getting the next instance.

Online Learning			
Statistical Learning Theory		Convex Optimization Theory	
Game Theory			
Online Learning with Full Feedback		Online Learning with Partial Feedback (Bandits)	
Online Supervised Learning		Stochastic Bandit	Adversarial Bandit
First-order Online Learning	Online Learning with Regularization	Stochastic Multi-armed Bandit	Adversarial Multi-armed Bandit
Second-order Online Learning	Online Learning with Kernels	Bayesian Bandit	Combinatorial Bandit
Prediction with Expert Advice	Online to Batch Conversion	Stochastic Contextual Bandit	Adversarial Contextual Bandit
Applied Online Learning		Online Active Learning	Online Semi-supervised Learning
Cost-Sensitive Online Learning	Online Collaborative Filtering	Selective Sampling	Online Manifold Regularization
Online Multi-task Learning	Online Learning to Rank	Active Learning with Expert Advice	Transductive Online Learning
Online Multi-view Learning	Distributed Online Learning	Online Unsupervised Learning (no feedback)	
Online Transfer Learning	Online Learning with Neural Networks	Online Clustering	Online Density Estimation
Online Metric Learning	Online Portfolio Selection	Online Dimension Reduction	Online Anomaly Detection

Figure 1: Taxonomy of Online Learning Techniques

### 1.3 Taxonomy

To help readers better understand the online learning literature as a whole, we attempt to construct a taxonomy of online learning methods and techniques, as summarized in Figure 1. In general, from a theoretical perspective, online learning methodologies are founded based on theory and principles from three major theory communities: learning theory, optimization theory, and game theory. From the perspective of specific algorithms, we can further group the existing online learning techniques into different categories according to their specific learning principles and problem settings. Specifically, according to the types of feedback information and the types of supervision in the learning tasks, online learning techniques can be classified into the following three major categories:

- **Online supervised learning:** This is concerned with supervised learning tasks where full feedback information is always revealed to a learner at the end of each online learning round. It can be further divided into two groups of studies: (i) “Online Supervised Learning” which forms the fundamental approaches and principles for Online Supervised Learning; and (ii) “Applied Online Learning” which constitute more non-traditional online supervised learning, where the fundamental approaches cannot be directly applied, and algorithms have been appropriately tailored to suit the non-traditional online learning setting.

- **Online learning with limited feedback:** This is concerned with tasks where an online learner receives partial feedback information from the environment during the online learning process. For example, consider an online multi-class classification task, at a particular round, the learner makes a prediction of class label for an incoming instance, and then receives the partial feedback indicating whether the prediction is correct or not, instead of the particular true class label explicitly. For such tasks, the online learner often has to make the online updates or decisions by attempting to achieve some tradeoff between the exploitation of disclosed knowledge and the exploration of unknown information with the environment.
- **Online unsupervised learning:** This is concerned with online learning tasks where the online learner only receives the sequence of data instances without any additional feedback (e.g., true class label) during the online learning tasks. Unsupervised online learning can be considered as a natural extension of traditional unsupervised learning for dealing with data streams, which is typically studied in batch learning fashion. Examples of unsupervised online learning include online clustering, online dimension reduction, and online anomaly detection tasks, etc. Unsupervised online learning has less restricted assumptions about data without requiring explicit feedback or label information which could be difficult or expensive to acquire.

This article will conduct a systematic review of existing work for online learning, especially for online supervised learning and online learning with partial feedback. Finally, we note that it is always very challenging to make a precise categorization of all the existing online learning work, and it is likely that the above proposed taxonomy may not fully cover all the existing online learning work in literature, though we have tried our best to cover as much as possible.

#### 1.4 Related Work and Further Reading

This paper attempts to make a comprehensive survey of online learning research work. In literature, there are some related books, PHD theses, and articles published over the past years dedicated to online learning (Cesa-Bianchi and Lugosi, 2006; Shalev-Shwartz, 2011), in which many of them also include rich discussions on related work on online learning. For example, the book titled “Prediction, Learning, and Games” (Cesa-Bianchi and Lugosi, 2006) gave a nice introduction about some niche subjects of online learning, particularly for online prediction with expert advice and online learning with partial feedback. Another work titled “Online Learning and Online Convex Optimization” (Shalev-Shwartz, 2011) gave a nice tutorial about basics of online learning and foundations of online convex optimization. In addition, there are also quite a few PHD theses dedicated to addressing different subjects of online learning (Kleinberg, 2005b; Shalev-Shwartz, 2007; Zhao, 2013; Li, 2013). Readers are also encouraged to read some older related books, surveys and tutorial notes about online learning and online algorithms (Fiat and Woeginger, 1998; Bottou, 1998b; Rakhlin, 2008; Blum, 1998; Albers, 2003). Finally, readers who are interested in applied online learning can explore some open-source toolboxes, including LIBOL (Hoi et al., 2014; Wu et al., 2017a) and Vowpal Wabbit (Langford et al., 2007).

## 2. Problem Formulations and Related Theory

Without loss of generality, we first give a formal formulation of a classic online learning problem, i.e., online binary classification, and then introduce basics of statistical learning theory, online convex optimization and game theory as the theoretical foundations for online learning techniques.

### 2.1 Problem Settings

Consider an online binary classification task, online learning takes place in a sequential way. On each round, a learner receives a data instance, and then makes a prediction of the instance, e.g., classifying it into some predefined categories. After making the prediction, the learner receives the true answer about the instance from the environment as a feedback. Based on the feedback, the learner can measure the loss suffered, depending on the difference between the prediction and the answer. Finally, the learner updates its prediction model by some strategy so as to improve predictive performance on future received instances.

Consider spam email detection as a running example of online binary classification, where the learner answers every question in binary: yes or no. The task is supervised binary classification from a machine learning perspective. More formally, we can formulate the problem as follows: consider a sequence of instances/objects represented in a vector space,  $\mathbf{x}_t \in \mathbb{R}^d$ , where  $t$  denotes the  $t$ -th round and  $d$  is the dimensionality, and we use  $y_t \in \{+1, -1\}$  to denote the true class label of the instance. The online binary classification takes place sequentially. On the  $t$ -th round, an instance  $\mathbf{x}_t$  is received by the learner, which then employs a binary classifier  $\mathbf{w}_t$  to make a prediction on  $\mathbf{x}_t$ , e.g.,  $\hat{y}_t = \text{sign}(\mathbf{w}_t^\top \mathbf{x}_t)$  that outputs  $\hat{y}_t = +1$  if  $\mathbf{w}_t^\top \mathbf{x}_t \geq 0$  and  $\hat{y}_t = -1$  otherwise. After making the prediction, the learner receives the true class label  $y_t$  and thus can measure the suffered loss, e.g., using the hinge-loss  $\ell_t(\mathbf{w}_t) = \max(0, 1 - y_t \mathbf{w}_t^\top \mathbf{x}_t)$ . Whenever the loss is nonzero, the learner updates the prediction model from  $\mathbf{w}_t$  to  $\mathbf{w}_{t+1}$  by applying some strategy on the training example  $(\mathbf{x}_t, y_t)$ . The procedure of Online Binary Classification is summarized in Algorithm 1.

---

**Algorithm 1:** Online Binary Classification process.

---

```

Initialize the prediction function as  $\mathbf{w}_1$ ;
for  $t = 1, 2, \dots, T$  do
    Receive instance:  $\mathbf{x}_t \in \mathbb{R}^d$ ;
    Predict  $\hat{y}_t = \text{sign}(\mathbf{w}_t^\top \mathbf{x}_t)$  as the label of  $\mathbf{x}_t$ ;
    Receive the true class label:  $y_t \in \{-1, +1\}$ ;
    Suffer loss:  $\ell_t(\mathbf{w}_t)$  which is a convex loss function on both  $\mathbf{w}_t^\top \mathbf{x}_t$  and  $y_t$ ;
    Update the prediction model  $\mathbf{w}_t$  to  $\mathbf{w}_{t+1}$ ;
end for
    
```

---

By running online learning over a sequence of  $T$  rounds, the number of mistakes made by the online learner can be measured as  $M_T = \sum_{t=1}^T \mathbb{I}(\hat{y}_t \neq y_t)$ . In general, the classic goal of an online learning task is to minimize the regret of the online learner's predictions

against the best fixed model in hindsight, defined as

$$R_T = \sum_{t=1}^T \ell_t(\mathbf{w}_t) - \min_{\mathbf{w}} \sum_{t=1}^T \ell_t(\mathbf{w}) \quad (1)$$

where the second term is the loss suffered by the optimal model  $\mathbf{w}^*$  that can only be known in hindsight after seeing all the instances and their class labels. From the theoretical perspective of regret minimization, if an online algorithm guarantees that its regret is sublinear as a function of  $T$ , i.e.,  $R_T = o(T)$ , it implies that  $\lim_{T \rightarrow \infty} R(T)/T = 0$  and thus on average the learner performs almost as well as the best fixed model in hindsight.

## 2.2 Statistical Learning Theory

Statistical learning theory, first introduced in the late 1960's, is one of key foundations for theoretical analysis of machine learning problems, especially for supervised learning. In literature, there are many comprehensive survey articles and books (Vapnik and Vapnik, 1998; Vapnik, 1999). In the following, we introduce some basic concept and framework.

### 2.2.1 EMPIRICAL ERROR MINIMIZATION

Assume instance  $\mathbf{x}_t$  is generated randomly from a fixed but unknown distribution  $P(\mathbf{x})$  and its class label  $y$  is also generated with a fixed but unknown distribution  $P(y|\mathbf{x})$ . The joint distribution of labeled data is  $P(\mathbf{x}, y) = P(\mathbf{x})P(y|\mathbf{x})$ . The goal of a learning problem is to find a prediction function  $f(\mathbf{x})$  that minimizes the expected value of the loss function:

$$R(f) = \int \ell(y, f(\mathbf{x})) dP(x, y)$$

which is also termed as the *True Risk* function. The solution  $f^* = \arg \min R(f)$  is the optimal predictor. In general, the true risk function  $R(f)$  cannot be computed directly because of the unknown distribution  $P(x, y)$ . In practice, we approximate the true risk by estimating the risk over a finite collection of instances  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$  drawn i.i.d., which is called the “Empirical Risk” or “Empirical Error”

$$R_{emp}(f) = \frac{1}{N} \sum_{n=1}^N \ell(y_n, f(\mathbf{x}_n))$$

The problem of learning via the Empirical Error Minimization (ERM) is to find a hypothesis  $f$  over a hypothesis space  $\mathcal{F}$  by minimizing the Empirical Error:

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} R_{emp}(f)$$

ERM is the theoretical base for many machine learning algorithms. For example, in the problem of binary classification, when assuming  $\mathcal{F}$  is the set of linear classifiers and the hinge loss is used, the ERM principle indicates that the best linear model  $\mathbf{w}$  can be trained by minimizing the following objective

$$R_{emp}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \max(0, 1 - y_n \mathbf{w}^\top \mathbf{x}_n)$$

### 2.2.2 ERROR DECOMPOSITION

The difference between the optimal predictor  $f^*$  and the empirical best predictor  $\hat{f}_n$  can be measured by the Excess Risk, which can be decomposed as follows:

$$R(\hat{f}_n) - R(f^*) = \left( R(\hat{f}_n) - \inf_{f \in \mathcal{F}} R(f) \right) + \left( \inf_{f \in \mathcal{F}} R(f) - R(f^*) \right)$$

where the first term is called the Estimation Error due to the finite amount of training samples that may not be enough to represent the unknown distribution, and the second term is called the Approximation Error due to the restriction of model class  $\mathcal{F}$  that may not be flexible enough to include the optimal predictor  $f^*$ . In general, the estimation error will be reduced when increasing the amount of training data, while the approximation error can be reduced by increasing the model complexity/capacity. However, the estimation error often increases when the model complexity grows, making it challenging for model selection.

### 2.3 Convex Optimization Theory

Many online learning problems can essentially be (re-)formulated as an Online Convex Optimization (OCO) task. In the following, we introduce some basics of OCO.

An online convex optimization task typically consists of two major elements: a convex set  $\mathcal{S}$  and a convex cost function  $\ell_t(\cdot)$ . At each time step  $t$ , the online algorithm decides to choose a weight vector  $\mathbf{w}_t \in \mathcal{S}$ ; after that, it suffers a loss  $\ell_t(\mathbf{w}_t)$ , which is computed based on a convex cost function  $\ell_t(\cdot)$  defined over  $\mathcal{S}$ . The goal of the online algorithm is to choose a sequence of decisions  $\mathbf{w}_1, \mathbf{w}_2, \dots$  such that the regret in hindsight can be minimized.

More formally, an online algorithm aims to achieve a low regret  $R_T$  after  $T$  rounds, where the regret  $R_T$  is defined as:

$$R_T = \sum_{t=1}^T \ell_t(\mathbf{w}_t) - \inf_{\mathbf{w}^* \in \mathcal{S}} \sum_{t=1}^T \ell_t(\mathbf{w}^*), \quad (2)$$

where  $\mathbf{w}^*$  is the solution that minimizes the convex objective function  $\sum_{t=1}^T \ell_t(\mathbf{w})$  over  $\mathcal{S}$ .

For example, consider an online binary classification task for training online Support Vector Machines (SVM) from a sequence of labeled instances  $(\mathbf{x}_t, y_t), t = 1, \dots, T$ , where  $\mathbf{x}_t \in \mathcal{R}^d$  and  $y_t \in \{+1, -1\}$ . One can define the loss function  $\ell(\cdot)$  as  $\ell_t(\mathbf{w}_t) = \max(0, 1 - y_t \mathbf{w}_t^\top \mathbf{x}_t)$  and the convex set  $\mathcal{S}$  as  $\{\mathbf{w} \in \mathcal{R}^d \mid \|\mathbf{w}\| \leq C\}$  for some constant parameter  $C$ . There are a variety of algorithms to solve this problem.

For a comprehensive treatment of this subject, readers are referred to the books in (Shalev-Shwartz, 2011; Hazan et al., 2016). Below we briefly review three major families of online convex optimization (OCO) methods, including first-order algorithms, second-order algorithms, and regularization based approaches.

#### 2.3.1 FIRST-ORDER METHODS

First order methods aim to optimize the objective function using the first order (sub) gradient information. Online Gradient Descent (OGD)(Zinkevich, 2003) can be viewed as an online version of Stochastic Gradient Descent (SGD) in convex optimization, and is one of the simplest and most popular methods for convex optimization.



At every iteration, based on the loss suffered on instance  $\mathbf{x}_t$ , the algorithm takes a step from the current model to update to a new model, in the direction of the gradient of the current loss function. This update gives us  $\mathbf{u} = \mathbf{w}_t - \eta_t \nabla \ell_t(\mathbf{w}_t)$ . The resulting update may push the model to lie outside the feasible domain. Thus, the algorithm projects the model onto the feasible domain, i.e.,  $\Pi_{\mathcal{S}}(\mathbf{u}) = \arg \min_{\mathbf{w} \in \mathcal{S}} \|\mathbf{w} - \mathbf{u}\|$  (where  $\Pi_{\mathcal{S}}$  denotes the projection operation). OGD is simple and easy to implement, but the projection step sometimes may be computationally intensive which depends on specific tasks. In theory (Zinkevich, 2003), OGD achieves sublinear regret  $O(\sqrt{T})$  for an arbitrary sequence of  $T$  convex cost functions (of bounded gradients), with respect to the best single decision in hindsight.

### 2.3.2 SECOND-ORDER METHODS.

Second-order methods aim to exploit second order information to speed up the convergence of the optimization. A popular approach is the Online Newton Step Algorithm. The Online Newton Step (Hazan et al., 2007a) can be viewed as an online analogue of the Newton-Raphson method in batch optimization. Like OGD, ONS also performs an update by subtracting a vector from the current model in each online iteration. While the vector subtracted by OGD is the gradient of the current loss function based on the current model, in ONS the subtracted vector is the inverse Hessian multiplied by the gradient, i.e.,  $A_t^{-1} \nabla \ell_t(\mathbf{w}_t)$  where  $A_t$  is related to the Hessian.  $A_t$  is also updated in each iteration as  $A_t = A_{t-1} + \nabla \ell_t(\mathbf{w}_t) \nabla \ell_t(\mathbf{w}_t)^\top$ . The updated model is projected back to the feasible domain as  $\mathbf{w}_{t+1} = \Pi_{\mathcal{S}}^{A_t}(\mathbf{w}_t - \eta A_t^{-1} \nabla \ell_t(\mathbf{w}_t))$ , where  $\Pi_{\mathcal{S}}^A(\mathbf{u}) = \arg \min_{\mathbf{w} \in \mathcal{S}} (\mathbf{w} - \mathbf{u})^\top A (\mathbf{w} - \mathbf{u})$ . Different from OGD where the projection is made under the Euclidean norm, ONS projects under the norm induced by the matrix  $A_t$ . Although ONS's time complexity  $O(n^2)$  is higher than OGD's  $O(n)$ , it guarantees a logarithmic regret  $O(\log T)$  under relatively weaker assumptions of exp-concave cost functions.

### 2.3.3 REGULARIZATION

Unlike traditional convex optimization, the aim of OCO is to optimize the regret. Traditional approaches (termed as Follow the Leader (FTL)) can be unstable, leading to high regret (e.g. linear regret) in the worst case (Hazan et al., 2016). This motivates the need to stabilize the approaches through regularization. Here we discuss the common regularization approaches.

*Follow-the-Regularized-Leader (FTRL)*. The idea of Follow-the-Regularized-Leader (FTRL) (Abernethy et al., 2008; Shalev-Shwartz and Singer, 2007) is to stabilize the prediction of the Follow-the-Leader (FTL) (Kalai and Vempala, 2005; Hannan, 1957) by adding a regularization term  $R(\mathbf{w})$  which is strongly convex, smooth and twice differentiable. The idea is to solve the following optimization problem in each iteration:

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in \mathcal{S}} \left[ \eta \sum_{s=1}^t \nabla \ell_s(\mathbf{w}_s)^\top \mathbf{w} + R(\mathbf{w}) \right]$$

where  $\mathcal{S}$  is the feasible convex set and  $\eta$  is the learning rate. In theory, the FTRL algorithm in general achieves a sublinear regret bound  $O(\sqrt{T})$ .

*Online Mirror Descent (OMD)*. OMD is an online version of the Mirror Descent (MD) method (Tseng, 2008; Duchi et al., 2010) in batch convex optimization. The OMD algorithm

behaves like OGD, in that it updates the model using a simple gradient rule. However, it generalizes OGD as it performs updates in the dual space. This duality is induced by the choice of the regularizer: the gradient of the regularization serves as a mapping from  $\mathbb{R}^d$  to itself. Due to this transformation by the regularizer, OMD is able to obtain better bounds in terms of the geometry of the space.

In general, OMD has two variants of algorithms: lazy OMD and active OMD. The lazy version keeps track of a point in Euclidean space and projects it onto the convex feasible domain only when making prediction, while the active version keeps a feasible model all the time, which is a direct generalization of OGD. Unlike OGD, the projection step in OMD is based on the Bregman Divergence  $\mathcal{B}_R$ , i.e.,  $\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in \mathcal{S}} \mathcal{B}_R(\mathbf{w} \parallel \mathbf{v}_{t+1})$ , where  $\mathbf{v}_{t+1}$  is the updated model after the gradient step. In general, the lazy OMD has the same regret bound as FTRL. The active OMD also has a similar regret bound. When  $R(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_2^2$ , OMD recovers OGD. If we use other functions as  $R$ , we can also recover some other interesting algorithms, such as the Exponential Gradient (EG) algorithm below.

*Exponential Gradient (EG).* Let  $R(\mathbf{w}) = \mathbf{w} \ln \mathbf{w}$  be the negative entropy function and the feasible convex domain be the simplex  $\mathcal{S} = \Delta_d = \{\mathbf{w} \in \mathbb{R}_+^d \mid \sum_i w_i = 1\}$ , then OMD will recover the Exponential Gradient (EG) algorithm (Kivinen and Warmuth, 1995). In this special case, the induced projection is the normalization by the  $L1$  norm, which indicates

$$w_{t+1,i} = \frac{w_{t,i} \exp[-\eta(\nabla \ell_t(\mathbf{w}_t))_i]}{\sum_j w_{t,j} \exp[-\eta(\nabla \ell_t(\mathbf{w}_t))_j]}$$

As a special case of OMD, the regret of EG is bounded by  $O(\sqrt{T})$ .

*Adaptive (Sub)-Gradient Methods.* In the previous algorithms, the regularization function  $R$  is always fixed and data independent, during the whole learning process. Adaptive (Sub)-Gradient (AdaGrad) algorithm (Duchi et al., 2011b) is an algorithm that can be considered as online mirror descent with adaptive regularization, i.e., the regularization function  $R$  can change over time. The regularizer  $R$  at the  $t$ -th step, is actually the function  $R(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_{A_t^{1/2}}^2 = \frac{1}{2} \mathbf{w}^\top A_t^{1/2} \mathbf{w}$ , which is constructed from the (sub)-gradients received before (and including) the  $t$ -th step. In each iteration the model is updated as:

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in \mathcal{S}} \left\| \mathbf{w} - [\mathbf{w}_t - \eta A_t^{-\frac{1}{2}} \nabla \ell_t(\mathbf{w}_t)] \right\|_{A_t^{\frac{1}{2}}}^2$$

where  $A_t$  is updated as:

$$A_t = A_{t-1} + \nabla \ell_t(\mathbf{w}_t) \nabla \ell_t(\mathbf{w}_t)^\top$$

We also note that there are also other emerging online convex optimization methods, such as Online Convex Optimization with long term constraints (Jenatton et al., 2016), which assumes that the constraints are only required to be satisfied in long term, and Online ADMM (Wang and Banerjee, 2012) which is an online version for the Alternating Direction Method of Multipliers (ADMM) (Gabay and Mercier, 1976; Boyd et al., 2011) and is particularly suitable for distributed optimization applications. The RESCALED-EXP algorithm (Cutkosky and Boahen, 2016), proposed recently, does not use any prior knowledge about the loss functions and does not require the tuning of learning rate.

## 2.4 Game Theory

Game theory is closely related to online learning. In general, an online prediction task can be formulated as a problem of learning to play a repeated game between a learner and an environment (Freund and Schapire, 1999b). Consider online classification as an example, during each iteration, the algorithm chooses one class from a finite number of classes and the environment reveals the true class label. Assume the environment is stable (e.g., i.i.d), i.e., not played by an adversary. The algorithm aims to perform as well as the best fixed strategy. The classic online classification problem thus can be modeled by the game theory under the simplest assumption, full feedback and a stable environment. More generally, various settings in game theory can be related to many other types of online learning problems. For example, the feedback may be partly observed, or the environment is not i.i.d. or can be operated by an adversary who aims to maximize the loss of the predictor. In this section, we will introduce some basic concepts about game theory and some fundamental theory of learning in games. We will focus on regret-based minimization procedures and limit our attention to finite strategic or normal form games. A more comprehensive study on this subject can be found in (Cesa-Bianchi and Lugosi, 2006; Nisan et al., 2007).

### 2.4.1 GAME PLAYING AND NASH EQUILIBRIUM

**$K$ -Player Normal-Form Games.** Consider a game with  $K$  players ( $1 < K < \infty$ ), where each player  $k \in \{1, \dots, K\}$  can take  $N_k$  possible actions. The players' actions can be represented by a vector  $\mathbf{i} = (i_1, \dots, i_K)$ , where  $i_k \in \{1, \dots, N_k\}$  denotes the action of player  $k$ . The loss suffered by the player  $k$  is denoted by  $\ell^{(k)}(\mathbf{i})$  since the loss is related to not only the action of player  $k$  but the action of all the other players. During each iteration of the game, each player tries to take actions in order to minimize its own loss.

Using a mixed strategy, player  $k$  takes actions based on a probability distribution  $\mathbf{p}^{(k)} = (p_1^{(k)}, \dots, p_{N_k}^{(k)})$  over the set of  $\{1, \dots, N_k\}$  actions. In particular, the actions played by all the  $K$  players can be denoted as a random vector  $\mathbf{I} = (I_1, \dots, I_K)$ , where  $I_k$  is the action played by player  $k$  which is a random variable taking value over the set of  $\{1, \dots, N_k\}$  actions distributed according to  $\mathbf{p}^{(k)}$ . The expected loss of player  $k$  can be computed as

$$\mathbb{E}\ell^{(k)}(\mathbf{I}) = \sum_{i_1=1}^{N_1} \cdots \sum_{i_K=1}^{N_K} p_{i_1}^{(1)} \times \cdots \times p_{i_K}^{(K)} \ell^{(k)}(i_1, \dots, i_K)$$

**Nash Equilibrium.** This is an important notion in game theory. In particular, a collective strategy of all players  $\mathbf{p}^{(1)} \times \cdots \times \mathbf{p}^{(K)}$  is called a *Nash equilibrium* if any mixed strategy among the  $K$  players  $\mathbf{p}^{(k)}$  is replaced by any new mixed strategy  $\mathbf{q}^{(k)}$  while all other  $K - 1$  players' mixed strategies make no change, we have

$$\mathbb{E}\ell^{(k)}(\mathbf{I}) \leq \mathbb{E}\ell^{(k)}(\mathbf{I}')$$

where  $\mathbf{I}'$  denotes the actions played by the  $K$  players using the new strategies. This definition means that in a Nash Equilibrium, no player can achieve a lower loss by only changing its own strategy if other players do not change. In a Nash Equilibrium, each player gets its own optimal strategy and has no incentive of changing its strategy. One can prove that every finite game has at least one Nash equilibrium, but a game may have multiple Nash equilibria depending on the structure of the game and the loss functions.

### 2.4.2 REPEATED TWO-PLAYER ZERO-SUM GAMES

A simple but important special class of  $K$ -Player Normal Form Games is the class of two-player zero-sum games where only one player plays against one opponent, i.e.,  $K = 2$ . *Zero-sum* means that for any action, the sum of losses of all players is zero. This indicates that the game is purely competitive and a player's loss results in another player's gain. In such games, the first player is often called the row player, and the second player is called the column player whose goal is to maximize the loss of the first player. To simplify notation, we consider the row player has  $N$  possible actions and the column player has  $M$  possible actions. We denote by  $L \in [0, 1]^{N \times M}$  where  $L(i, j)$  is the loss of the row player taking action  $i$  while the column player chooses action  $j$ , and the mixed strategies for the row and column players denoted by  $\mathbf{p} = (p_1, \dots, p_N)$  and  $\mathbf{q} = (q_1, \dots, q_M)$ , respectively. For the two mixed strategies  $\mathbf{p}$  and  $\mathbf{q}$ , the expected loss for the row player (which is equivalent to the expected gain of the column player) can be computed by

$$L(\mathbf{p}, \mathbf{q}) = \sum_{i=1}^N \sum_{j=1}^M p(i)q(j)L(i, j)$$

A pair of mixed strategies  $(\mathbf{p}, \mathbf{q})$  is a Nash equilibrium if and only if

$$L(\mathbf{p}, \mathbf{q}') \leq L(\mathbf{p}, \mathbf{q}) \leq L(\mathbf{p}', \mathbf{q}), \quad \forall \mathbf{p}', \forall \mathbf{q}'$$

One natural solution to the two-player zero-sum games is to follow the minimax solution. In particular, for the row player using some strategy  $\mathbf{p}$ , the worst-case loss is at most  $\max_{\mathbf{q}} L(\mathbf{p}, \mathbf{q})$  if the column player makes the decision after seeing  $\mathbf{p}$ . Therefore, the worst-case optimal strategy (also called the minimax optimal strategy) for the row player is  $\mathbf{p}^* = \arg \min_{\mathbf{p}} \max_{\mathbf{q}} L(\mathbf{p}, \mathbf{q})$ . Similarly, the maximin optimal strategy for the column player is  $\mathbf{q}^* = \arg \max_{\mathbf{q}} \min_{\mathbf{p}} L(\mathbf{p}, \mathbf{q})$ . The pair of  $(\mathbf{p}^*, \mathbf{q}^*)$  is called a minimax solution of the game. Surprisingly there is no difference between  $\min_{\mathbf{p}} \max_{\mathbf{q}} L(\mathbf{p}, \mathbf{q})$  and  $\max_{\mathbf{q}} \min_{\mathbf{p}} L(\mathbf{p}, \mathbf{q})$ , which is known as the von Neumann's minimax theorem, a fundamental result of game theory.

**Theorem 1** (*von Neumann's minimax theorem*) *In a two-player zero-sum game, when two players follow the strategies of the minimax solution, they reach the same optimal value*

$$V^* = \min_{\mathbf{p}} \max_{\mathbf{q}} L(\mathbf{p}, \mathbf{q}) = \max_{\mathbf{q}} \min_{\mathbf{p}} L(\mathbf{p}, \mathbf{q})$$

$V^*$  is called the value of the game which is unique for a two-player zero-sum game. A pair of mixed strategies  $(\mathbf{p}, \mathbf{q})$  is a Nash equilibrium if and only if it achieves the value of game.

We can now relate game theory to online learning as the problem of learning to play repeated two-player zero-sum games. In the context of online learning, the row player is also called the *learner* and the column player is called the *environment*. The repeated game playing between the row player and the column player is treated as a sequence of  $T$  rounds of interactions between the learner and the environment. On each round  $t = 1, \dots, T$ ,

- the learner chooses a mixed strategy  $\mathbf{p}_t$ ;
- the environment chooses a mixed strategy  $\mathbf{q}_t$  (may be chosen by the knowledge  $\mathbf{p}_t$ );
- the learner observes the losses  $L(i, \mathbf{q}_t) \quad \forall i \in [N]$

In general, the goal of the learner is to minimize the cumulative loss, i.e.,  $\sum_{t=1}^T L(\mathbf{p}_t, \mathbf{q}_t)$ .

### 3. Online Supervised Learning

#### 3.1 Overview

In this section, we survey a family of “online supervised learning” algorithms which define the fundamental approaches and principles for online learning methodologies toward supervised learning tasks Shalev-Shwartz (2011); Rakhlin et al. (2010).

We first discuss linear online learning methods, where a target model is a linear function. More formally, consider an input domain  $\mathcal{X}$  and an output domain  $\mathcal{Y}$  for a learning task, we aim to learn a hypothesis  $f : \mathcal{X} \mapsto \mathcal{Y}$ , where the target model  $f$  is linear. For example, consider a typical linear binary classification task, our goal is to learn a linear classifier  $f : \mathcal{X} \mapsto \{+1, -1\}$  as follows:  $f(\mathbf{x}_t; \mathbf{w}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x}_t)$ , where  $\mathcal{X}$  is typically a  $d$ -dimensional vector space  $\mathbb{R}^d$ ,  $\mathbf{w} \in \mathcal{X}$  is a weight vector specified for the classifier to be learned, and  $\text{sgn}(z)$  is an indicator function that outputs  $+1$  when  $z > 0$  and  $-1$  otherwise. We review two major types of linear online learning algorithms: first-order online learning and second-order online learning algorithms. Following this, we discuss Prediction with expert advice, and Online Learning with Regularization. This is followed by reviewing nonlinear online learning using kernel based methods. We discuss a variety of kernel-based online learning approaches, their computational challenges, and several approximation strategies for efficient learning. We end this section by discussing the theory for converting using online learning algorithms to learn a batch model that can generalize well.

#### 3.2 First-order Online Learning

In the following, we survey a family of first-order linear online learning algorithms, which exploit the first order information of the model during learning process.

##### 3.2.1 PERCEPTRON

Perceptron (Rosenblatt, 1958; Agmon, 1954; Novikoff, 1963) is the oldest algorithm for online learning. Algorithm 2 gives the Perceptron algorithm for online binary classification.

---

**Algorithm 2:** Perceptron
 

---

```

INIT:  $\mathbf{w}_1 = 0$ 
for  $t = 1, 2, \dots, T$  do
    Given an incoming instance  $\mathbf{x}_t$ , predict
     $\hat{y}_t = f_t(\mathbf{x}_t) = \text{sign}(\mathbf{w}_t \cdot \mathbf{x}_t)$ ;
    Receive the true class label  $y_t \in \{+1, -1\}$ ;
    if  $\hat{y}_t \neq y_t$  then
         $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_t \mathbf{x}_t$ ;
    end if
end for
    
```

---

In theory, by assuming the data is separable with some margin  $\gamma$ , the Perceptron algorithm makes at most  $\left(\frac{R}{\gamma}\right)^2$  mistakes, where the margin  $\gamma$  is defined as  $\gamma = \min_{t \in [T]} |\mathbf{x}_t \cdot \mathbf{w}^*|$

and  $R$  is a constant such that  $\forall t \in [T], \|\mathbf{x}_t\| \leq R$ . The larger the margin  $\gamma$  is, the tighter the mistake bound will be.

In literature, many variants of Perceptron algorithms have been proposed. One simple modification is the “normalized Perceptron” algorithm that differs only in the updating rule as follows:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + y_t \frac{\mathbf{x}_t}{\|\mathbf{x}_t\|}$$

The mistake bound of the “normalized Perceptron” algorithm can be improved from  $(\frac{R}{\gamma})^2$  to  $(\frac{1}{\gamma})^2$  for the separable case due to the normalization effect.

### 3.2.2 WINNOWER

Unlike the Perceptron algorithm that uses additive updates, Winnower (Littlestone, 1988) employs multiplicative updates. The problem setting is slightly different from the Perceptron:  $\mathcal{X} = \{0, 1\}^d$  and  $y \in \{0, 1\}$ . The goal is to learn a classifier  $f(x_1, \dots, x_n) = x_{i_1} \vee \dots \vee x_{i_k}$  called monotone disjunction, where  $i_k \in 1, \dots, d$ . The separating hyperplane for this classifier is given by  $x_{i_1} + \dots + x_{i_k}$ . The Winnower algorithm is outlined in Algorithm 3.

---

#### Algorithm 3: Winnower

---

**INIT:**  $\mathbf{w}_1 = \mathbf{1}^d$ , constant  $\alpha > 1$  (e.g.,  $\alpha = 2$ )  
**for**  $t = 1, 2, \dots, T$  **do**  
 Given an instance  $\mathbf{x}_t$ , predict  $\hat{y}_t = \mathbb{I}_{\mathbf{w}_t \cdot \mathbf{x}_t \geq \theta}$  (outputs 1 if statement holds and 0 otherwise);  
 Receive the true class label  $y_t \in \{1, 0\}$ ;  
**if**  $\hat{y}_t = 1, y_t = 0$  **then**  
 set  $w_i = 0$  for all  $x_{t,i} = 1$  (“elimination” or “demotion”),  
**end if**  
**if**  $\hat{y}_t = 0, y_t = 1$  **then**  
 set  $w_i = \alpha w_i$  for all  $x_{t,i} = 1$  (“promotion”).  
**end if**  
**end for**

---

The Winnower algorithm has a mistake bound  $\alpha k(\log_\alpha \theta + 1) + n/\theta$  where  $\alpha > 1$  and  $\theta \geq 1/\alpha$  and the target function is a  $k$ -literal monotone disjunction.

### 3.2.3 PASSIVE-AGGRESSIVE ONLINE LEARNING (PA)

This is a popular family of first-order online learning algorithms which generally follows the principle of margin-based learning (Crammer et al., 2006). Specifically, given an instance  $\mathbf{x}_t$  at round  $t$ , PA formulates the updating optimization as follows:

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 \quad \text{s.t.} \quad \ell_t(\mathbf{w}) = 0 \quad (3)$$

where  $\ell_t(\mathbf{w}) = \max(0, 1 - y_t \mathbf{w} \cdot \mathbf{x}_t)$  is the hinge loss. The above resulting update is passive whenever the hinge loss is zero, i.e.,  $\mathbf{w}_{t+1} = \mathbf{w}_t$  whenever  $\ell = 0$ . In contrast, whenever the

loss is nonzero, the approach will force  $\mathbf{w}_{t+1}$  aggressively to satisfy the constraint regardless of any step-size; the algorithm is thus named as “Passive-Aggressive” (PA) (Crammer et al., 2006). Specifically, PA aims to keep the updated classifier  $\mathbf{w}_{t+1}$  stay close to the previous classifier (“passiveness”) and ensure every incoming instance to be classified correctly by the updated classifier (“aggressiveness”). The regular PA algorithm assumes training data is always separable, which may not be true for noisy training data in real-world applications. To overcome such limitation, two variants of PA relax the assumption as:

$$\begin{aligned}
 \text{PA - I : } \mathbf{w}_{t+1} &= \arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 + C\xi \\
 \text{subject to } \ell_t(\mathbf{w}) &\leq \xi \text{ and } \xi \geq 0 \\
 \text{PA - II : } \mathbf{w}_{t+1} &= \arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 + C\xi^2 \\
 \text{subject to } \ell_t(\mathbf{w}) &\leq \xi
 \end{aligned} \tag{4}$$

where  $C$  is a positive parameter to balance the tradeoff between “passiveness” (first regularization term) and “aggressiveness” (second slack-variable term). By solving the three optimization tasks, we can derive the closed-form updating rules of three PA algorithms:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \tau_t y_t \mathbf{x}_t, \quad \tau_t = \begin{cases} \ell_t / \|\mathbf{x}_t\|^2 & \text{(PA)} \\ \min\{C, \ell_t / \|\mathbf{x}_t\|^2\} & \text{(PA-I)} \\ \frac{\ell_t}{\|\mathbf{x}_t\|^2 + \frac{1}{2C}} & \text{(PA-II)} \end{cases}$$

It is important to note a major difference between PA and Perceptron algorithms. Perceptron makes an update only when there is a classification mistake. However, PA algorithms aggressively make an update whenever the loss is nonzero (even if the classification is correct). In theory (Crammer et al., 2006), PA algorithms have comparable mistake bounds as the Perceptron algorithms, but empirically PA algorithms often outperform Perceptron significantly. The PA algorithms are outlined in Algorithm 4.

---

**Algorithm 4:** Passive Aggressive Algorithms

---

**INIT:**  $\mathbf{w}_1$ , Aggressiveness Parameter  $C$ ;  
**for**  $t = 1, 2, \dots, T$  **do**  
 Receive  $\mathbf{x}_t \in \mathbb{R}^d$ , predict  $\hat{y}_t$  using  $\mathbf{w}_t$ ;  
 Suffer loss  $\ell_t(\mathbf{w}_t)$ ;  
 Set  $\tau = \begin{cases} \ell_t / \|\mathbf{x}_t\|^2 & \text{(PA)} \\ \min\{C, \ell_t / \|\mathbf{x}_t\|^2\} & \text{(PA-I)} \\ \frac{\ell_t}{\|\mathbf{x}_t\|^2 + \frac{1}{2C}} & \text{(PA-II)} \end{cases}$   
 Update  $\mathbf{w}_{t+1} = \mathbf{w}_t + \tau_t y_t \mathbf{x}_t$ ;  
**end for**

---

### 3.2.4 ONLINE GRADIENT DESCENT (OGD)

Many online learning problems can be formulated as an online convex optimization task, which can be solved by applying the OGD algorithm. Consider the online binary classifica-

tion as an example, where we use the hinge loss function, i.e.,  $\ell_t(\mathbf{w}) = \max(0, 1 - y_t \mathbf{w} \cdot \mathbf{x}_t)$ . By applying the OGD algorithm, we can derive the updating rule as follows:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \eta_t y_t \mathbf{x}_t \quad (5)$$

where  $\eta_t$  is the learning rate (or step size) parameter. The OGD algorithm is outlined in Algorithm 5, where any generic convex loss function can be used.  $\Pi_{\mathcal{S}}$  is the projection function to constrain the updated model to lie in the feasible domain.

---

**Algorithm 5:** Online Gradient Descent

---

**INIT:**  $\mathbf{w}_1$ , convex set  $\mathcal{S}$ , step size  $\eta_t$ ;  
**for**  $t = 1, 2, \dots, T$  **do**  
     Receive  $\mathbf{x}_t \in \mathbb{R}^d$ , predict  $\hat{y}_t$  using  $\mathbf{w}_t$ ;  
     Suffer loss  $\ell_t(\mathbf{w}_t)$ ;  
     Update  $\mathbf{w}_{t+1} = \Pi_{\mathcal{S}}(\mathbf{w}_t - \eta_t \nabla \ell_t(\mathbf{w}_t))$   
**end for**

---

OGD and PA share similar updating rules but differ in that OGD often employs some predefined learning rate scheme while PA chooses the optimal learning rate  $\tau_t$  at each round (but subject to a predefined cost parameter  $C$ ). In literature, different OGD variants have been proposed to improve either theoretical bounds or practical issues, such as adaptive OGD (Hazan et al., 2007b), and mini-batch OGD (Dekel et al., 2012), amongst others.

### 3.2.5 OTHER FIRST-ORDER ALGORITHMS

In literature, there are also some other first-order online learning algorithms, such as Approximate Large Margin Algorithms (ALMA) (Gentile, 2001) which is a large margin variant of the p-norm Perceptron algorithm, and the Relaxed Online Maximum Margin Algorithm (ROMMA) (Li and Long, 2002). Many of these algorithms often follow the principle of large margin learning. The metaGrad algorithm (van Erven and Koolen, 2016) tries to adapt the learning rate automatically for faster convergence.

## 3.3 Second-Order Online Learning

Unlike the first-order online learning algorithms that only exploit the first order derivative information of the gradient for the online optimization tasks, second-order online learning algorithms exploit both first-order and second-order information in order to accelerate the optimization convergence. Despite the better learning performance, second-order online learning algorithms often fall short in higher computational complexity. In the following we present a family of popular second-order online learning algorithms.

### 3.3.1 SECOND ORDER PERCEPTRON (SOP)

SOP algorithm (Cesa-Bianchi et al., 2005a) is able to exploit certain geometrical properties of the data which are missed by the first-order algorithms.

For better understanding, we first introduce the whitened Perceptron algorithm, which strictly speaking, is not an online learning method. Assuming that the instances  $\mathbf{x}_1, \dots, \mathbf{x}_T$



are preliminarily available, we can get the correlation matrix  $M = \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top$ . The whitened Perceptron algorithm is simply the standard Perceptron run on the transformed sequence  $(M^{-1/2} \mathbf{x}_1, y_1), \dots, (M^{-1/2} \mathbf{x}_T, y_T)$ . By reducing the correlation matrix of the transformed instances, the whitened Perceptron algorithm can achieve significantly better mistake bound.

SOP can be viewed as an online variant of the whitened Perceptron algorithm. In on-line setting the correlation matrix  $M$  can be approximated by the previously seen instances. SOP is outlined in Algorithm 6

---

**Algorithm 6:** SOP
 

---

**INIT:**  $\mathbf{w}_1 = 0$ ,  $X_0 = \emptyset$ ,  $\mathbf{v}_0 = 0$ ,  $k = 1$   
**for**  $t = 1, 2, \dots, T$  **do**  
     Given an incoming instance  $\mathbf{x}_t$ , set  $S_t = [X_{k-1} \ \mathbf{x}_t]$ ,  
     predict  $\hat{y}_t = f_t(\mathbf{x}_t) = \text{sign}(\mathbf{w}_t \cdot \mathbf{x}_t)$ , where  $\mathbf{w}_t = (aI_n + S_t S_t^\top)^{-1} \mathbf{v}_{k-1}$   
     Receive the true class label  $y_t \in \{+1, -1\}$ ;  
     **if**  $\hat{y}_t \neq y_t$  **then**  
          $\mathbf{v}_k = \mathbf{v}_{k-1} + y_t \mathbf{x}_t$ ,  $X_k = S_t$ ,  $k = k + 1$ .  
     **end if**  
**end for**

---

Here  $a \in \mathbb{R}^+$  is a parameter that guarantees the existence of the matrix inverse.

### 3.3.2 CONFIDENCE WEIGHTED LEARNING (CW)

The CW algorithm (Dredze et al., 2008) is motivated by the following observation: the frequency of occurrence of different features may differ a lot in an online learning task. (For example) The parameters of binary features are only updated when the features occur. Thus, the frequent features typically receive more updates and are estimated more accurately compared to rare features. However, no distinction is made between these feature types in most online algorithms. This indicates that the lack of second order information about the frequency or confidence of the features can hurt the learning.

In the CW setting, we model the linear classifier with a Gaussian distribution, i.e.,  $\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ , where  $\boldsymbol{\mu} \in \mathbb{R}^d$  is the mean vector and  $\Sigma \in \mathbb{R}^{d \times d}$  is the covariance matrix. When making a prediction, the prediction confidence  $M = \mathbf{w} \cdot \mathbf{x}$  also follows a Gaussian distribution:  $M \sim \mathcal{N}(\boldsymbol{\mu}_M, \Sigma_M)$ , where  $\boldsymbol{\mu}_M = \boldsymbol{\mu} \cdot \mathbf{x}$  and  $\Sigma_M = \mathbf{x}^\top \Sigma \mathbf{x}$ .

Similar to the PA update strategy, the update rule in round  $t$  can be obtained by solving the following convex optimization problem:

$$(\boldsymbol{\mu}_{t+1}, \Sigma_{t+1}) = \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^d} D_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu}, \Sigma) \| \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t)) \quad \text{s.t.} \quad \Pr[y_t M_t \geq 0] \geq \eta \quad (6)$$

The objective function means that the new distribution should stay close to the previous distribution so that the classifier does not forget the information learnt from previous instances, where the distance between the two distributions is measured by the KL divergence. The constraint means that the new classifier should classify the new instance  $\mathbf{x}_t$  correctly with probability higher than a predefined threshold parameter  $\eta \in (0, 1)$ .

Note that this is only the basic form of confidence weighted algorithms and has several drawbacks. 1) Similar to the hard margin PA algorithm, the constraint forces the new instance to be correctly classified, which makes this algorithm very sensitive to noise. 2) The constraint is in a probability form. It is easy to solve a problem with the constraint  $g(\boldsymbol{\mu}_M, \Sigma_M) < 0$ . However, a problem with a probability form constraint is only solvable when the distribution is known. Thus, this method faces difficulty in generalizing to other online learning tasks where the constraint does not follow a Gaussian distribution.

### 3.3.3 ADAPTIVE REGULARIZATION OF WEIGHT VECTORS (AROW)

AROW (Crammer et al., 2009b) is a variant of CW that is designed for non-separable data. This algorithm adopts the same Gaussian distribution assumption on classifier vector  $\mathbf{w}$  while the optimization problem is different. By recasting the CW constraint as regularizers, the optimization problem can be formulated as:

$$\mathcal{C}(\boldsymbol{\mu}, \Sigma) = \text{D}_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu}, \Sigma) || \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t)) + \lambda_1 \ell(y_t, \boldsymbol{\mu} \cdot \mathbf{x}_t) + \lambda_2 \mathbf{x}_t^\top \Sigma \mathbf{x}_t \quad (7)$$

where  $\ell(y_t, \boldsymbol{\mu} \cdot \mathbf{x}_t) = (\max(0, 1 - y_t \boldsymbol{\mu} \cdot \mathbf{x}_t))^2$  is the squared-hinge loss. During each iteration, the update rule is obtained by solving the optimization problem:

$$(\boldsymbol{\mu}_{t+1}, \Sigma_{t+1}) = \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^d} (\mathcal{C}(\boldsymbol{\mu}, \Sigma))$$

which balances the three desires. First, the parameters should not change radically on each round, since the current parameters contain information about previous examples (first term). Second, the new mean parameters should predict the current example with low loss (second term). Finally, as we see more examples, our confidence in the parameters should generally grow (third term).  $\lambda_1$  and  $\lambda_2$  are two positive parameters that control the weight of the three desires.

Besides the robustness to noisy data, another important advantage of AROW is its ability to be easily generalized to other online learning tasks, such as Confidence Weighted Online Collaborative Filtering algorithm (Lu et al., 2013) and Second-Order Online Feature Selection (Wu et al., 2014).

### 3.3.4 SOFT CONFIDENCE WEIGHTED LEARNING (SCW)

This is a variant of CW learning in order to deal with non-separable data (Wang et al., 2012b, 2016a). Different from AROW which directly adds loss and confidence regularization, and thus loses the adaptive margin property, SCW exploits adaptive margin by assigning different margins for different instances via a probability formulation. Consequently, SCW tends to be more efficient and effective.

Specifically, the constraint of CW can be rewritten as  $y_t(\boldsymbol{\mu} \cdot \mathbf{x}_t) \geq \phi \sqrt{\mathbf{x}_t^\top \Sigma \mathbf{x}_t}$ . Thus, the loss function can be defined as:  $\ell(\mathcal{N}(\boldsymbol{\mu}, \Sigma); (\mathbf{x}_t, y_t)) = \max(0, \phi \sqrt{\mathbf{x}_t^\top \Sigma \mathbf{x}_t} - y_t(\boldsymbol{\mu} \cdot \mathbf{x}_t))$ . The original CW optimization can be rewritten as:

$$\begin{aligned} (\boldsymbol{\mu}_{t+1}, \Sigma_{t+1}) = \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^d} & \text{D}_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu}, \Sigma) || \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t)) \\ & \text{subject to } \ell(\mathcal{N}(\boldsymbol{\mu}, \Sigma); (\mathbf{x}_t, y_t)) = 0 \end{aligned}$$

Inspired by soft-margin PA variants, SCW generalized CW into two soft-margin formulations:

$$\begin{aligned}(\boldsymbol{\mu}_{t+1}, \Sigma_{t+1}) &= \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^d} D_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu}, \Sigma) || \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t)) + C\ell(\mathcal{N}(\boldsymbol{\mu}, \Sigma); (\mathbf{x}_t, y_t)) \\(\boldsymbol{\mu}_{t+1}, \Sigma_{t+1}) &= \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^d} D_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu}, \Sigma) || \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t)) + C\ell^2(\mathcal{N}(\boldsymbol{\mu}, \Sigma); (\mathbf{x}_t, y_t))\end{aligned}$$

where  $C \in \mathbb{R}^+$  is a parameter controls the aggressiveness of this algorithm, similar to the  $C$  in PA algorithm. The two algorithms are termed “SCW-I” and “SCW-II”.

### 3.3.5 OTHER SECOND-ORDER ALGORITHMS

The confidence weighted idea also works for other online learning tasks such as multi-class classification (Crammer et al., 2009a), active learning (Dredze and Crammer, 2008) and structured-prediction (Mejer and Crammer, 2010). There are many other online learning algorithms that adopt second order information: IELLIP (Yang et al., 2009b) assumes the objective classifier  $\mathbf{w}$  lies in an ellipsoid and incrementally updates the ellipsoid based on the current received instance. Other approaches include New variant of Adaptive Regularization (NAROW) (Orabona and Crammer, 2010) and the Normal Herding method via Gaussian Herding (NHERD) (Crammer and Lee, 2010). Recently, Sketched Online Newton (Luo et al., 2016) made significant improvements to speed-up second order online learning.

## 3.4 Prediction with Expert Advice

This is an important online learning subject (Roughgarden and Schrijvers, 2017) with many applications. A general setting is as follows. A learner has  $N$  experts to choose from, denoted by integers  $1, \dots, N$ . At each time step  $t$ , the learner decides on a distribution  $\mathbf{p}_t$  over the experts, where  $p_{t,i} \geq 0$  is the weight of each expert  $i$ , and  $\sum_{i=1}^N p_{t,i} = 1$ . Each expert  $i$  then suffers some loss  $\ell_{t,i}$  according to the environment. The overall loss suffered by the learner is  $\sum_{i=1}^N p_{t,i} \ell_{t,i} = \mathbf{p}_t^\top \boldsymbol{\ell}_t$ , i.e., the weighted average loss of the experts with respect to the distribution chosen by the learner.

Typically we assume that the loss suffered by any expert is bounded. Specifically, we assume  $\ell_{t,i} \in [0, 1]$  without loss of generality. Besides this condition, no assumption is made on the form of the loss, or about how they are generated. Suppose the cumulative losses experienced by each expert and the forecaster are calculated respectively as follows:

$$L_{t,i} = \sum_{s=1}^t \ell_{s,i}, \quad L_t = \sum_{s=1}^t \mathbf{p}_s^\top \boldsymbol{\ell}_s.$$

The loss difference between the forecaster and the expert is known as the “regret”, i.e.,

$$R_{t,i} = L_t - L_{t,i}, \quad i = 1, \dots, N.$$

The goal of learning the forecaster is to make the regret with respect to each expert as small as possible, which is equivalent to minimizing the overall regret, i.e.,

$$R_T = \max_{1 \leq i \leq N} R_{T,i} = L_T - \min_{1 \leq i \leq N} L_{T,i}$$

In general, online prediction with expert advice aims to find an ideal forecaster to achieve a vanishing per-round regret, a property known as the *Hannan-consistency* (Hannan, 1957), i.e.,

$$R_T = o(T) \Leftrightarrow \lim_{T \rightarrow \infty} \frac{1}{T} \left( L_T - \min_{1 \leq i \leq N} L_{T,i} \right)$$

An online learner satisfying the above is called a Hannan-consistent forecaster (Cesa-Bianchi and Lugosi, 2006). Next we review some representative algorithms for prediction with expert advice.

#### 3.4.1 WEIGHTED MAJORITY ALGORITHMS

The weighted majority algorithm (WM) is a simple but widely studied algorithm that makes a binary prediction based on a series of expert advices (Littlestone and Warmuth, 1994, 1989). The simplest version is shown in Algorithm 7, where  $\beta \in (0, 1)$  is a user specified discount rate parameter.

---

##### **Algorithm 7:** Weighted Majority

---

**INIT:** Initialize the weights  $p_1, p_2, \dots, p_N$  of all experts to  $1/N$ .  
**for**  $t = 1, 2, \dots, T$  **do**  
     Get the prediction  $x_1, \dots, x_N$  from  $N$  experts.  
     Output 1 if  $\sum_{i:x_i=1} p_i \geq \sum_{i:x_i=0} p_i$ ; otherwise output 0.  
     Receive the true value; if the  $i$ -th expert made a mistake, then  $p_i = p_i * \beta$   
**end for**

---

#### 3.4.2 RANDOMIZED MULTIPLICATIVE WEIGHTS ALGORITHMS

This algorithm works under the same assumption that the expert advices are all binary (Arora et al., 2012b). While the prediction is random, the algorithm gives the prediction 1 with probability of  $\gamma = \frac{\sum_{i:x_i=1} p_i}{\sum_i p_i}$  and 0 with probability of  $1 - \gamma$ .

#### 3.4.3 HEDGE ALGORITHM

The Hedge algorithm (Freund and Schapire, 1997) is perhaps the most well-known approach for online prediction with expert advice, which can be viewed as a direct generalization of Littlestone and Warmuth's weighted majority algorithm (Littlestone and Warmuth, 1994, 1989). The working of Hedge algorithm is shown in Algorithm 8. The algorithm maintains a weight vector whose value at time  $t$  is denoted  $\mathbf{w}_t = (w_{t,1}, \dots, w_{t,N})$ . At all times, all weights are nonnegative. All of the weights of the initial weight vector  $\mathbf{w}_1$  must be nonnegative and sum to one, which can be considered as a prior over the set of experts. If it is believed that one expert performs the best, it is better to assign it more weight. If no prior is known, it is better to set all the initial weights equally, i.e.,  $w_{1,i} = 1/N$  for all  $i$ . The algorithm uses the normalized distribution to make prediction, i.e.,  $\mathbf{p}_t = \mathbf{w}_t / \sum_{i=1}^N w_{t,i}$ . After the loss  $\ell_t$  is disclosed, the weight vector  $\mathbf{w}_t$  is updated using a multiplicative rule  $w_{t+1,i} = w_{t,i} \beta^{\ell_{t,i}}$ ,  $\beta \in [0, 1]$ , which implies that the weight of expert  $i$  will exponentially decrease with the loss  $\ell_{t,i}$ . In theory, the Hedge algorithm is proved to be Hannan consistent.

---

**Algorithm 8:** Hedge Algorithm
 

---

**INIT:**  $\beta \in [0, 1]$ , initial weight vector  $\mathbf{w}_1 \in [0, 1]^N$  with  $\sum_{i=1}^N w_{1,i} = 1$   
**for**  $t = 1, 2, \dots, T$  **do**  
     set distribution  $\mathbf{p}_t = \frac{\mathbf{w}_t}{\sum_{i=1}^N w_{t,i}}$ ;  
     Receive loss  $\ell_t \in [0, 1]^N$  from environment;  
     Suffer loss  $\mathbf{p}_t^\top \ell_t$ ;  
     Update the new weight vector to  $w_{t+1,i} = w_{t,i} \beta^{\ell_{t,i}}$   
**end for**

---

### 3.4.4 EWAF ALGORITHMS

Besides Hedge, there are some other algorithms for online prediction with expert advice under more challenging settings, including exponentially weighted average forecaster (EWAF) and Greedy Forecaster (GF) (Cesa-Bianchi and Lugosi, 2006). We will mainly discuss EWAF, which is shown in Algorithm 9

---

**Algorithm 9:** EWAF
 

---

**INIT:** a poll of experts  $f_i$ ,  $i = 1, \dots, N$  and  $L_{0,i} = 0$ ,  $i = 1, \dots, N$ , and learning rate  $\eta$   
**for**  $t = 1, 2, \dots, T$  **do**  
     The environment chooses the next outcome  $y_t$  and the expert advice  $\{f_{t,i}\}$ ;  
     The expert advice is revealed to the forecaster  
     The forecaster chooses the prediction  $\hat{p}_t = \frac{\sum_{i=1}^N \exp(-\eta L_{t-1,i}) f_{t,i}}{\sum_{i=1}^N \exp(-\eta L_{t-1,i})}$   
     The environment reveals the outcome  $y_t$ ;  
     The forecaster incurs loss  $\ell(\hat{p}_t, y_t)$  and;  
     Each expert incurs loss  $\ell(f_{t,i}, y_t)$   
     The forecaster update the cumulative loss  $L_{t,i} = L_{t-1,i} + \ell(f_{t,i}, y_t)$   
**end for**

---

The difference between EWAF and Hedge is that the loss in Hedge is the inner product between the distribution and the loss suffered by each expert, while for EWAF, the loss is between the prediction and the true label, which can be much more complex.

### 3.4.5 PARAMETER-FREE ONLINE LEARNING

A category of prediction with expert advice deals with learning without user specified learning rate. It is a difficult task to set a learning rate prior to the learning procedure. To address this issue, parameter-free online algorithms were proposed. Among the early efforts, Chaudhuri et al. (2009) proposed a variant of Hedge Algorithm without the use of a learning rate. The proposed method achieved optimal regret matching the best bounds of all the previous algorithms (with optimally-tuned parameters). Chernov and Vovk (2010) improved this bound, but did not have a closed-form solution. There were further extensions which derived data-dependent bounds too (Luo and Schapire, 2015; Koolen and Van Erven, 2015).

### 3.5 Online Learning with Regularization

Traditional online learning methods learn a classifier  $\mathbf{w} \in \mathbb{R}^d$  where the magnitude of each element  $|\mathbf{w}^j|$  weights the importance of each feature, which are often non-zero. When dealing with high dimensional data, traditional online learning methods suffer from expensive computational time and space costs. This drawback is often addressed using regularization by performing Sparse online learning, which aims to exploit the sparsity property with real-world high-dimensional data. Specifically, a batch sparse learning problem can be formalized as:

$$P(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell_t(\mathbf{w}) + \phi_s(\mathbf{w})$$

where  $\phi_s$  is a sparsity-inducing regularizer. For example, when choosing  $\phi_s = \lambda \|\mathbf{w}\|_0$ , it is equivalent to imposing a hard constraint on the number of nonzero elements in  $\mathbf{w}$ . Instead of choosing  $\ell_0$ -norm which is hard to be optimized, a more commonly used regularizer is  $\ell_1$ -norm, i.e.,  $\phi_s = \lambda \|\mathbf{w}\|_1$ , which can induce sparsity of the weight vector but does not explicitly constrain the number of nonzero elements. The following reviews some popular sparse online learning methods.

#### 3.5.1 TRUNCATED GRADIENT DESCENT

A straightforward idea to sparse online learning is to modify Online Gradient Descent and round small coefficients of the weight vector to 0 after every  $K$  iterations:

$$\mathbf{w}_{t+1} = T_0(\mathbf{w}_t - \eta \nabla \ell_t(\mathbf{w}_t), \theta)$$

where the function  $T_0(\mathbf{v}, \theta)$  performs an element-wise rounding on the input vector: if the  $j$ -th element  $v^j$  is smaller than the threshold  $\theta$ , set  $v^j = 0$ . Despite its simplicity, this method struggles to provide satisfactory performance because the aggressive rounding strategy may ignore many useful weights which may be very small due to low frequency of appearance.

Motivated by addressing the above limitation, the Truncated Gradient Descent (TGD) method (Langford et al., 2009) explores a less aggressive version of the truncation function:

$$\begin{aligned} \mathbf{w}_{t+1} &= T_1(\mathbf{w}_t - \eta \nabla \ell_t(\mathbf{w}_t), \eta g_i, \theta) \\ \text{where } T_1(v^j, \alpha, \theta) &= \begin{cases} \max(0, v^j - \alpha) & \text{if } v^j \in [0, \theta] \\ \min(0, v^j + \alpha) & \text{if } v^j \in [-\theta, 0] \\ v^j & \text{otherwise} \end{cases} \end{aligned}$$

where  $g_i > 0$  is a parameter that controls the level of aggressiveness of the truncation. By exploiting sparsity, TGD achieves efficient time and space complexity that is linear with respect to the number of nonzero features and independent of the dimensionality  $d$ . In addition, it is proven to enjoy a regret bound of  $O(\sqrt{T})$  for convex loss functions when setting  $\eta = O(1/\sqrt{T})$ .

#### 3.5.2 FORWARD LOOKING SUBGRADIENTS (FOBOS)

Consider the objective function in the  $t$ -th iteration of a sparse online learning task as  $\ell_t(\mathbf{w}) + r(\mathbf{w})$ , FOBOS (Duchi and Singer, 2009) assumes  $f_t$  is a convex loss function (dif-

ferentiable), and  $r$  is a sparsity-inducing regularizer (non-differentiable). FOBOS updates the classifier in the following two steps:

- (1) Perform Online Gradient Descent:  $\mathbf{w}_{t+\frac{1}{2}} = \mathbf{w}_t - \eta_t \nabla \ell_t(\mathbf{w}_t)$
- (2) Project the solution in (i) such that the projection stays close to the interim vector  $\mathbf{w}_{t+\frac{1}{2}}$  and (ii) has a low complexity due to  $r$ :

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w}} \left\{ \frac{1}{2} \|\mathbf{w} - \mathbf{w}_{t+\frac{1}{2}}\|^2 + \eta_{t+\frac{1}{2}} r(\mathbf{w}) \right\}$$

When choosing  $\ell_1$ -norm as the regularizer, the above optimization can be solved with the closed-form solution for each coordinate:

$$w_{t+1}^j = \text{sgn}(w_{t+\frac{1}{2}}^j) \left[ |w_{t+\frac{1}{2}}^j| - \eta_{t+\frac{1}{2}} \right]_+$$

The FOBOS algorithm with  $\ell_1$ -norm regularizer can be viewed as a special case of TGD, where the truncation threshold  $\theta = \infty$ , and the truncation frequency  $K = 1$ . When  $\eta_{t+\frac{1}{2}} = \eta_{t+1}$  and  $\eta_t = O(1/\sqrt{t})$ , this algorithm also achieves  $O(\sqrt{T})$  regret bound.

### 3.5.3 REGULARIZED DUAL AVERAGING (RDA)

Motivated by the theory of dual-averaging techniques (Nesterov, 2009), the RDA algorithm (Xiao, 2009) updates the classifier by:

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w}} \left\{ \bar{\mathbf{g}}_t^\top \mathbf{w} + \Psi(\mathbf{w}) + \frac{\beta_t}{t} h(\mathbf{w}) \right\}$$

where  $\Psi(\mathbf{w})$  is the original sparsity-inducing regularizer, i.e.,  $\Psi(\mathbf{w}) = \lambda \|\mathbf{w}\|_1$ ;  $h(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$  is an auxiliary strongly convex function and  $\bar{\mathbf{g}}_t$  is the averaged gradients of all previous iterations, i.e.,  $\bar{\mathbf{g}} = \frac{1}{t} \sum_{\tau=1}^t \nabla \ell_\tau(\mathbf{w}_\tau)$ . Setting the step size  $\beta_t = \gamma \sqrt{t}$ , one can derive the closed-form solution:

$$w_{t+1}^j = \begin{cases} 0 & \text{if } |\bar{g}_t^j| < \lambda \\ -\frac{\sqrt{t}}{\gamma} (\bar{g}_t^j - \lambda \text{sgn}(\bar{g}_t^j)) & \text{otherwise} \end{cases}$$

To further pinpoint the differences between RDA and FOBOS, we rewrite FOBOS in the same notation as RDA:

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w}} \left\{ \mathbf{g}_t^\top \mathbf{w} + \Psi(\mathbf{w}) + \frac{1}{2\alpha_t} \|\mathbf{w} - \mathbf{w}_t\|_2^2 \right\}$$

Specifically, RDA differs from FOBOS in several aspects. First, RDA uses the averaged gradient instead of the current gradient. Second,  $h(\mathbf{w})$  is a global proximal function instead of its local Bregman divergence. Third, the coefficient for  $h(\mathbf{w})$  is  $\beta_t/t = \gamma/\sqrt{t}$  which is  $1/\alpha_t = O(\sqrt{t})$  in FOBOS. Fourth, the truncation of RDA is a constant  $\lambda$ , while the truncation in FOBOS  $\eta_{t+\frac{1}{2}}$  decrease with a factor  $\sqrt{t}$ . Clearly, RDA uses a more aggressive truncation threshold, thus usually generates significantly more sparse solutions. RDA also ensures the  $O(\sqrt{T})$  regret bound.

### 3.5.4 ADAPTIVE REGULARIZATION

One major issue with both FOBOS and RDA is that the auxiliary strongly convex function  $h(\mathbf{w})$  may not fully exploit the geometry information of underlying data distribution. Instead of choosing  $h(\mathbf{w})$  as an  $\ell_2$ -norm  $\frac{1}{2}\|\mathbf{w}\|^2$  in RDA or a Mahalanobis norm  $\|\cdot\|_{H_t}$  in FOBOS, (Duchi et al., 2011a) proposed a data-driven adaptive regularization for  $h(\mathbf{w})$ , i.e.,

$$h_t(\mathbf{w}) = \frac{1}{2} \mathbf{w}^\top H_t \mathbf{w}$$

where  $H_t = (\sum_{\tau=1}^t \mathbf{g}_\tau \mathbf{g}_\tau^\top)^{\frac{1}{2}}$  accumulates the second order info from the previous instances over time. Replacing the previous  $h(\mathbf{w})$  in both RDA and FOBOS by the temporal adaptation function  $h_t(\mathbf{w})$ , (Duchi et al., 2011a) derived two generalized algorithms (Ada-RDA and Ada-FOBOS) with the solutions as follows respectively.

Ada-RDA:

$$w_{t+1}^j = \begin{cases} 0 & \text{if } |\bar{g}_t^j| < \lambda \\ -\frac{t}{\beta H_{t,jj}} (\bar{g}_t^j - \lambda \text{sgn}(\bar{g}_t^j)) & \text{otherwise} \end{cases}$$

Ada-FOBOS:

$$w_{t+1}^j = \text{sgn}(w_t^i - \frac{\alpha_t}{H_{t,jj}} g_t^j) \left[ |w_t^i - \frac{\alpha_t}{H_{t,jj}} g_t^j| - \frac{\alpha_t \lambda}{H_{t,ii}} \right]_+ \quad (8)$$

In the above,  $H_t$  is approximated by a diagonal matrix since computing the root of a matrix is computationally impractical in high-dimensional data.

### 3.5.5 ONLINE FEATURE SELECTION

Online feature selection (Hoi et al., 2012; Wang et al., 2014c; Kale et al., 2017; Wu et al., 2017b) is closely related to sparse online learning in that they both aim to learn an efficient classifier for very high dimensional data. However, the sparse learning algorithms aim to minimize the  $\ell_1$  regularized loss, while the feature selection algorithms are motivated to explicitly address the feature selection issue and thus impose a hard constraint on the number of non-zero elements in classifier. Because of these similarities, they share some common strategies such as truncation and projection.

### 3.5.6 OTHERS

Two stochastic methods were proposed in (Shalev-Shwartz and Tewari, 2011) for  $\ell_1$ -regularized loss minimization. The Stochastic Coordinate Descent (SCD) algorithm randomly selects one coordinate from  $d$  dimensions and update this single coordinate with the gradient of the total loss of all instances. The Stochastic Mirror Descent Made Sparse (SMIDAS) algorithm combines the idea of truncating the gradient with mirror descent algorithm, i.e., truncation is performed on the vector in dual space. The disadvantage of the two algorithms is that their computational complexity depends on the dimensionality  $d$ . Besides, the two algorithms are designed in batch learning setting, i.e., they assume all instances are known prior to the learning task. Besides, there are also some recent sparse online learning algorithms proposed (Wang et al., 2014a, 2015a), which combine the ideas of sparse learning, second order online learning, and cost-sensitive classification together to make the online algorithms scalable for high-dimensional class-imbalanced learning tasks.



### 3.6 Online Learning with Kernels

We now survey a family of “Kernel-based Online Learning” algorithms for learning a nonlinear target function, where the nonlinearity is induced by kernels. We take the typical nonlinear binary classification task as an example. Our goal is to learn a nonlinear classifier  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  from a sequence of labeled instances  $(\mathbf{x}_t, y_t), t = 1, \dots, T$ , where  $\mathbf{x}_t \in \mathbb{R}^d$  and  $y_t \in \{+1, -1\}$ . We build the classification rule as:  $\hat{y}_t = \text{sgn}(f(\mathbf{x}_t))$ , where  $\hat{y}_t$  is the predicted class label. We measure the classification confidence of certain instance  $\mathbf{x}_t$  by  $|f(\mathbf{x}_t)|$ . Similar to the linear case, for an online classification task, one can define the hinge loss function  $\ell(\cdot)$  for the  $t$ -th instance using the classifier at the  $t$ -th iteration:

$$\ell((\mathbf{x}_t, y_t); f_t) = \max(0, 1 - y_t f_t(\mathbf{x}_t))$$

Formally speaking, an online nonlinear learner aims to achieve the lowest regret  $R(T)$  after time  $T$ , where the regret function  $R(T)$  is defined as follows:

$$R(T) = \sum_{t=1}^T \ell_t(f_t) - \inf_f \sum_{t=1}^T \ell_t(f), \quad (9)$$

where  $\ell_t(\cdot)$  is the loss for the classification of instance  $(\mathbf{x}_t, y_t)$ , which is short for  $\ell((\mathbf{x}_t, y_t); \cdot)$ . We denote by  $f^*$  the optimal solution of the second term, i.e.,  $f^* = \arg \min_f \sum_{t=1}^T \ell_t(f)$ .

In the following, we first introduce online kernel methods and then survey a family of scalable online kernel learning algorithms organized into two major categories: (i) budget online kernel learning using *budget maintenance* strategies and (ii) budget online kernel learning using *functional approximation* strategies. Then we briefly introduce some approaches for online learning with multiple kernels. Without loss of generality, we will adopt the above online binary classification setting for the discussions in this section.

#### 3.6.1 ONLINE KERNEL METHODS

We refer to the output  $f$  of the learning algorithm as a *hypothesis* and denote the set of all possible hypotheses by  $\mathcal{H} = \{f | f : \mathbb{R}^d \rightarrow \mathbb{R}\}$ . Here  $\mathcal{H}$  a Reproducing Kernel Hilbert Space (**RKHS**) endowed with a kernel function  $\kappa(\cdot, \cdot) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  (Vapnik and Vapnik, 1998) implementing the inner product  $\langle \cdot, \cdot \rangle$  such that: 1)  $\kappa$  has the reproducing property  $\langle f, \kappa(\mathbf{x}, \cdot) \rangle = f(\mathbf{x})$  for  $\mathbf{x} \in \mathbb{R}^d$ ; 2)  $\mathcal{H}$  is the closure of the span of all  $\kappa(\mathbf{x}, \cdot)$  with  $\mathbf{x} \in \mathbb{R}^d$ , that is,  $\kappa(\mathbf{x}, \cdot) \in \mathcal{H} \forall \mathbf{x} \in \mathcal{X}$ . The inner product  $\langle \cdot, \cdot \rangle$  induces a norm on  $f \in \mathcal{H}$  in the usual way:  $\|f\|_{\mathcal{H}} := \langle f, f \rangle^{\frac{1}{2}}$ . We denote by  $\mathcal{H}_{\kappa}$  an RKHS with explicit dependence on kernel  $\kappa$ . Throughout the analysis, we assume  $\kappa(\mathbf{x}, \mathbf{x}) \leq X^2, \forall \mathbf{x} \in \mathbb{R}^d, X \in \mathbb{R}^+$  is a constant.

The task of training the model of SVM  $f(\mathbf{x})$  in batch is formulated as the optimization:

$$\min_{f \in \mathcal{H}_{\kappa}} \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2 + \frac{1}{T} \sum_{t=1}^T \ell(f(\mathbf{x}_t); y_t) \quad (10)$$

where  $\lambda > 0$  is a regularization parameter used to control model complexity. According to the Representer Theorem (Schölkopf et al., 2001), the optimal solution of the above convex optimization problem lies in the span of  $T$  kernels, i.e., those centered on the training points. Consequently, the goal of a typical online kernel learning algorithm is to learn

the kernel-based predictive model  $f(\mathbf{x})$  for classifying a new instance  $\mathbf{x} \in \mathbb{R}^d$  as follows:  $f(\mathbf{x}) = \sum_{t=1}^T \alpha_t \kappa(\mathbf{x}_t, \mathbf{x})$ , where  $T$  is the number of processed instances,  $\alpha_t$  denotes the coefficient of the  $t$ -th instance, and  $\kappa(\cdot, \cdot)$  denotes the kernel function. We define support vector (SV) as the instance whose coefficient  $\alpha$  is nonzero. Thus, we rewrite the previous classifier as  $f(\mathbf{x}) = \sum_{i \in \mathcal{SV}} \alpha_i \kappa(\mathbf{x}_i, \mathbf{x})$ , where  $\mathcal{SV}$  is the set of SV's and  $i$  is its index. We use the notation  $|\mathcal{SV}|$  to denote the SV set size. In literature, different online kernel methods have been proposed. We begin by introducing the simplest one, that is, the kernelized Perceptron algorithm.

**Kernelized Perceptron.** This extends the Perceptron algorithm using the kernel trick. Algorithm 10 outlines the Kernelized Perceptron algorithm (Freund and Schapire, 1999a).

---

**Algorithm 10:** Kernelized Perceptron

---

```

INIT:  $f_0 = 0$ 
for  $t = 1, 2, \dots, T$  do
    Given an incoming instance  $\mathbf{x}_t$ , predict  $\hat{y}_t = \text{sgn}(f_t(\mathbf{x}_t))$ ;
    Receive the true class label  $y_t \in \{+1, -1\}$ ;
    if  $\hat{y}_t \neq y_t$  then
         $\mathcal{SV}_{t+1} = \mathcal{SV}_t \cup (\mathbf{x}_t, y_t)$ ,  $f_{t+1} = f_t + y_t \kappa(\mathbf{x}_t, \cdot)$ ;
    end if
end for
    
```

---

The algorithm works similar to the standard Perceptron algorithm, except that the inner product, i.e.,  $f_t(\mathbf{x}_t) = \sum_i \alpha_i \mathbf{x}_i^\top \mathbf{x}_t$ , is replaced by a kernel function in the kernel Perceptron, i.e.,  $f_t(\mathbf{x}_t) = \sum_i \alpha_i \kappa(\mathbf{x}_i, \mathbf{x}_t)$ .

**Kernelized OGD.** This extends the OGD algorithm with kernels (Kivinen et al., 2004), as shown in Algorithm 11. Here,  $\eta_t > 0$  is the learning rate parameter, and  $\ell'_t$  is used to denote the derivative of loss function with respect to the classification score  $f_t(\mathbf{x}_t)$ .

---

**Algorithm 11:** Kernelized OGD

---

```

INIT:  $f_0 = 0$ 
for  $t = 1, 2, \dots, T$  do
    Given an incoming instance  $\mathbf{x}_t$ , predict  $\hat{y}_t = \text{sgn}(f_t(\mathbf{x}_t))$ ;
    Receive the true class label  $y_t \in \{+1, -1\}$ ;
    if  $\ell_t(f_t) > 0$  then
         $\mathcal{SV}_{t+1} = \mathcal{SV}_t \cup (\mathbf{x}_t, y_t)$ ,  $f_{t+1} = f_t - \eta_t \nabla \ell_t(f_t(\mathbf{x}_t)) = f_t - \eta_t \ell'_t \kappa(\mathbf{x}_t, \cdot)$ ;
    end if
end for
    
```

---

**Other Related Work.** The kernel trick implies that the inner product between any two instances can be replaced by a kernel function, i.e.,  $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i)^\top \Phi(\mathbf{x}_j)$ ,  $\forall i, j$ , where  $\Phi(\mathbf{x}_t) \in R^D$  denotes the feature mapping from the original space to a new  $D$ -dimensional space which can be infinite. Using the kernel trick, many existing linear online learning algorithms can be easily extended to their kernelized variants, such as the kernelized Perceptron

and kernelized OGD as well as kernel PA variants (Crammer et al., 2006). However, some algorithms that use complex update rules are non-trivial to be converted into kernelized versions, such as Confidence-Weighted algorithms (Dredze et al., 2008). Moreover, some online kernel learning methods also attempt to make more effective updates at each iteration. For example, Double Updating Online Learning (DUOL) (Zhao et al., 2009, 2011a; Zhao and Hoi, 2012) improves the efficacy of traditional online kernel learning methods by not only updating the weight of the newly added SV, but also the weight for one existing SV. Finally, we note one major challenge of online kernel method is the computational efficiency and scalability due to the curse of kernelization (Wang et al., 2012d). In the following, we will discuss two types of techniques to scale up kernel-based online learning methods.

### 3.6.2 SCALABLE ONLINE KERNEL LEARNING VIA BUDGET MAINTENANCE

Despite enjoying the clear advantage of accuracy performance over linear models, online kernel learning falls short in some critical drawbacks, in which one critical issue is the growing unbounded number of support vectors with increasing computational and space complexity over time. To address this challenge, a family of algorithms, termed “budget online kernel learning”, have been proposed to bound the number of SV’s with a fixed budget  $B = |\mathcal{SV}|$  using diverse budget maintenance strategies whenever the budget overflows. The general framework for budgeting strategies is shown in Algorithm 12. Most existing budget online kernel methods maintain the budget by three strategies: (i) SV Removal, (ii) SV Projection, and (iii) SV Merging. We briefly review each of them below.

---

**Algorithm 12:** Budget Online Kernel Learning
 

---

```

INIT:  $f_0 = 0$ 
for  $t = 1, 2, \dots, T$  do
    Given an incoming instance  $\mathbf{x}_t$ , predict  $\hat{y}_t = \text{sgn}(f_t(\mathbf{x}_t))$ ;
    Receive the true class label  $y_t \in \{+1, -1\}$ ;
    if update is needed then
        update the classifier from  $f_t$  to  $f_{t+\frac{1}{2}}$  and  $\mathcal{SV}_{t+\frac{1}{2}} = \mathcal{SV}_t \cup (\mathbf{x}_t, y_t)$ 
    end if
    if  $|\mathcal{SV}_{t+\frac{1}{2}}| > B$  then
        Update Support Vector Set from  $\mathcal{SV}_{t+\frac{1}{2}}$  to  $\mathcal{SV}_{t+1}$  such that  $|\mathcal{SV}_{t+1}| = B$ 
        Update the classifier from  $f_{t+\frac{1}{2}}$  to  $f_{t+1}$ 
    end if
end for
    
```

---

**SV Removal.** This strategy maintains the budget by a simple and efficient way: 1) update the classifier by adding a new SV whenever necessary (depending on the prediction mistake/loss); 2) if the SV size exceeds the budget, discard one of existing SV’s and update the classifier accordingly. To achieve this, we need to address the following concerns: (i) how to update the classifier and (ii) how to choose one of existing SV’s for removal.

The first step of updating classifiers depends on which online learning method is used. For example, the Perceptron algorithm has been used in RBP (Cavallanti et al., 2007),

Forgetron (Dekel et al., 2005), and Budget Perceptron (Crammer et al., 2003). The OGD algorithm has been adopted by BOGD (Zhao et al., 2012) and BSGD+ removal (Wang et al., 2012d), while PA has been used by BPA-S (Wang and Vucetic, 2010).

The second step of SV removal is to find one of existing SV's, denoted as  $(\mathbf{x}_{del}, y_t)$ , to be removed by minimizing the impact of the resulting classifier. One simple way is to randomly discard one of existing SV's uniformly with probability  $\frac{1}{B}$ , as adopted by RBP (Cavallanti et al., 2007) and BOGD (Zhao et al., 2012). Instead of choosing randomly, another way as used in "Forgetron" (Dekel et al., 2005) is to discard the oldest SV by assuming an older SV is less representative for the distribution of fresh training data streams. Despite enjoying the merits of simplicity and high efficiency, these methods are often too simple to achieve satisfactory learning results.

To optimize the performance, some approaches have tried to perform exhaustive search in deciding the best SV for removal. For instance, the Budget Perceptron algorithm (Crammer et al., 2003) searches for one SV that is classified with high confidence by the classifier:

$$y_{del}(f_{t+\frac{1}{2}}(\mathbf{x}_{del}) - \alpha_{del}\kappa(\mathbf{x}_{del}, \mathbf{x}_{del})) > \beta$$

where  $\beta > 0$  is a fixed tolerance parameter. BPA-S shares the similar idea of exhaustive search. For every  $r \in [B]$ , a candidate classifier  $f^r = f_{t+\frac{1}{2}} - \alpha_r\kappa(\mathbf{x}_r, \cdot)$  is generated by discarding the  $r$ -th SV from  $f_{t+\frac{1}{2}}$ . By comparing the  $B$  candidate classifiers, the algorithm selects the one that minimizes the current objective function of PA:

$$f_{t+1} = \operatorname{argmin}_{r \in [B]} \frac{1}{2} \|f^r - f_t\|_{\mathcal{H}}^2 + C\ell_t(f^r)$$

where  $C > 0$  is the regularization parameter of PA to balance aggressiveness and passiveness. Comparing the principles of different SV removal strategies, we observe that a simple rule may not always generate satisfactory accuracy, while an exhaustive search often incurs non-trivial computational overhead, which again may limit the application to large-scale problems. When deploying a solution in practice, one would need to balance the trade-off between effectiveness and efficiency.

**SV Projection.** SV Projection strategy first appeared in (Orabona et al., 2009) where two new algorithms, Projectron and Projectron++, were proposed, which significantly outperformed the previous SV removal based algorithms such RBP and Forgetron. The SV projection method follows the setting of SV removal and identifies a support vector for removal during the update of the model. It then chooses a subset of  $\mathcal{SV}$  as the projection base, which will be denoted by  $\mathcal{P}$ . Following this, a linear combination of kernels in  $\mathcal{P}$  is used to approximate the removed SV. The procedure of finding the optimal linear combination can be formulated as a convex optimization of minimizing the projection error:

$$\beta = \operatorname{argmin}_{\beta \in \mathbb{R}^{|\mathcal{P}|}} E_{proj} = \operatorname{argmin}_{\beta \in \mathbb{R}^{|\mathcal{P}|}} \|\alpha_{del}\kappa(\mathbf{x}_{del}, \cdot) - \sum_{i \in \mathcal{P}} \beta_i \kappa(\mathbf{x}_i, \cdot)\|_{\mathcal{H}}^2$$

Finally, the classifier is updated by combining this result with the original classifier:

$$f_{t+1} = f_{t+\frac{1}{2}} - \alpha_{del}\kappa(\mathbf{x}_{del}, \cdot) + \sum_{i \in \mathcal{P}} \beta_i \kappa(\mathbf{x}_i, \cdot)$$

There are several algorithms adopting the projection strategy, for example Projectron, Projectron++, BPA-P, BPA-NN (Wang and Vucetic, 2010) and BSGD+Project (Wang et al., 2012d). These methods differ in a few aspects. First, the update rules are based on different online learning algorithms. Generally speaking, PA based and OGD based tend to outperform Perceptron based algorithms because of their effective update. Second, the choice of discarded SV is different. Since projection itself is relative slow, exhaustive search based algorithms (BPA-NN, BPA-P) are extremely time consuming. Thus algorithms with simple selecting rules are preferred (Projectron, Projectron++, BSGD+Project). Third, the choice of projection base set  $\mathcal{P}$  is different. In Projectron, Projectron++, BPA-P and BSGD+Project, the discarded SV is projected onto the whole SV set, i.e.  $\mathcal{P} = \mathcal{SV}$ . While in BPA-NN,  $\mathcal{P}$  is only a small subset of  $\mathcal{SV}$ , made up of the nearest neighbors of the discarded SV  $(\mathbf{x}_{del}, y_{del})$ . In general, a larger projection base set implies a more complicated optimization problem and thus more time costs. The research direction of SV projection based budget learning is to find a proper way of selecting  $\mathcal{P}$  so that the algorithm achieves the minimized projection error with a relative small projection base set.

**SV Merging.** Wang et al. (2012d) proposed a SV merging method called “BSGD+Merge” which replaces the sum of two SV’s  $\alpha_m \kappa(\mathbf{x}_m, \cdot) + \alpha_n \kappa(\mathbf{x}_n, \cdot)$  by a newly created SV  $\alpha_z \kappa(\mathbf{z}, \cdot)$ , where  $\alpha_m$ ,  $\alpha_n$  and  $\alpha_z$  are the corresponding coefficients of  $\mathbf{x}_m$ ,  $\mathbf{x}_n$  and  $\mathbf{z}$ . Following the previous discussion, the goal of online budget learning through SV merging strategy is to find the optimal  $\alpha_z \in \mathbb{R}$  and  $\mathbf{z} \in \mathbb{R}^d$  that minimizes the gap between  $f_{t+1}$  and  $f_{t+\frac{1}{2}}$ .

As it is relatively complicated to optimize the two terms simultaneously, the optimization is divided into two steps. First, assuming the coefficient of  $\mathbf{z}$  is  $\alpha_m + \alpha_n$ , this algorithm tries to create the optimal SV that minimizes the merging error as follows

$$\min_{\mathbf{z}} \|(\alpha_m + \alpha_n) \kappa(\mathbf{z}, \cdot) - (\alpha_m \kappa(\mathbf{x}_m, \cdot) + \alpha_n \kappa(\mathbf{x}_n, \cdot))\|$$

The solution is  $\mathbf{z} = h\mathbf{x}_m + (1-h)\mathbf{x}_n$ , where  $0 < h < 1$  is a real number that can be found by a line search method. This solution indicates that the optimal created SV lies on the line connecting  $\mathbf{x}_m$  and  $\mathbf{x}_n$ . After obtaining the optimal created SV  $\mathbf{z}$ , the next step is to find the optimal coefficient  $\alpha_z$ , which can be formulated as

$$\min_{\alpha_z} \|(\alpha_z \kappa(\mathbf{z}, \cdot) - (\alpha_m \kappa(\mathbf{x}_m, \cdot) + \alpha_n \kappa(\mathbf{x}_n, \cdot)))\|.$$

The solution becomes  $\alpha_z = \alpha_m \kappa(\mathbf{x}_m, \mathbf{z}) + \alpha_n \kappa(\mathbf{x}_n, \mathbf{z})$ . The remaining problem is which two SV’s  $\mathbf{x}_m$  and  $\mathbf{x}_n$  should be merged. The ideal solution is to find the optimal pair with the minimal merging error through exhaustive search, which however requires  $O(B^2)$  time complexity. How to find the optimal pair efficiently remains an open challenge.

**Summary.** Among various algorithms of budget online kernel learning using budget maintenance, the key differences are their updating rules and budget maintenance strategies. Table 1 gives a summary of different algorithms and their properties. In addition to the previous budget kernel learning algorithms, there are also some other works in online kernel learning. For example, some studies (Lu et al., 2016a; Zhang et al., 2012) introduce the idea of sparse kernel learning to reduce the number of SV’s in the online-to-batch-conversion problem, where an online algorithm can be used to train a kernel model efficiently for the batch setting (See Section 3.7).

Table 1: Comparisons of different budget online kernel learning algorithms.

Algorithms	Update Strategy	Budget Strategy	Update Time	Space
Stopton (Orabona et al., 2009)	Perceptron	Stop	$O(1)$	$O(B)$
Tighter Perceptron (Weston et al., 2005)	Perceptron	Removal	$O(B^2)$	$O(B)$
Tightest Perceptron (Wang and Vucetic, 2009)	Perceptron	Removal	$O(B^2)$	$O(B)$
Budget Perceptron (Crammer et al., 2003)	Perceptron	Removal	$O(B^2)$	$O(B)$
RBP (Cavallanti et al., 2007)	Perceptron	Removal	$O(B)$	$O(B)$
Forgetron (Dekel et al., 2005)	Perceptron	Removal	$O(B)$	$O(B)$
BOGD (Zhao et al., 2012)	OGD	Removal	$O(B)$	$O(B)$
BPA-S (Wang and Vucetic, 2010)	PA	Removal	$O(B)$	$O(B)$
BSGD+removal (Wang et al., 2012d)	OGD	Removal	$O(B)$	$O(B)$
Projectron (Orabona et al., 2009)	Perceptron	Projection	$O(B^2)$	$O(B^2)$
Projectron++ (Orabona et al., 2009)	Perceptron	Projection	$O(B^2)$	$O(B^2)$
BPA-P (Wang and Vucetic, 2010)	PA	Projection	$O(B^3)$	$O(B^2)$
BPA-NN (Wang and Vucetic, 2010)	PA	Projection	$O(B)$	$O(B)$
BSGD+projection (Wang et al., 2012d)	OGD	Projection	$O(B^2)$	$O(B^2)$
BSGD+merging (Wang et al., 2012d)	OGD	Merging	$O(B)$	$O(B)$

### 3.6.3 SCALABLE ONLINE KERNEL LEARNING VIA FUNCTIONAL APPROXIMATION

In contrast to the previous budget online kernel learning methods using budget maintenance strategies to guarantee efficiency and scalability, another emerging and promising strategy is to explore functional approximation techniques for achieving scalable online kernel learning (Wang et al., 2013b; Lu et al., 2015).

The key idea is to construct a kernel-induced feature representation  $\mathbf{z}(\mathbf{x}) \in \mathbb{R}^D$  such that the inner product of instances in the new feature space can effectively approximate the kernel function:

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) \approx \mathbf{z}(\mathbf{x}_i)^\top \mathbf{z}(\mathbf{x}_j)$$

Using the above approximation, the predictive model with kernels can be rewritten as follows:

$$f(\mathbf{x}) = \sum_{i=1}^B \alpha_i \kappa(\mathbf{x}_i, \mathbf{x}) \approx \sum_{i=1}^B \alpha_i \mathbf{z}(\mathbf{x}_i)^\top \mathbf{z}(\mathbf{x}) = \mathbf{w}^\top \mathbf{z}(\mathbf{x})$$

where  $\mathbf{w} = \sum_{i=1}^B \alpha_i \mathbf{z}(\mathbf{x}_i)$  denotes the weight vector to be learned in the new feature space.

As a consequence, solving a regular online kernel classification task can be turned into a linear online classification task on the new feature space derived from the kernel approximation. For example, the methods of online kernel learning with kernel approximation in (Wang et al., 2013b; Lu et al., 2015) integrate some existing online learning algorithms (e.g., OGD) with kernel approximation techniques (Lin and Chen, 2011; Sonnenburg and Franc, 2010; Chang et al., 2010) to derive scalable online kernel learning algorithms, including Fourier Online Gradient Descent (FOGD) that explores random Fourier features for kernel functional approximation (Rahimi and Recht, 2007), and Nyström Online Gradient Descent (NOGD) that explores Nyström low-rank matrix approximation methods for approximating large-scale kernel matrix (Williams and Seeger, 2001). A recent work, Dual Space Gradient Descent (Le et al., 2016; Nguyen et al., 2017) updates the model as the RBP algorithm, but also builds an FOGD model using the discarded SV's. The final prediction is the combination of the two models.

### 3.6.4 ONLINE MULTIPLE KERNEL LEARNING

Traditional online kernel methods usually assume a predefined good kernel is given prior to the online learning task. Such approaches could be restricted since it is often hard to choose a good kernel prior to the learning task. To overcome the drawback, Online Multiple Kernel Learning (OMKL) aims to combine multiple kernels automatically for online learning tasks without fixing any predefined kernel. In the following, we begin by introducing some basics of batch Multiple Kernel Learning (MKL) (Sonnenburg et al., 2006).

Given a training set  $\mathcal{D} = \{(\mathbf{x}_t, y_t), t = 1, \dots, T\}$  where  $\mathbf{x}_t \in \mathbb{R}^d$ ,  $y_t \in \{-1, +1\}$ , and a set of  $m$  kernel functions  $\mathcal{K} = \{\kappa_i : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}, i = 1, \dots, m\}$ . MKL learns a kernel-based prediction function by identifying an optimal combination of the  $m$  kernels, denoted by  $\theta = (\theta_1, \dots, \theta_m)$ , to minimize the margin-based classification error, which can be cast into the optimization below:

$$\min_{\theta \in \Delta} \min_{f \in \mathcal{H}_{K(\theta)}} \frac{1}{2} \|f\|_{\mathcal{H}_{K(\theta)}}^2 + C \sum_{t=1}^n \ell(f(\mathbf{x}_t), y_t) \quad (11)$$

where  $\Delta = \{\theta \in \mathbb{R}_+^m | \theta^\top \mathbf{1}_m = 1\}$ ,  $K(\theta)(\cdot, \cdot) = \sum_{i=1}^m \theta_i \kappa_i(\cdot, \cdot)$ ,  $\ell(f(\mathbf{x}_t), y_t) = \max(0, 1 - y_t f(\mathbf{x}_t))$ . In the above, we use notation  $\mathbf{1}_T$  to represent a vector of  $T$  dimensions with all its elements being 1. It can also be cast into the following mini-max optimization problem:

$$\min_{\theta \in \Delta} \max_{\alpha \in \Xi} \left\{ \alpha^\top \mathbf{1}_T - \frac{1}{2} (\alpha \circ \mathbf{y})^\top \left( \sum_{i=1}^m \theta_i K^i \right) (\alpha \circ \mathbf{y}) \right\} \quad (12)$$

where  $K^i \in \mathbb{R}^{T \times T}$  with  $K_{j,l}^i = \kappa_i(\mathbf{x}_j, \mathbf{x}_l)$ ,  $\Xi = \{\alpha | \alpha \in [0, C]^T\}$ , and  $\circ$  defines the element-wise product between two vectors. The above batch MKL optimization has been extensively studied (Rakotomamonjy et al., 2008; Xu et al., 2008), but an efficient solution remains an open challenge.

Some efforts of online MKL studies (Jie et al., 2010; Martins et al., 2011) have attempted to solve batch MKL optimization via online learning. Unlike these approaches that are mainly concerned in optimizing the optimal kernel combination as regular MKL, another framework of *Online Multiple Kernel Learning* (OMKL) (Jin et al., 2010; Hoi et al., 2013; Yang et al., 2012) is focused on exploring effective online combination of multiple kernel classifiers via a significantly more efficient and scalable way. Specifically, the OMKL in (Jin et al., 2010; Hoi et al., 2013) learns a kernel-based prediction function by selecting a subset of predefined kernel functions in an online learning fashion, which is in general more challenging than typical online learning because both the kernel classifiers and the subset of selected kernels are unknown, and more importantly the solutions to the kernel classifiers and their combination weights are correlated. (Hoi et al., 2013) proposed novel algorithms based on the fusion of two types of online learning algorithms, i.e., the *Perceptron* algorithm that learns a classifier for a given kernel, and the *Hedge* algorithm (Freund and Schapire, 1997) that combines classifiers by linear weights. Some stochastic selection strategies were also proposed by randomly selecting a subset of kernels for combination and model updating to further improve the efficiency. These methods were later extended for regression (Sahoo et al., 2014), learning from data with time-sensitive patterns (Sahoo et al., 2016a) and imbalanced data streams (Sahoo et al., 2016b). In addition, there have been budgeting approaches to make OMKL scalable (Lu et al., 2018).

### 3.7 Online to Batch Conversion

Let us denote by  $\mathcal{A}$  an online learning algorithm for the purpose of training a binary classifier from a sequence of training examples. On each round, the algorithm receives an instance  $\mathbf{x}_t \in \mathbb{R}^d$ , the algorithm chooses a vector  $\mathbf{w}_t \in \mathcal{S} \subseteq \mathbb{R}^d$  to predict the class label of the instance, i.e.,  $\hat{y}_t = \text{sgn}(\mathbf{w}_t^\top \mathbf{x}_t)$ . After that, the environment responds by disclosing the true label  $y_t$  and some convex loss function  $\ell(\mathbf{w}; (\mathbf{x}_t, y_t))$ , and the algorithm suffers a loss  $\ell_t(\mathbf{w}_t)$  at the end of this round. For such a setting, consider a sequence of  $T$  rounds  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$ , the online algorithm aims to minimize the following regret

$$\text{Reg}_{\mathcal{A}}(T) = \sum_{t=1}^T \ell(\mathbf{w}_t; (\mathbf{x}_t, y_t)) - \min_{\mathbf{w} \in \mathcal{S}} \sum_{t=1}^T \ell(\mathbf{w}; (\mathbf{x}_t, y_t))$$

However, for batch training setting, we are more interested in finding a model  $\hat{\mathbf{w}}$  with good generalization ability, i.e., we want to achieve a small excess risk defined as

$$R(\hat{\mathbf{w}}) - \min_{\mathbf{w} \in \mathcal{S}} R(\mathbf{w})$$

where the generalization risk is  $R(\mathbf{w}) = \mathbb{E}_{(\mathbf{x}, y)}[\ell(\mathbf{w}; (\mathbf{x}, y))]$ , and  $(\mathbf{x}, y)$  satisfies a fixed unknown distribution. Therefore, we would like to study the generalization performance of online algorithms through the Online to Batch Conversion (Cesa-Bianchi et al., 2004; Wang et al., 2012c; Kakade and Tewari, 2009; Levy, 2017), where the conversion relates the regret of the online algorithm to its generalization performance.

#### 3.7.1 A GENERAL CONVERSION THEORY

We now consider the generalization ability of online learning under the assumption that the loss  $\ell(\mathbf{w}; (\mathbf{x}, y))$  is strongly convex, which is reasonable as many loss functions (e.g., square loss) are strongly convex, and even if some loss (e.g., hinge loss) is not strongly convex, we can impose some regularization term (e.g.,  $\frac{1}{2} \|\cdot\|$ ) to achieve strong convexity. We denote the dual norm of  $\|\cdot\|$  as  $\|\cdot\|_*$ , where  $\|\mathbf{v}\|_* = \sup_{\|\mathbf{w}\| \leq 1} \mathbf{v}^\top \mathbf{w}$ . Let  $Z = (\mathbf{x}, y)$  be a random variable taking values in some space  $\mathcal{Z}$ . Our goal is to minimize  $R(\mathbf{w}) = \mathbb{E}_Z[\ell(\mathbf{w}; Z)]$  over  $\mathbf{w} \in \mathcal{S}$ . Specifically, we assume the loss  $\ell : \mathcal{S} \times \mathcal{Z} \rightarrow [0, B]$  satisfies the following assumption:

**Assumption LIST**(Lipschitz and STrongly convex assumption) For all  $z \in \mathcal{Z}$ , the function  $\ell_z(\mathbf{w}) = \ell(\mathbf{w}; z)$  is convex in  $\mathbf{w}$  and satisfies:

1.  $\ell_z$  has Lipschitz constant  $L$  w.r.t. the norm  $\|\cdot\|$ , i.e.,  $|\ell_z(\mathbf{w}) - \ell_z(\mathbf{w}')| \leq L \|\mathbf{w} - \mathbf{w}'\|$ .
2.  $\ell_z$  is  $\lambda$ -strongly convex w.r.t.  $\|\cdot\|$ , i.e.,  $\forall \theta \in [0, 1], \forall \mathbf{w}, \mathbf{w}' \in \mathcal{S}$ ,

$$\ell_z(\theta \mathbf{w} + (1 - \theta) \mathbf{w}') \leq \theta \ell_z(\mathbf{w}) + (1 - \theta) \ell_z(\mathbf{w}') - \frac{\lambda}{2} \theta (1 - \theta) \|\mathbf{w} - \mathbf{w}'\|^2.$$

For this kind of loss function, we consider an online learning setting where  $Z_1, \dots, Z_T$  are given sequentially in i.i.d. We then have

$$\mathbb{E}[\ell(\mathbf{w}; Z_t)] = \mathbb{E}[\ell(\mathbf{w}; (\mathbf{x}_t, y_t))] := R(\mathbf{w}), \quad \forall t, \mathbf{w} \in \mathcal{S}.$$



Now consider an online learning algorithm  $\mathcal{A}$ , which is initialized as  $\mathbf{w}_1$ . Whenever  $Z_t$  is given, model  $\mathbf{w}_t$  is updated to  $\mathbf{w}_{t+1}$ . Let  $E_t[\cdot]$  denote conditional expectation w.r.t.  $Z_1, \dots, Z_t$ , we have  $E_t[\ell(\mathbf{w}_t; Z_t)] = R(\mathbf{w}_t)$ . Using the above assumptions and the Freedman's inequality leads to the following theorem for the generalization ability of online learning

**Theorem 1** (*Cesa-Bianchi et al. (2004)*) *Under the assumption LIST, we have the following inequality, with probability at least  $1 - 4\delta \ln T$ ,*

$$\frac{1}{T} \sum_{t=1}^T R(\mathbf{w}_t) - R(\mathbf{w}_*) \leq \frac{\text{Reg}_{\mathcal{A}}(T)}{T} + 4\sqrt{\frac{L^2 \ln \frac{1}{\delta}}{\lambda}} \frac{\sqrt{\text{Reg}_{\mathcal{A}}(T)}}{T} + \max\left(\frac{16L^2}{\lambda}, 6B\right) \frac{\ln \frac{1}{\delta}}{T},$$

where  $\mathbf{w}_* = \arg \min_{\mathbf{w} \in \mathcal{S}} R(\mathbf{w})$ . Further, using Jensen's inequality,  $\frac{1}{T} \sum_t R(\mathbf{w}_t)$  can be replaced by  $R(\bar{\mathbf{w}}_T)$  where  $\bar{\mathbf{w}}_T = \frac{1}{T} \sum_t \mathbf{w}_t$ .

If the assumption LIST is satisfied by  $\ell_z(\mathbf{w})$ , then the Online Gradient Descent (OGD) algorithm that generates  $\mathbf{w}_1, \dots, \mathbf{w}_T$  has the following regret  $\text{Reg}_{\mathcal{A}}(T) \leq \frac{L^2}{2\lambda}(1 + \ln T)$ . Plugging this inequality back into the theorem and using  $(1 + \ln T)/(2T) \leq \ln T/T$ ,  $\forall T \geq 3$  gives the following Corollary.

**Corollary 2** *Suppose assumption LIST holds for  $\ell_z(\mathbf{w})$ . Then the Online Gradient Descent (OGD) algorithm that generates  $\mathbf{w}_1, \dots, \mathbf{w}_T$  and finally outputs  $\bar{\mathbf{w}}_T = \frac{1}{T} \sum_t \mathbf{w}_t$ , satisfies the following inequality for its generalization ability, with probability at least  $1 - 4\delta \ln T$ ,*

$$R(\bar{\mathbf{w}}_T) - R(\mathbf{w}_*) \leq \frac{L^2 \ln T}{\lambda T} + \sqrt{\ln \frac{1}{\delta}} \frac{4L^2 \sqrt{\ln T}}{\lambda T} + \max\left(\frac{16L^2}{\lambda}, 6B\right) \frac{\ln \frac{1}{\delta}}{T},$$

for any  $T \geq 3$ , where  $\mathbf{w}_* = \arg \min_{\mathbf{w} \in \mathcal{S}} R(\mathbf{w})$ .

### 3.7.2 OTHER CONVERSION THEORIES

Online to batch conversion has been studied in literature (Littlestone, 1989; Cesa-Bianchi et al., 2004; Zhang, 2005; Cesa-Bianchi and Gentile, 2008; Wang et al., 2012c). For general convex loss functions, Cesa-Bianchi et al. (2004) proved the following generalization ability of online learning algorithm with probability at least  $1 - \delta$

$$R(\bar{\mathbf{w}}_T) \leq \frac{1}{T} \sum_{t=1}^T \ell(\mathbf{w}_t; z_t) + \sqrt{\frac{2}{T} \ln \frac{1}{\delta}} = \frac{\text{Reg}_{\mathcal{A}}(T)}{T} + \min_{\mathbf{w} \in \mathcal{S}} \frac{1}{T} \sum_{t=1}^T \ell(\mathbf{w}; z_t) + \sqrt{\frac{2}{T} \ln \frac{1}{\delta}},$$

where the loss  $\ell \leq 1$ . Zhang (2005) is another work that explicitly goes by the exponential moment method to drive sharper concentration results. In addition, Cesa-Bianchi and Gentile (2008) improved their initial generalization bounds using Bernstein's inequality by assuming  $\ell(\cdot) \leq 1$ , and proves the following inequality with probability at least  $1 - \delta$

$$R(\hat{\mathbf{w}}) \leq \frac{1}{T} \sum_{t=1}^T \ell(\mathbf{w}_t; z_t) + O\left(\frac{\ln(T^2/\delta)}{T} + \sqrt{\frac{1}{T} \sum_{t=1}^T \ell(\mathbf{w}_t; z_t) \frac{\ln(T^2/\delta)}{T}}\right).$$

where  $\hat{\mathbf{w}}$  is selected from  $\mathbf{w}_1, \dots, \mathbf{w}_T$ , by minimizing a specifically designed penalized empirical risk. In particular, the generalization risk converges to  $\frac{1}{T} \sum_{t=1}^T \ell(\mathbf{w}_t; z_t)$  at rate  $O(\sqrt{\ln T^2/T})$  and vanishes at rate  $O(\ln T^2/T)$  whenever the loss  $\sum_{t=1}^T \ell(\mathbf{w}_t; z_t)$  is  $O(1)$ .

## 4. Applied Online Learning for Supervised Learning

### 4.1 Overview

In this section, we survey the most representative algorithms for a group of non-traditional online learning tasks, wherein the supervised online algorithms cannot be used directly. These algorithms are motivated by new problem settings and applications which follow the traditional online setting, where the data arrives in a sequential manner. However, there was a need to develop new algorithms which were suited to these scenarios. Our review includes cost-sensitive online learning, online multi-task learning, online multi-view learning, online transfer learning, online metric learning, online collaborative filtering, online learning structured prediction, distributed online learning, online learning with neural networks, and online portfolio selection.

### 4.2 Cost-Sensitive Online Learning

In a supervised classification task, traditional online learning methods are often designed to optimize mistake rate or equivalently classification accuracy. However, it is well-known that classification accuracy becomes a misleading metric when dealing with class-imbalanced data which is common for many real-world applications, such as anomaly detection, fraud detection, intrusion detection, etc. To address this issue, cost-sensitive online learning (Chen et al., 2017) represents a family of online learning algorithms that are designed to take care of different misclassification costs of different classes in a class-imbalanced classification task. Next, we briefly survey these algorithms.

**Perceptron Algorithms with Uneven Margin (PAUM)** PAUM (Li et al., 2002) is a cost-sensitive extension of Perceptron (Rosenblatt, 1958) and the Perceptron with Margins (PAM) algorithms (Krauth and M  zard, 1987). Perceptron makes an update only when there is a mistake, while PAM tends to make more aggressive updates by checking the margin instead of mistake. PAM makes an update whenever  $y_t \mathbf{w}_t^\top \mathbf{x}_t \leq \tau$ , where  $\tau \in \mathbb{R}^+$  is a fixed parameter controlling the aggressiveness. To deal with class imbalance, PAUM extends PAM via an uneven margin setting, i.e., employing different margin parameters for the two classes:  $\tau_+$  and  $\tau_-$ . Consequently, the update becomes  $y_t \mathbf{w}_t^\top \mathbf{x}_t \leq \tau_{y_t}$ . By properly adjusting the two parameters, PAUM achieves cost-sensitive updating effects for different classes. One of major limitations with PAUM is that it does not directly optimize a predefined cost-sensitive measure, thus, it does not fully resolve the cost-sensitive challenge.

**Cost-sensitive Passive Aggressive (CPA)** CPA (Crammer et al., 2006) was proposed as a cost-sensitive variant of the PA algorithms. It was originally designed for multi-class classification by the following prediction rule:  $\hat{y}_t = \arg \max_y (\mathbf{w}_t^\top \Phi(\mathbf{x}_t, y))$ , where  $\Phi$  is a feature mapping function that maps  $\mathbf{x}_t$  to a new feature according to the class  $y$ . For simplicity, we restrict the discussion on the binary classification setting. Using  $\Phi(\mathbf{x}, y) = \frac{1}{2} y \mathbf{x}$ , we will map the formulas to our setting. The prediction rule is:  $\hat{y}_t = \text{sgn}(\mathbf{w}_t^\top \mathbf{x}_t)$ . We define the cost-sensitive loss as

$$\ell(\mathbf{w}, \mathbf{x}, y) = \mathbf{w} \cdot \Phi(\mathbf{x}, \hat{y}) - \mathbf{w} \cdot \Phi(\mathbf{x}, y) + \sqrt{\rho(y, \hat{y})},$$

where  $\rho(y_1, y_2)$  is the function define to distinguish the different cost of different kind misclassifications and we have assumed  $\rho(y, y) = 0$ . When being converted to binary setting,

the loss becomes

$$\ell(\mathbf{w}, \mathbf{x}, y) = \begin{cases} 0 & y_t = \hat{y} \\ |\mathbf{w}^\top \mathbf{x}| + \sqrt{\rho(y, \hat{y})} & y_t \neq \hat{y} \end{cases}$$

The mistake depends on the prediction confidence and the loss type. We omit the detailed update steps since it follows the similar optimization as PA learning as discussed before. Similar to PAUM, this algorithm also is limited in that it does not optimize a cost-sensitive measure directly.

**Cost-Sensitive Online Gradient Descent (CSOGD)** Unlike traditional OGD algorithms that often optimize accuracy, CSOGD (Wang et al., 2014b, 2012a; Zhao et al., 2015) applies OGD to directly optimize two cost-sensitive measures:

- (1) maximizing the weighted sum of *sensitivity* and *specificity*, i.e.,  $sum = \eta_p \times sensitivity + \eta_n \times specificity$ , where the two weights satisfy  $0 \leq \eta_p, \eta_n \leq 1$  and  $\eta_p + \eta_n = 1$ .
- (2) minimizing the weighted *misclassification cost*, i.e.,  $cost = c_p \times M_p + c_n \times M_n$ , where  $M_p$  and  $M_n$  are the number of false negatives and false positives respectively,  $0 \leq c_p, c_n \leq 1$  are the cost parameters for positive and negative classes, respectively, and we assume  $c_p + c_n = 1$ .

The objectives can be equivalently reformulated into the following objective:

$$\sum_{y_t=+1} \rho \mathbf{I}_{(y_t \mathbf{w} \cdot \mathbf{x}_t < 0)} + \sum_{y_t=-1} \mathbf{I}_{(y_t \mathbf{w} \cdot \mathbf{x}_t < 0)}$$

where we set  $\rho = \frac{\eta_p T_n}{\eta_n T_p}$  when maximizing the weighted sum,  $T_p$  and  $T_n$  are the number of positive and negative instances respectively; when minimizing the weighted misclassification cost, we instead set  $\rho = \frac{c_p}{c_n}$ . The objective is however non-convex, making it hard to optimize directly. Instead of directly optimizing the non-convex objective, we attempt to optimize a convex surrogate. Specifically, we replace the indicator function  $\mathbf{I}_{(\cdot)}$  by a convex surrogate, and attempt to optimize either one of the following modified hinge-loss functions at each online iteration:

$$\ell^I(\mathbf{w}; (\mathbf{x}, y)) = \max(0, \rho * \mathbf{I}_{(y=1)} + \mathbf{I}_{(y=-1)} - y(\mathbf{w} \cdot \mathbf{x}))$$

$$\ell^{II}(\mathbf{w}; (\mathbf{x}, y)) = (\rho * \mathbf{I}_{(y=1)} + \mathbf{I}_{(y=-1)}) * \max(0, 1 - y(\mathbf{w} \cdot \mathbf{x}))$$

One can then derive cost-sensitive ODG (CSOGD) algorithms by applying OGD to optimize either one of the above loss functions. The detailed algorithms can be found in (Wang et al., 2014b). Two recent works extend the problem setting to cost-sensitive classification of multi-class problem (Wang et al., 2016b; Zhang et al., 2016b). Further there are efforts to do cost-sensitive online learning with kernels (Zhao and Hoi, 2013; Hu et al., 2015).

**Online AUC Maximization** Instead of optimizing accuracy, some online learning studies have attempted to directly optimize the Area Under the ROC curve (AUC), i.e.,

$$\text{AUC}(\mathbf{w}) = \frac{\sum_{i=1}^{T_+} \sum_{j=1}^{T_-} \mathbb{I}_{\mathbf{w} \cdot \mathbf{x}_i^+ > \mathbf{w} \cdot \mathbf{x}_j^-}}{T_+ T_-} = 1 - \frac{\sum_{i=1}^{T_+} \sum_{j=1}^{T_-} \mathbb{I}_{\mathbf{w} \cdot \mathbf{x}_i^+ \leq \mathbf{w} \cdot \mathbf{x}_j^-}}{T_+ T_-}$$

where  $\mathbf{x}^+$  is a positive instance,  $\mathbf{x}^-$  is a negative instance,  $T_+$  is the total number of positive instances and  $T_-$  is the total number of negative instances. AUC measures the probability for a randomly drawn positive instance to have a higher decision value than a randomly sampled negative instance, and it is widely used in many applications. Optimizing AUC online is however very challenging.

First of all, in the objective, the term  $\sum_{i=1}^{T_+} \sum_{j=1}^{T_-} \mathbb{I}_{\mathbf{w} \cdot \mathbf{x}_i^+ \leq \mathbf{w} \cdot \mathbf{x}_j^-}$  is non-convex. A common way is to replace the indicator function by a convex surrogate, e.g., a hinge loss function

$$\ell(\mathbf{w}, \mathbf{x}_i^+ - \mathbf{x}_j^-) = \max\{0, 1 - \mathbf{w}(\mathbf{x}_i^+ - \mathbf{x}_j^-)\}$$

Consequently, the goal of online AUC maximization is equivalent to minimizing the accumulated loss  $\mathcal{L}_t(\mathbf{w})$  over all previous iterations, where the loss at the  $t$ -th iteration is

$$\mathcal{L}_t(\mathbf{w}) = \mathbb{I}_{y_t=1} \sum_{\tau=1}^{t-1} \mathbb{I}_{y_\tau=-1} \ell(\mathbf{w}, \mathbf{x}_t - \mathbf{x}_\tau) + \mathbb{I}_{y_t=-1} \sum_{\tau=1}^{t-1} \mathbb{I}_{y_\tau=1} \ell(\mathbf{w}, \mathbf{x}_\tau - \mathbf{x}_t)$$

The above takes the sum of the pairwise hinge loss between the current instance  $(\mathbf{x}_t, y_t)$  and all the received instances with the opposite class  $-y_t$ . Despite being convex, it is however impractical to directly optimize the above objective in online setting since one would need to store all the received instances and thus lead to growing computation and memory cost.

The Online AUC Maximization method in (Zhao et al., 2011b) proposed a novel idea of exploring *reservoir sampling* techniques to maintain two buffers,  $B_+$  and  $B_-$  of size  $N_+$  and  $N_-$ , which aim to store a sketch of historical instances. Specifically, when receiving instance  $(\mathbf{x}_t, y_t)$ , it will be added to buffer  $B_{y_t}$  whenever it is not full, i.e.  $|B_{y_t}| < N_{y_t}$ . Otherwise,  $\mathbf{x}_t$  randomly replaces one instance in the buffer with probability  $\frac{N_{y_t}}{N_{y_t}^{t+1}}$ , where  $N_{y_t}^{t+1}$  is the total number of instances with class  $y_t$  received so far. Reservoir sampling is able to guarantee the instances in the buffers maintain an unbiased sampling of the original full dataset. As a result, the loss  $\mathcal{L}_t(\mathbf{w})$  can be approximated by only considering the instances in the buffers, and the classifier  $\mathbf{w}$  can be updated by either OGD or PA algorithms.

*Others.* To improve the study in (Zhao et al., 2011b), a number of following studies have attempted to make improvements from different aspects (Yi Ding, 2017). For example, the study in (Wang et al., 2012c) generalized online AUC maximization as online learning with general pairwise loss functions, and offered new generalization bounds for online AUC maximization algorithms similar to (Zhao et al., 2011b). The bounds were further improved by (Kar et al., 2013) which employs a generic decoupling technique to provide Rademacher complexity-based generalization bounds. In addition, the work in (Gao et al., 2013) overcomes the buffering storage cost by developing a regression-based algorithm which only needs to maintain the first and second-order statistics of training data in memory, making the resulting storage requirement independent from the training size. The very recent work in (Ding et al., 2015) presented a new second-order AUC maximization method by improving the convergence using the adaptive gradient algorithm. The stochastic online AUC maximization (SOLAM) algorithm (Ying et al., 2016) formulates the online AUC maximization as a stochastic saddle point problem and greatly reduces the memory cost.

### 4.3 Online Multi-task Learning

Multi-task Learning (Caruana, 1998) is an approach that learns a group of related machine learning tasks together. By considering the relationship between different tasks, multi-task learning algorithms are expected to achieve better performance than algorithms that learn each task individually. Batch multi-task learning problems are usually solved by transfer learning methods (Pan and Yang, 2010) which transfer the knowledge learnt from one task to another similar tasks. In Online Multi-task Learning (OML) (Dekel et al., 2006; Li et al., 2010a, 2014), however, the tasks have to be solved in parallel with instances arriving sequentially, which makes the problem more challenging.

During time  $t$ , each of the task  $i \in \{1, \dots, K\}$  receives an instance  $\mathbf{x}_{i,t} \in \mathbb{R}^{d_i}$ , where  $d_i$  is the feature dimension of task  $i$ . The algorithm then makes a prediction for each task  $i$  based on the current model  $\mathbf{w}_{i,t}$  as  $\hat{y}_{i,t} = \text{sign}(\mathbf{w}_{i,t}^\top \mathbf{x}_{i,t})$ . After making the prediction, the true labels  $y_{i,t}$  are revealed and we get a loss function vector  $\ell_{i,t} \in \mathbb{R}_+^K$ . Finally, the models are updated by considering the loss vector and task relationship.

A straightforward baseline algorithm is to parallel update all the classifiers  $\mathbf{w}_i, i \in \{1, \dots, K\}$ . OML algorithm should utilize the relationships between tasks to achieve higher accuracy compared with the baseline. The multitask Perceptron algorithm (Cavallanti et al., 2010) is a pioneering work of OML that considers the inter-task relationship. Assuming that a matrix  $A \in \mathbb{R}^{K \times K}$  is known and fixed, we can update the model  $\mathbf{w}_i$  when an instance  $\mathbf{x}_{j,t}$  for task  $j$  is received as follows:

$$\mathbf{w}_{i,t+1} = \mathbf{w}_{i,t} + y_{j,t} A_{j,i}^{-1} \mathbf{x}_{j,t}$$

Other approaches in (Saha et al., 2011; Wang et al., 2013a) learned to optimize the relationship matrix, which also offers the flexibility of using a time-varying relationship.

Apart from learning the relationship explicitly, another widely used approach in OML field is to add some structure regularization terms to the original objective function (Evgeniou and Pontil, 2004; Yang et al., 2010; Kumar and Daumé III, 2012). For example, we may assume that each model is made up of two parts, a shared part across all tasks  $\mathbf{w}_0$  and an individual part  $\mathbf{v}_i$ , i.e.,  $\mathbf{w}_i = \mathbf{w}_0 + \mathbf{v}_i$  where the common part helps to take advantage of the task similarity. Now the regularized loss becomes

$$\sum_{i=1}^K (\ell_{i,t} + \|\mathbf{v}_i\|_2^2) + \lambda \|\mathbf{w}_0\|_2^2$$

It was also improved using more complex inter-task relationship (Murugesan et al., 2016).

### 4.4 Online Multi-view Learning

Multi-view learning deals with problems where data are collected from diverse domains or obtained from various feature extractors. By exploring features from different views, multi-view learning algorithms are usually more effective than single-view learning. In literature, there many surveys that offer comprehensive summary of state-of-the-art methods in multi-view learning in batch setting (Sun, 2013; Xu et al., 2013; Li et al., 2016c), while few works tried to address this problem in online setting.

**Two-view PA** We first introduce a seminal work, the two-view online passive aggressive learning (Two-view PA) algorithm (Nguyen et al., 2012), which is motivated by the famous single-view PA algorithm (Crammer et al., 2006) and the two-view SVM algorithm (Farquhar et al., 2006) in batch setting.

During each iteration  $t$ , the algorithm receives an instance  $(\mathbf{x}_t^A, \mathbf{x}_t^B, y_t)$ , where  $\mathbf{x}_t^A \in \mathbb{R}^n$  is the feature vector in the first view,  $\mathbf{x}_t^B \in \mathbb{R}^m$  is for the second view and  $y_t \in \{1, -1\}$  is the label. The goal is to learn two classifiers  $\mathbf{w}^A \in \mathbb{R}^n$  and  $\mathbf{w}^B \in \mathbb{R}^m$ , each for one view, and make accuracy prediction with their combination

$$\hat{y}_t = \text{sign}(\mathbf{w}^A \cdot \mathbf{x}_t^A + \mathbf{w}^B \cdot \mathbf{x}_t^B).$$

Thus the hinge loss at iteration  $t$  is redefined as

$$\ell_t(\mathbf{w}_t^A, \mathbf{w}_t^B) = \max(0, 1 - \frac{1}{2}y_t(\mathbf{w}_t^A \cdot \mathbf{x}_t^A + \mathbf{w}_t^B \cdot \mathbf{x}_t^B))$$

In the single-view PA algorithm, the objective function in each iteration is a balance between two desires: minimizing the loss function at the current instance and minimizing the change made to the classifier. While to utilize the special information in the multi-view data, an additional term that measures the agreement between two terms is added. Thus, the optimization is as follows,

$$\begin{aligned} (\mathbf{w}_{t+1}^A, \mathbf{w}_{t+1}^B) = \arg \min_{\mathbf{w}^A, \mathbf{w}^B} & \frac{1}{2} \|\mathbf{w}^A - \mathbf{w}_t^A\|_2^2 + \frac{1}{2} \|\mathbf{w}^B - \mathbf{w}_t^B\|_2^2 \\ & + C\ell_t(\mathbf{w}_t^A, \mathbf{w}_t^B) + \gamma|y_t\mathbf{w}^A \cdot \mathbf{x}_t^A - y_t\mathbf{w}^B \cdot \mathbf{x}_t^B| \end{aligned}$$

where  $\gamma$  and  $C$  are weight parameters. Fortunately, this optimization problem has a closed form solution.

**Other related works:** Other than solving classification tasks, online multi-view learning has been explored for solving similarity learning or distance metric learning, such as Online multimodal deep similarity learning (Wu et al., 2013) and online multi-modal distance metric learning (Wu et al., 2016).

## 4.5 Online Transfer Learning

Transfer learning aims to address the machine learning tasks of building models in a new target domain by taking advantage of information from another existing source domain through knowledge transfer. Transfer learning is important for many applications where training data in a new domain may be limited or too expensive to collect. There are two different problem settings, *homogeneous* setting where the target domain shares the same feature space as the old/source one, and *heterogeneous* setting where the feature space of the target domain is different from that of the source domain. Although several surveys on transfer learning are available (Sousa et al.; Pan and Yang, 2010), most of the referred algorithms are in batch setting.

Online Transfer Learning (OTL) algorithms aim to learn a classifier  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  from a well-trained classifier  $h : \mathbb{R}^{d'} \rightarrow \mathbb{R}$  in the source domain and a group of sequentially arriving instances  $\mathbf{x}_t \in \mathbb{R}^d, t = 1, \dots, T$  in the target domain. For conciseness, we will use the previous notations for the online classification task. We first introduce a pioneer work of OTL (Zhao et al., 2014; Zhao and Hoi, 2010).

**Homogeneous Setting** One key challenge of this task is to address the concept drifting issue that often occurs in this scenario. The algorithm in homogeneous setting ( $d = d'$ ) is based on the ensemble learning approach. At time  $t$ , an instance  $\mathbf{x}_t$  is received. The algorithm makes a prediction based on the weighted average of the classifier in the source domain  $h(\mathbf{x}_t)$  and the current classifier in the target domain  $f_t(\mathbf{x}_t)$ ,

$$\hat{y}_t = \text{sgn}(w_{1,t}\Pi(h(\mathbf{x}_t)) + w_{2,t}\Pi(f_t(\mathbf{x}_t)) - \frac{1}{2})$$

where  $w_{1,t} > 0, w_{2,t} > 0$  are the weights for the two functions and need to be updated during each iteration.  $\Pi$  is a normalization function, i.e.  $\Pi(a) = \max(0, \min(1, \frac{a+1}{2}))$ .

In addition to updating the function  $f_t$  by using some online learning algorithms, the weights  $w_{1,t}$  and  $w_{2,t}$  should also be updated. One suggested scheme is

$$w_{1,t+1} = C_t w_{1,t} \exp(-\eta \ell^*(h)), \quad w_{2,t+1} = C_t w_{2,t} \exp(-\eta \ell^*(f_t))$$

where  $C_t$  is a normalization term to keep  $w_{1,t+1} + w_{2,t+1} = 1$  and  $\ell^*(g) = (\Pi(g(\mathbf{x}_t)) - \Pi(y_t))^2$ .

**Heterogeneous Setting** Since heterogeneous OTL is generally very challenging, we consider one simpler case where the feature space of the source domain is a subset of that of the target domain. Without loss of generality, we assume the first  $d'$  dimensions of  $\mathbf{x}_t$  represent the old features, denoted as  $\mathbf{x}_t^{(1)} \in \mathbb{R}^{d'}$ . The other dimensions form a feature vector  $\mathbf{x}_t^{(2)} \in \mathbb{R}^{d-d'}$ . The key idea is to adopt a co-regularization principle of online learning two classifiers  $f_t^{(1)}$  and  $f_t^{(2)}$  simultaneously from the two views, and predict an unseen example on the target domain by

$$\hat{y}_t = \text{sgn}\left(\frac{1}{2}f_t^{(1)}(\mathbf{x}_t^{(1)}) + \frac{1}{2}f_t^{(2)}(\mathbf{x}_t^{(2)})\right)$$

The function from source domain  $h(\mathbf{x}^{(1)})$  is used to initialize  $f_t^{(1)}$ . The update strategy at time  $t$  is

$$(f_{t+1}^{(1)}, f_{t+1}^{(2)}) = \arg \min_{f^{(1)}, f^{(2)}} \frac{\gamma_1}{2} \|f^{(1)} - f_t^{(1)}\|_{\mathcal{H}}^2 + \frac{\gamma_2}{2} \|f^{(2)} - f_t^{(2)}\|_{\mathcal{H}}^2 + C \ell_t$$

where  $\gamma_1, \gamma_2$  and  $C$  are positive parameters and  $\ell_t$  is the hinge loss.

**Other Related Work** Multi-source Online Transfer Learning (MSOTL) (Ge et al., 2013; Tommasi et al., 2012) solves a more challenging problem where  $k$  classifiers  $h_1, \dots, h_k$  are provided by  $k$  sources. The goal is to learn the optimal combination of the  $k$  classifiers and the online updated classifier  $f_t$ . A naive solution is to construct a new  $d + k$  dimensional feature representation  $\mathbf{x}'_t = [\mathbf{x}_t, h_1(\mathbf{x}_t), \dots, h_k(\mathbf{x}_t)]$  and the online classifier in this new feature space. An extension of MSOTL (Ge et al., 2014) aims to deal with transfer learning problem under two disadvantageous assumptions, negative transfer where instead of improving performance, transfer learning from highly irrelevant sources degrades the performance on the target domain, and imbalanced distributions where examples in one class dominate. The Co-transfer Learning algorithm (Bhatt et al., 2014, 2012) considers the transfer learning problem not only in multi-source setting but also in the scenario where a large group of instances are unlabeled.

#### 4.6 Online Metric Learning

Distance metric learning (DML) (Yang and Jin, 2006) or similarity learning (Hao et al., 2017b) is an important problem in machine learning, which enjoys many real-world applications, such as image retrieval, classification and clustering. The goal of classic DML is to seek a distance matrix  $A \in \mathbb{R}^{d \times d}$  that defines the Mahalanobis distance between any two instances  $\mathbf{x}_i \in \mathbb{R}^d$  and  $\mathbf{x}_j \in \mathbb{R}^d$

$$d_A(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)^\top A (\mathbf{x}_i - \mathbf{x}_j) = \|W\mathbf{x}_i - W\mathbf{x}_j\|_2^2$$

Typically, matrix  $A \succeq 0$  is required to be symmetric positive semi-definite, i.e., there exist a matrix  $W \in \mathbb{R}^{d \times d}$  such that  $A = W^\top W$ . It is often hard to collect training data with the exact true values of distances. Therefore, there are two types of problem settings for online DML: 1) Pairwise data, where at each round  $t$  the learner receives a pair of instances  $(\mathbf{x}_t^1, \mathbf{x}_t^2)$  and a label  $y_t$  which is  $+1$  if the pair is similar and  $-1$  otherwise; 2) Triple data, where at each round  $t$  the learner receives a triple  $(\mathbf{x}_t, \mathbf{x}_t^+, \mathbf{x}_t^-)$ , with the feedback that  $d_A(\mathbf{x}_t, \mathbf{x}_t^+) > d_A(\mathbf{x}_t, \mathbf{x}_t^-)$ . The goal of online learning is to minimize the accumulated loss during the whole learning process  $\sum_{t=1}^T \ell_t(A)$ , where  $\ell_t$  is the loss suffered from imperfect prediction at round  $t$ . When evaluating the output model for online-to-batch-conversion, we may use the metric in information retrieval to evaluate the actual performance, such as mean average precision (mAP) or precision-at-top- $k$ .

Below, we briefly introduce a few representative work for DML in online setting.

**Pseudo-metric Online Learning (POLA)** The POLA algorithm (Shalev-Shwartz et al., 2004) learns the distance matrix  $A$  from a stream of pairwise data. The loss at time  $t$  is an adaptation of the hinge loss

$$\ell_t(A, b) = \max\{0, 1 - y_t(b - d_A(\mathbf{x}_t^1, \mathbf{x}_t^2))\}$$

where  $b$  is the adaptive threshold value for similarity and will be updated incrementally along with matrix  $A$ . We denote  $(A, b) \in \mathbb{R}^{d^2+1}$  as the new variable to learn. The update strategy mainly follows the PA approach

$$(A_{t+\frac{1}{2}}, b_{t+\frac{1}{2}}) = \arg \min_{(A, b)} \|(A, b) - (A_t, b_t)\|_2^2 \quad \text{s.t. } \ell_t(A, b) = 0$$

The solution  $(A_{t+\frac{1}{2}}, b_{t+\frac{1}{2}})$  ensures correct prediction to current pair while makes the minimal change to the previous model. Then, the algorithm projects this solution to the feasible space  $\{(A, b) : A \succeq 0, b \geq 1\}$  to obtain the updated model  $(A_{t+1}, b_{t+1})$ . Like PA, one can generalize POLA to soft-margin variants to be robust to noise.

Another similar work named Online Regularized Metric Learning (Jin et al., 2009) is simpler due to the adoption of fixed threshold, and adopts the following loss function

$$\ell_t(A) = \max(0, b - y_t(1 - d_A(\mathbf{x}_t^1, \mathbf{x}_t^2)))$$

whose gradient is

$$\nabla \ell_t(A) = y_t(\mathbf{x}_t^1 - \mathbf{x}_t^2)(\mathbf{x}_t^1 - \mathbf{x}_t^2)^\top$$

At time  $t$ , if the prediction is incorrect, the algorithm updates the matrix  $A$  by projecting the OGD updated matrix into the positive definite space.



**Information Theoretic Metric Learning (ITML).** In the above algorithms, distances between two matrices  $A_t$  and  $A$  are usually defined using the Frobenius norm, i.e.  $\|A_t - A\|_F^2$ . The ITML algorithm (Davis et al., 2007; Jain et al., 2009) adopts a different definition from an information theoretic perspective. Given a Mahalanobis distance parameterized by  $A$ , its corresponding multivariate Gaussian distribution is  $p(\mathbf{x}, A) = \frac{1}{Z} \exp(-\frac{1}{2}d_A(\mathbf{x}, \boldsymbol{\mu}))$ . The difference between matrices is defined as the KL divergence between two distributions. Assuming all distributions have the same mean, the KL divergence can be calculated as

$$\text{KL}(p(\mathbf{x}; A), p(\mathbf{x}, A_t)) = \text{tr}(AA_t^{-1}) - \log \det(AA_t^{-1}) - d$$

Similar to the PA update strategy, during time  $t$  the matrix is updated by optimizing

$$A_{t+1} = \arg \min_{A \succeq 0} \text{KL}(p(\mathbf{x}; A), p(\mathbf{x}, A_t)) + \eta \ell_t(A)$$

where  $\eta > 0$  is a regularization parameter. This optimization has a closed-form solution.

**Online Algorithm for Scalable Image Similarity Learning (OASIS)** The OASIS algorithm learns a bilinear similarity matrix  $W \in \mathbb{R}^d$  from a stream of triplet data, where the bilinear similarity measure between two instances is defined as

$$S_W(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top W \mathbf{x}_j$$

During time  $t$ , one triplet  $(\mathbf{x}_t, \mathbf{x}_t^+, \mathbf{x}_t^-)$  is received. Ideally, we expect the  $\mathbf{x}_t$  is more similar to  $\mathbf{x}_t^+$  than to  $\mathbf{x}_t^-$ , i.e.  $S_W(\mathbf{x}_t, \mathbf{x}_t^+) > S_W(\mathbf{x}_t, \mathbf{x}_t^-)$ . Similar to the PA algorithm, for a large margin, the loss function is defined as the hinge loss

$$\ell_t(W) = \max\{0, 1 - S_W(\mathbf{x}_t, \mathbf{x}_t^+) + S_W(\mathbf{x}_t, \mathbf{x}_t^-)\}$$

The optimization problem to solve for updating  $W_t$  is

$$W_{t+1} = \arg \min_W \frac{1}{2} \|W - W_t\|_F^2 + C\xi \text{ s.t. } \ell_t(W) \leq \xi \text{ and } \xi \geq 0$$

where  $C$  is the parameter controlling the trade-off. The OASIS algorithm differs from the previous work in several aspects. First, it does not require the similarity matrix  $W$  to be positive semi-definite and thus saves computational cost for the projection step. The bilinear similarity matrix may be better than the Mahalanobis distance for some applications. Third, the triplet data may be easier to collect in some applications.

Most of the previous methods assume a linear proximity function, (Xia et al., 2014) overcomes the limitaiton using a kernelized approach for metric learning. Another approach in (Gao et al., 2014, 2017) performs sparse online metric learning for very high dimensional data. There is also some work that solves the online similarity learning in an active learning setting (Hao et al., 2015b), which significantly reduces the cost of collecting labeled data. Some work in (Chen et al., 2015, 2016a) also applied techniques of the online similarity learning for real-world applications of mobile application recommendation and tagging. The series of work in (Xia et al., 2013; Wu et al., 2016) proposed the online multi-modal distance metric learning algorithms which learn distance metrics in multiple modalities, enabling multimedia information retrieval applications.

#### 4.7 Online Collaborative Filtering

Collaborative Filtering (CF) (Heckel and Ramchandran, 2017) is an important learning technique for recommender systems. Different from content-based filtering techniques, CF algorithms usually require minimal knowledge about the features of items or users apart from the previous preferences. The fundamental assumption of CF is that if two users rate many items similarly, they expect to share common preference on some other items. Several survey papers in (Shi et al., 2014; Su and Khoshgoftaar, 2009) gave detailed reviews of regular CF techniques. However, most of them assume batch settings. Below we introduce basics of CF and then review several popular online algorithms for CF tasks.

An online CF algorithm works on a sequence of observed ratings given by  $n$  users to  $m$  items. At time  $t \in \{1, 2, \dots, T\}$ , the algorithm receives the index of a user  $u^{(t)} \in \{1, 2, \dots, n\}$  and the index of an item  $i^{(t)} \in \{1, 2, \dots, m\}$  and makes a prediction of the rating  $\hat{r}_{u,i}^{(t)} \in \mathbb{R}$  based on the knowledge of the previous ratings. Then the real rating  $r_{u,i}^{(t)} \in \mathbb{R}$  is revealed and the algorithm updates the model based on the loss suffered from the imperfect prediction, denoted as  $\ell(\hat{r}_{u,i}^{(t)}, r_{u,i}^{(t)})$ . The goal of online CF is to minimize the Root Mean Square Error (RMSE) or Mean Absolute Error (MAE) along the whole learning process, defined as follows:

$$\text{RMSE} = \sqrt{\frac{1}{T} \sum_{t=1}^T (r_{u,i}^{(t)} - \hat{r}_{u,i}^{(t)})^2}, \quad \text{MAE} = \frac{1}{T} \sum_{t=1}^T |r_{u,i}^{(t)} - \hat{r}_{u,i}^{(t)}|$$

CF techniques are generally categorized into two types: memory-based methods and model-based methods. Below briefly introduces some popular algorithms in each category.

**Memory-Based CF Methods.** This follows the instance-based learning paradigm:

1. Calculate the similarity score between any pairs of items. For example, the cosine similarity between item  $i$  and item  $j$  is defined as,

$$S_{i,j} = \frac{\sum_{u \in \mathcal{U}_i \cap \mathcal{U}_j} r_{ui} \cdot r_{uj}}{\sqrt{\sum_{u \in \mathcal{U}_i} r_{ui}^2 \sum_{u \in \mathcal{U}_j} r_{uj}^2}}$$

where  $\mathcal{U}_i$  denotes the set of users that have rated item  $i$ .

2. For each item  $i$ , find its  $k$  nearest neighbor set  $\mathcal{N}_i$  based on the similarity score.
3. Predict the rating  $r_{u,i}$  as the weighted average of ratings from user  $u$  to the neighbors of item  $j$ , where the weight is proportional to the similarity.

We name the above described algorithm as item-based CF, while similarly, the predictions may also be calculated as the weighted average of ratings from similar users, which is called user-based CF method. Memory-based CF methods were used for some early generation recommendation systems, but very few is online learning approach. One reason is because of data sparsity, as the similarity score  $S_{i,j}$  is only available when there is at least one common user that rates the two items  $i$  and  $j$ , which might be unrealistic during the beginning stage. Another challenge is the large time consumption when updating the large number of similarity scores incrementally with the arrival of new ratings. The Online Evolutionary Collaborative Filtering (Liu et al., 2010) algorithm provides an efficient similarity score updating method to address this problem.

**Model-Based CF Methods** Memory-based online CF methods suffer two limitations, i.e., sensitivity due to data sparsity and inefficiency for similarity score update. To address these issues, extensive work has been focused on model-based CF algorithms. One of the most successful approaches is the matrix factorization methodology (Koren et al., 2009), which assumes the rating by a user to an item is determined by  $k$  potential features,  $k \ll n, m$ . Thus each user  $u$  can be represented by a vector  $\mathbf{u}_u \in \mathbb{R}^k$ , and each item  $i$  can be represented by a vector  $\mathbf{v}_i \in \mathbb{R}^k$ . The rating  $r_{u,i}$  can then be approximated by the dot product of the corresponding user vector and item vector, i.e.,  $\hat{r}_{u,i} = \mathbf{u}_u^\top \mathbf{v}_i$ . The CF problem can then be represented by the following optimization problem:

$$\arg \min_{U \in \mathbb{R}^{k \times n}, V \in \mathbb{R}^{k \times m}} \sum_{t=1}^T \ell(\mathbf{u}_u^{(t)} \cdot \mathbf{v}_i^{(t)}, r_{u,i}^{(t)})$$

where the loss function is defined to optimize certain evaluation metric:

$$\ell_{rmse}(\hat{r}_{u,i}, r_{u,i}) = (r_{u,i} - \hat{r}_{u,i})^2 \quad \ell_{mae}(\hat{r}_{u,i}, r_{u,i}) = |r_{u,i} - \hat{r}_{u,i}|$$

The regularized loss at time  $t$  is

$$\mathcal{L}_t = \lambda \|\mathbf{u}_u^{(t)}\|_2^2 + \lambda \|\mathbf{v}_i^{(t)}\|_2^2 + \ell(\mathbf{u}_u^{(t)} \cdot \mathbf{v}_i^{(t)}, r_{u,i}^{(t)})$$

where  $\lambda > 0$  is the regularization parameter. A straightforward CF approach is to apply OGD on the regularized loss function (Abernethy et al., 2007),

$$\mathbf{u}_u^{(t+1)} = (1 - 2\eta\lambda)\mathbf{u}_u^{(t)} - \eta \frac{\partial \ell(\mathbf{u}_u^{(t)} \cdot \mathbf{v}_i^{(t)}, r_{u,i}^{(t)})}{\partial \mathbf{u}_u^{(t)}}, \quad \mathbf{v}_i^{(t+1)} = (1 - 2\eta\lambda)\mathbf{v}_i^{(t)} - \eta \frac{\partial \ell(\mathbf{u}_u^{(t)} \cdot \mathbf{v}_i^{(t)}, r_{u,i}^{(t)})}{\partial \mathbf{v}_i^{(t)}}$$

where  $\eta > 0$  is the learning rate. Later, several improved algorithms are proposed, such as Online Multi-Task Collaborative Filtering algorithm (Wang et al., 2013a), Dual-Averaging Online Probabilistic Matrix Factorization algorithm (Yuan-Xiang et al., 2014), Adaptive Gradient Online Probabilistic Matrix Factorization algorithm (Yuan-Xiang et al., 2014) and Second-order Online Collaborative Filtering algorithm (Lu et al., 2013; Liu et al., 2016b). These algorithms adopt more advanced update strategies beyond OGD and thus can achieve faster adaptation for rapid user preference changes in real-world recommendation tasks.

Besides the algorithms introduced, there are many online CF methods that explore other challenging tasks. First, in many applications, both features of users and items are available and thus need to be considered for better prediction. This generalized CF problem can be solved by using tensor product kernel functions. For example, the Online Low-rank with Features algorithm (Abernethy et al., 2007) addresses this problem in online setting. However, it only adopts the linear kernel for efficiency. Perhaps, better performance might be achieved if online budget learning algorithms are adopted. Second, most CF algorithms are based on a regression model, which is mainly concerned with the accuracy of rating prediction, while there are some applications where ranking prediction might be much more important. Two algorithms based on OGD and Dual Averaging approaches are proposed to address this problem by replacing the regression-based loss with the ranking-based loss (Ling et al., 2012). Third, for very large-scale applications, when the model has to be

learnt using parallel computing, conventional OGD update is not suitable because of the possible conflict in updating the user/item vectors. The Streaming Distributed Stochastic Gradient Descent algorithm (Ali et al., 2011) provides an operable approach to addresses this problem. Finally, the CF methods for Google News recommendation (Das et al., 2007) is a combination of memory-based and model-based algorithms. Last but not least, to address the sparsity problem and imbalance of rating data, (Liu et al., 2017) incorporate content information via latent dirichlet allocation into online CF.

#### 4.8 Online Learning to Rank

Learning to rank is an important family of machine learning techniques for information retrieval (Trotman, 2005; Cao et al., 2007; Hang, 2011; Zoghi et al., 2017; Shah et al., 2017). Different from classification problems where instances are classified as either “relevant” or “not relevant”, learning to rank aims to produce a permutation of a group of unseen instances which is similar to the knowledge acquired from the previously seen rankings. To evaluate the performance of ranking algorithms, metrics for information retrieval such as Mean Average Precision (MAP), Normalized Discounted Cumulative Gain (NDCG) and Precision-At-Top- $k$  are most popular.

Unlike traditional learning to rank methods which are often based on batch learning (Hang, 2011), we mainly focus on reviewing existing learning to rank methods in online settings (Wang et al., 2015b; Wan et al., 2015), where instances are observed sequentially. Learning to rank techniques are generally categorized into two approaches: pointwise and pairwise. We will introduce some of the most representative algorithms in each category.

**Pointwise Approach:** We first introduce a simple Perceptron-based algorithm, the Prank (Crammer et al., 2001; Crammer and Singer, 2005), which provides a straightforward view of the commonly used problem setting for pointwise learning to rank approaches.

To define the online learning to rank problem setting formally, we have a finite set of ranks  $\mathcal{Y} = \{1, \dots, k\}$  from which a rank  $y \in \mathcal{Y}$  is assigned to an instance  $\mathbf{x} \in \mathbb{R}^d$ . During time  $t$ , an instance  $\mathbf{x}_t$  is received and the algorithm makes a prediction  $\hat{y}_t$  based on the current model  $H_t : \mathbb{R}^d \rightarrow \mathcal{Y}$ . Then the true rank  $y_t$  is revealed and the model is updated based on the loss  $\ell(\hat{y}_t, y_t)$ . The loss, for instance, can be defined as  $\ell(\hat{y}_t, y_t) = |\hat{y}_t - y_t|$ . The goal of the online learning to rank task is to minimize the accumulated loss along the whole learning process  $\sum_{t=1}^T \ell(\hat{y}_t, y_t)$ .

The ranking rule of Prank algorithm consists of the combination of Perceptron weight  $\mathbf{w} \in \mathbb{R}^d$  and a threshold vector  $\mathbf{c} \in \{\mathbb{R}, \infty\}^d$ , whose elements are in nondecreasing order i.e.,  $c^1 \leq c^2 \leq \dots \leq c^k = \infty$ . Like the Perceptron algorithm, the rank prediction is determined by the value of the inner product  $\mathbf{w}_t^\top \mathbf{x}_t$ ,

$$\hat{y}_t = \min_{r \in \{1, \dots, k\}} \{r : \mathbf{w}_t^\top \mathbf{x}_t < c_t^r\}$$

We can expand the target rank  $y_t$  to a vector  $\mathbf{y}_t = \{+1, \dots, +1, -1, \dots, -1\} \in \mathbb{R}^k$ . For  $r = 1, \dots, k$ ,  $y_t^r = -1$  if  $y_t < r$ , and  $y_t^r = 1$  otherwise. Thus, for a correct prediction,  $\mathbf{y}_t^\top (\mathbf{w}_t^\top \mathbf{x}_t - \mathbf{c}_t^r) > 0$  holds for all  $r \in \mathcal{Y}$ . When a mistake appears  $\hat{y}_t \neq y_t$ , there is subset  $\mathcal{M}$  of  $\mathcal{Y}$  where  $\mathbf{y}_t^\top (\mathbf{w}_t^\top \mathbf{x}_t - \mathbf{c}_t^r) > 0$  does not hold. The update rule is to move the corresponding

thresholds for ranks in  $\mathcal{M}$  and the weight vector toward each other:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \left( \sum_{r \in \mathcal{M}} y_t^r \right) \mathbf{x}_t, \quad \text{and} \quad c_{t+1}^r = c_t^r - y_t^r, \forall r \in \mathcal{M}$$

In theory, the elements in threshold vector  $\mathbf{c}$  are always in nondecreasing order and the total number of mistakes made during the learning process is bounded.

Online Aggregate Prank-Bayes Point Machine (OAP-BPM) (Harrington, 2003) is an extension of the Prank algorithm by approximating the Bayes point. Specifically, the OAP-BPM algorithm generates  $N$  diverse solutions of  $\mathbf{w}$  and  $\mathbf{c}$  during each iteration and combines them for a better final solution. We denote  $H_{j,t}$  as the  $j$ -th solution at time  $t$ . The algorithm samples  $N$  Bernoulli variables  $b_{j,t} \in \{0, 1\}, j = \{1, \dots, N\}$  independently. If  $b_{j,t} = 1$ , The  $j$ -th solution is updated using the Prank algorithm according to the current instance,  $H_{j,t+1} = \text{Prank}(H_{j,t}, (\mathbf{x}_t, y_t))$ . Otherwise, no update is conducted to the  $j$ -th solution. The solution  $\mathbf{w}_{t+1}$  and  $\mathbf{c}_{t+1}$  is the average over  $N$  solutions. This work shows better generalization performance than the basic Prank algorithm.

**Pairwise Approach:** One simple method is to address the ranking problem by transforming it to a classification problem (Herbrich et al., 1999). In a more challenging problem setting, where no accurate rank  $y$  is available when collecting the data, only pairwise instances are provided. At time  $t$ , a pair of instances  $(\mathbf{x}_t^1, \mathbf{x}_t^2)$  are received with the knowledge that  $\mathbf{x}_t^1$  is ranked before  $\mathbf{x}_t^2$  or the inverse case, and the aim is to find a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  that fits the instance pairs, i.e.,  $f(\mathbf{x}^1) > f(\mathbf{x}^2)$  when it is known  $\mathbf{x}_t^1$  is ranked before  $\mathbf{x}_t^2$  or otherwise  $f(\mathbf{x}^1) < f(\mathbf{x}^2)$ . When the function is linear, the problem can be rewritten as  $\mathbf{w}^\top (\mathbf{x}^1 - \mathbf{x}^2) > 0$  when  $\mathbf{x}^1$  is in front and otherwise  $\mathbf{w}^\top (\mathbf{x}^1 - \mathbf{x}^2) < 0$ , where  $\mathbf{w} \in \mathbb{R}^d$  is the weight vector. This problem can easily be solved by using a variety of online classification algorithms (Online Gradient Descent (Chapelle and Keerthi, 2010) for example).

#### 4.9 Distributed Online Learning

Distributed online learning (Zhang et al., 2017) has become increasing popular due to the explosion in size and complexity of datasets. Similar to the mini-batch online learning, during each iteration,  $K$  instances are received and processed simultaneously. Usually, each node processes one of the instances and updates its local model. These nodes communicate with each other to make their local model consistent. When designing a distributed algorithm, besides computational time cost and accuracy, another important issue to consider is the communication load between nodes. This is because in real world systems with limited network capacity and large communication burden result in long latency.

Based on the network structure, distributed online learning algorithms can be classified into two groups: *centralized* and *decentralized* algorithms. A centralized network is made up of 1 master node and  $K - 1$  worker nodes, where the workers can only communicate with the master node. By gathering and distributing information across the network, it is not difficult for distributed algorithms to reach a global consensus (Boyd et al., 2011). In decentralized networks, however, there is no master and each node can only communicate with its neighbors (Agarwal et al., 2010; Mota et al., 2013). Although the algorithms are more complex, decentralized learning is more popular because of the robustness of network structure.

We can also group the distributed learning algorithms by *synchronized* and *asynchronous* working modes. Synchronized algorithms are easy to design and enjoy better theoretical bounds but the speed of the whole network is limited by the slowest node. Asynchronous learning algorithms, on the other hand, are complex and usually have worse theoretical bounds. The advantage is its faster processing speed (Smyth et al., 2009).

#### 4.10 Online Learning with Neural Networks

In addition to kernel-based online learning approaches, another rich family of nonlinear online learning algorithms follows the general idea of neural network based learning approaches (Williams and Zipser, 1989; Platt, 1991; Carpenter et al., 1991; LeCun et al., 1998; Liang et al., 2006; Polikar et al., 2001). For example, the Perceptron algorithm could be viewed as the simplest form of online learning with neural networks (but it is not nonlinear due to its trivial network). Despite many extensive studies for online learning (or incremental learning) with neural networks, many of existing studies in this field fall short due to some critical drawbacks, including the lack of theoretical analysis for performance guarantee, heuristic algorithms without solid justification, and computational too expensive to achieve efficient and scalable online learning. Due to the large body of related work, it is impossible to examine every piece of work in this area. In the following, we review several of the most popularly cited related papers and discuss their key ideas for online learning with neural networks.

A series of related work has explored online convex optimization methods for training classical neural network models (Bottou, 1998a), such as the Multi-layer Perceptron (MLP). For example, online/stochastic gradient descent has been extensively studied for training neural networks in sequential/online learning settings, such as the efficient back-propagation algorithm using SGD (LeCun et al., 1998). These works are mainly motivated to accelerate the training of batch learning tasks instead of solving online learning tasks directly and seldom give theoretical analysis.

In addition to the above, we also briefly review other highly cited works that address online/incremental learning with neural networks. For example, the study in (Williams and Zipser, 1989) presented a novel learning algorithm for training fully recurrent neural networks for temporal supervised learning which can continually run over time. However, the work is limited in lacking theoretical analysis and performance guarantee, and the solution could be quite computationally expensive. The work in (Platt, 1991) presented a Resource-Allocating Network (RAN) that learns a two-layer network by a strategy for allocating new units whenever an unusual pattern occurs and a learning rule for refining the network using gradient descent. Although the algorithm was claimed to run in online learning settings, it may suffer poor scalability as the model complexity would grow over time. The study in (Carpenter et al., 1991) proposed a new neural network architecture called “ARTMAP” that autonomously learns to classify arbitrarily many, arbitrarily ordered vectors into recognition categories based on predictive success. This supervised learning system was built from a pair of ART modules that are capable of self organizing stable recognition categories in response to arbitrary sequences of input patterns. Although an online learning simulation has been done with ART (adaptive resonance theory), the solution is not optimized directly for online learning tasks and there is also no theoretical analysis. The work in (Liang et al.,

2006) proposed an online sequential extreme learning machine (OS-ELM) which explores an online/sequential learning algorithm for training single hidden layer feedforward networks (SLFNs) with additive or radial basis function (RBF) hidden nodes in a unified framework. The limitation also falls short in some heuristic approaches and lacking theoretical analysis. Last but not least, there are also quite many studies in the field which claim that they design neural network solutions to work online, but essentially they are not truly online learning. They just adapt some batch learning algorithms to work efficiently for sequential learning environments, such as the series of Learning++ algorithms and their variants (Polikar et al., 2001; Elwell and Polikar, 2011). Recently, Hedge Backpropagation (Sahoo et al., 2018) was proposed to learn deep neural networks in the online setting with the aim to address slow convergence of deep networks through dynamic depth adaptation.

#### 4.11 Online Portfolio Selection

On-line Portfolio Selection (OLPS) is a natural application of online learning for sequential decisions of selecting a portfolio of stocks for optimizing certain metrics, e.g. cumulative wealth, risk adjusted returns, etc. (Cesa-Bianchi and Lugosi, 2006; Li and Hoi, 2014, 2015; Li et al., 2016a). Consider a financial market with  $m$  assets, in which we have to allocate our wealth. At every time period (or iteration), the price of the  $m$  stocks changes by a factor of  $\mathbf{x}_t \in \mathbb{R}_+^m$ . This vector is also called the price relative vector.  $x_{t,i}$  denotes the ratio of the closing price of asset  $i$  at time  $t$  to the last closing price at time  $t-1$ . Thus, an investment in asset  $i$  changes by a factor of  $x_{t,i}$  in period  $t$ . At the beginning of time period  $t$  the investment is specified by a portfolio vector  $\mathbf{b}_t \in \delta_m$  where  $\delta_m = \{\mathbf{b} : \mathbf{b} \succeq 0, \mathbf{b}^\top \mathbf{1} = 1\}$ . The portfolio is updated in every time-period based on a specific strategy, and produces a sequence of mappings:

$$\mathbf{b}_1 = \frac{1}{m}, \quad \mathbf{b}_t : \mathbb{R}_+^{m(t-1)} \rightarrow \delta_m, \quad t = 2, 3, \dots, T$$

where  $T$  is the maximum length of the investment horizon. To make a decision for constructing a portfolio at time  $t$ , the entire historical information from  $\mathbf{x}_1, \dots, \mathbf{x}_{t-1}$  is available. The theoretical framework starts with a wealth of  $S_0 = 1$ , and at the end of every time period, the wealth changes as  $S_t = S_{t-1} \times (\mathbf{b}_t^\top \mathbf{x}_t)$ .

Most efforts in OLPS make a few (possibly unrealistic) assumptions, including no transaction costs, perfectly liquid market, and no impact cost (the portfolio selection strategy does not affect the market). Besides the traditional benchmarking approaches, the approaches for OLPS can be categorized into *Follow-the-winner*, *Follow-the-loser*, *Pattern Matching*, and *Meta-Learning* approaches (Li and Hoi, 2014).

The benchmark approaches, as the name suggests, are simple baseline methods whose performance can be used to benchmark the performance of proposed algorithms. Common baselines are Buy and Hold (BAH) strategy, Best Stock and Constant Rebalanced Portfolio (CRP). The idea of BAH is to start with a portfolio with equal investment in each asset, and never rebalance it. Best Stock is the performance of the asset with the highest returns at the end of the investment horizon. CRP (Kelly Jr, 2011) is a fixed portfolio allocation to which the portfolio is rebalanced to at the end of every period, and Best-CRP is the CRP which obtains the highest returns at the end of the investment horizon. It should be noted that Best Stock and Best-CRP strategies can only be executed in hindsight.

Follow-the-winner approaches adhere to the principle of increasing the relative portfolio allocation weight of the stocks that have performed well in the past. Many of the approaches are directly inspired by Convex Optimization theory, including Universal Portfolios (Cover, 2011), Exponential Gradient (Helmbold et al., 1998), Follow the Leader (Gaivoronski and Stella, 2000) and Follow the Regularized Leader (Agarwal et al., 2006). In contrast to follow-the-winner, there is a set of approaches that aim to follow-the-loser, with the belief that asset prices have a tendency to revert back to a mean, i.e., if the asset price falls, it is likely to rise up in the next time-period. These are also called mean-reversion strategies. The early efforts in this category included Anti Correlation (Borodin et al., 2004) which designed a strategy by betting making statistical bets on positive-lagged correlation and negative auto-correlation; and Passive-Aggressive Mean Reversion (PAMR) (Li et al., 2012), which extended the Online Passive Aggressive Algorithms (Crammer et al., 2006) to update the portfolio to an optimal "loser" portfolio - by selecting a portfolio that would have made an optimal loss in the last observed time-period. A similar idea was used to extend confidence-weighted online learning to develop Confidence-Weighted Mean Reversion (Li et al., 2011b, 2013). The idea of PAMR was extended to consider multi-period asset returns, which led to the development of Online Moving Average Reversion (OLMAR) (Li and Hoi, 2012; Li et al., 2015) and Robust Median Reversion (RMR) (Huang et al., 2013) strategies.

Another popular set of approaches is the Pattern-Matching approaches, which aim to find patterns (they may be able to exploit both follow-the-winner and follow-the-loser) for optimal sequential decision making. Most of these approaches are non-parametric. Exemplar approaches include (Gyorfi and Schafer, 2003; Györfi et al., 2008; Li et al., 2011a). Finally, meta-learning algorithms for portfolio selection aim to rebalance the portfolio on the basis of the expert advice. There are a set of experts that output a portfolio vector, and the meta-learner uses this information to obtain the optimal portfolio. In general, the meta-learner adheres to the follow-the-winner principle to identify the best performing experts. Popular approaches in this category include Aggregating Algorithms (Vovk and Watkins, 1998), Fast Universalization Algorithm (Akcoglu et al., 2004) and Follow the Leading History (Hazan and Seshadhri, 2009). Besides these approaches, there are also efforts in portfolio selection with aims to optimize the returns accounting for transaction costs. The idea is to incorporate the given transaction cost into the optimization objective (Ormos and Urbán, 2013; Huang et al., 2015; Li et al., 2017a).

A closely related area to Online Portfolio Selection is Online Learning for Time Series Prediction. Time series analysis and prediction (George, 1994; Clements and Hendry, 1998; Chatfield, 2000; Sapankevych and Sankar, 2009) is a classical problem in machine learning, statistics, and data mining. The typical problem setting of time series prediction is as follows: a learner receives a temporal sequence of observations,  $x_1, \dots, x_t$ , and the goal of the learner is to predict the future observations (e.g.,  $x_{t+1}$  or onwards) as accurately as possible. In general, machine learning methods for time series prediction may also be divided into linear and non-linear, univariate and multivariate, and batch and online. Some time series prediction tasks may be resolved by adapting an existing batch learning algorithm using sliding window strategies. Recently there have been some emerging studies for exploring online learning algorithms for time series prediction (Anava et al., 2013; Liu et al., 2016a).



## 5. Bandit Online Learning

### 5.1 Overview

Bandit online learning, a.k.a. the “Multi-armed Bandit (MAB) problem (Robbins, 1985; Katehakis and Veinott Jr, 1987; Vermorel and Mohri, 2005; Bubeck and Cesa-Bianchi, 2012; Gittins et al., 2011), is an important branch of online learning where a learner makes sequential decisions by receiving only partial feedback from the environment each time.

MAB problems are online learning tasks for sequential decisions with a trade-off between exploration and exploitation. Specifically, on each round, a player chooses one out of  $K$  actions, the environment then reveals the payoff of the player’s action, and the goal of the learner is to maximize the total payoff obtained during online learning process. A fundamental challenge of MAB is to address the exploration-exploitation tradeoff (Audibert et al., 2009), i.e., balancing between the *exploitation* of actions that gave highest payoffs in the past and the *exploration* of new actions that might give higher payoffs in the future.

MAB problems are collectively called “bandit” problems. Historically, the name “bandit” is referred to the scenario of playing slot machines in a casino, where a player faces a number of slot machines to insert coins (one slot machine is called one-armed bandit in American slang), the expected reward of each machine might be different, and the player’s goal is to maximize the reward by repeatedly choosing where to insert the next coin.

MAB problems can be roughly divided into two major categories: *stochastic MAB* and *adversarial MAB*. The former assumes a stochastic environment where rewards (or “losses” equivalently) are i.i.d. and independent from each other, while the later removes the stochastic assumption where the rewards (or losses) can be arbitrary or more formally “adversarial” by adapting to the past decisions. Note that the notion of “reward” and “loss” are symmetric and can be translated from one to the other equivalently. To be consistent, if not mentioned specifically, we will use “loss” instead of “reward” for the rest discussion.

We now introduce the formal procedure of the MAB problem. Formally, a  $K$ -armed bandit problem takes place in a sequence of rounds with length  $T \in \mathcal{N}$ , where  $T$  is often unknown at the beginning (typically a learner is called an “anytime” algorithm if  $T$  is unknown in advance). At the  $t$ -th round, the player chooses one out of  $K$  actions  $I_t \in [K] = \{1, \dots, K\}$  using some strategy. After that, the environment reveals the loss  $\ell_t(I_t)$  of the action to the forecaster. The goal is to minimize the total loss over the  $T$  rounds. In theory, we are interested in analyzing the behavior of the learner, typically by comparing the performance of its actions against with some optimal strategy. More formally, we can define the “regret” as the difference between the cumulative loss of the best fixed arm by an optimal strategy and that of the player after playing  $T$  rounds  $I_1, \dots, I_T$  as

$$R_T = \sum_{t=1}^T \ell_t(I_t) - \min_{i \in [K]} \sum_{t=1}^T \ell_t(i)$$

As both loss  $\ell_t(i)$  and action  $I_t$  could be stochastic, we can define the expected regret as

$$\mathbb{E}[R_T] = \mathbb{E} \left[ \sum_{t=1}^T \ell_t(I_t) - \min_{i \in [K]} \sum_{t=1}^T \ell_t(i) \right]$$

where the expectation is taken with respect to the random draw of both losses and learner’s actions.

## 5.2 Stochastic Bandits

For stochastic MAB problem, each arm  $i \in [k]$  corresponds to an unknown distribution  $P_i$  on  $[0, 1]$ , and the losses  $\ell_t(i)$  are independently drawn from the distribution  $P_i$  corresponding to the selected arm. Let us denote by  $\mu_i$  the mean of the distribution  $P_i$ , and define

$$\mu^* = \min_{i \in [k]} \mu_i \quad \text{and} \quad i^* \in \arg \min_{i \in [k]} \mu_i$$

The expected regret can be rewritten as

$$\mathbb{E}[R_T] = \mathbb{E}\left[\sum_{t=1}^T \mu_{I_t}\right] - T \min_{i \in [K]} \mu_i = \mathbb{E}\left[\sum_{t=1}^T \mu_{I_t}\right] - T\mu^* = \mathbb{E}\left[\sum_{t=1}^T (\mu_{I_t} - \mu^*)\right]$$

Further let  $N_i(s) = \sum_{t=1}^s \mathbb{I}(I_t = i)$  denote the number of times the player selected arm  $i$  on the first  $s$  rounds and  $\Delta_i = \mu_i - \mu^*$  be the suboptimality parameter of arm  $i$ , we can simplify the expected regret as follows:

$$\mathbb{E}[R_T] = \mathbb{E}\left[\sum_{t=1}^T \Delta_{I_t}\right] = \sum_{i=1}^K \Delta_i \mathbb{E}[N_i(T)].$$

### 5.2.1 STOCHASTIC MULTI-ARMED BANDIT

In this section we mainly introduce two well-known algorithms for stochastic MAB.

**$\epsilon$ -Greedy.** The first simplest algorithm for stochastic MAB is called the  $\epsilon$ -Greedy rule (Sutton and Barto, 1998). The idea is to with probability  $1 - \epsilon$  play the the current best arm of the highest average reward, and with probability  $\epsilon$  play a random arm, where parameter  $\epsilon > 0$  is a constant value in  $(0, 1)$ . Algorithm 13 gives a summary of this algorithm. However,

---

#### Algorithm 13: $\epsilon$ -Greedy

---

**INPUT:** parameter  $\epsilon > 0$   
**INIT:** empirical means  $\mu_i = 0, \forall i \in [K]$   
**for**  $t = 1, 2, \dots, T$  **do**  
    with probability  $1 - \epsilon$  play the current best arm  $i_t = \arg \min_{i \in [K]} \mu_i$   
    with probability  $\epsilon$  play a random arm  
    receive  $\ell_t(i_t)$   
    update the empirical means  $\mu_{i_t} = (\mu_{i_t} + \ell_t(i_t)) / (N_{i_t}(t) + 1)$   
**end for**

---

the constant exploration probability  $\epsilon$  results in a linear growth in the regret. One way to fix it is to decrease the value of  $\epsilon$  over time and let it go to zero at a certain rate. For example, an improved  $\epsilon_t$ -greedy algorithm is to follow the epsilon-decreasing strategy by defining  $\epsilon_t$  at round  $t$  as

$$\epsilon_t = \min \left\{ 1, \frac{cK}{d^2 t} \right\}$$

where  $c > 0$  and  $d \in (0, 1)$ . When  $0 < d \leq \min_{i: \mu_i < \mu^*} \Delta_i < 1$  and  $T > \frac{cK}{d}$ , this improved  $\epsilon_t$ -greedy algorithm can achieve the logarithm regret  $O(\ln T)$ .

**UCB.** Another well-known algorithm for stochastic MAB is the Upper Confidence Bound (UCB) algorithm (Auer et al., 2002a), a strategy that simultaneously performs exploration and exploitation using a heuristic principle of *optimism in face of uncertainty*. The intuition is that, despite lacking knowledge about what actions are best, we will try to construct an optimistic guess as to how good the expected payoff/loss of each action is, and choose the action with the best guess. If our guess is correct, we will be able to exploit that action and incur little regret; but if our guess is wrong, then our optimistic guess will quickly decrease and we will then be compelled to switch to a different action, which therefore is able to balance the exploration-exploitation tradeoff.

Formally, the “optimism” comes in the form of Upper Confidence Bound (UCB). In particular, the idea is to calculate the confidence intervals of the averages, which is a region around our estimates such that the true value falls within with high probability, and repeatedly shrink the confidence bounds such that the average will become more reliable. Algorithm 14 gives a summary of the UCB algorithm.

---

**Algorithm 14:** UCB

---

**INPUT:** parameter  $\epsilon > 0$   
**INIT:** empirical means  $\mu_i = 0, \forall i \in [K]$   
**for**  $t = 1, 2, \dots, T$  **do**  
     play the arm  $i_t = \arg \min_i (\mu_i - \sqrt{\frac{2 \ln t}{N_i(t)}})$   
     receive  $\ell_t(i_t)$   
     update the empirical means  $\mu_{i_t} = (\mu_{i_t} + \ell_t(i_t)) / (N_{i_t}(t) + 1)$   
**end for**

---

In theory, by running the UCB algorithm over  $T$  rounds, the expected regret is

$$\mathbb{E}[R_T] \leq 8 \sum_{i: \mu_i < \mu^*} \left( \frac{\ln T}{\Delta_i} \right) + \left( 1 + \frac{\pi^2}{3} \right) \left( \sum_{j=1}^K \Delta_j \right)$$

The above is a specific worst-case bound on the expected regret (Auer et al., 2002a). More concisely, one can show that the expected regret of UCB is at most  $O(\sqrt{KT \ln T})$ . Auer et al. (2002a) also gave some variants of improved UCB algorithms. Improved algorithms and regret bound were also given in (Abbasi-Yadkori et al., 2011).

### 5.2.2 BAYESIAN BANDITS

Bayesian methods have been explored for studying bandit problems from the beginning of this field. One of the most well-known and classic algorithm in Bayesian Bandits is Thompson Sampling (Thompson, 1933), which is considered one of the oldest algorithm to address the exploration-exploitation trade-off for bandit problems. Recent years have seen a lot of interests in analyzing both empirical performance (Chapelle and Li, 2011) and theoretical properties of Thompson Sampling for bandit problems (Agrawal and Goyal, 2012). In the following, we introduce a Bayesian setting of Bernoulli bandit and then discusses the Thompson Sampling algorithm.

Consider a standard  $K$ -armed Bernoulli bandit, each action corresponds to the choice of an arm, and the reward of the  $i$ -th arm is either 0 or 1 which follows a Bernoulli distribution

with mean  $\mu_i$ , i.e., the probability of success for arm  $i$  (reward=1) is  $\mu_i$ . The algorithm maintains Bayesian priors on the Bernoulli means  $\mu_i$ 's. At the beginning, the Thompson Sampling algorithm initializes each arm  $i$  to have prior  $Beta(1, 1)$  on  $\mu_i$ , since  $Beta(1, 1)$  is the uniform distribution on  $(0, 1)$ . On round  $t$ , after observing  $S_i(t)$  successes (reward=1) and  $F_i(t)$  failures (reward=0) on  $k_i(t) = S_i(t) + F_i(t)$  times of playing arm  $i$ , the algorithm updates the distribution on  $\theta_i$  as  $Beta(S_i(t) + 1, F_i(t) + 1)$ . The algorithm then samples from these posterior distributions of the  $\theta_i$ 's, and plays an arm according to the probability of its mean being the largest. Algorithm 15 gives a summary of Thompson Sampling algorithm for  $K$ -armed Bernoulli bandits problems.

---

**Algorithm 15:** Thompson Sampling
 

---

**INIT:**  $S_i(1) = 0, F_i(1) = 0, \forall i = 1, \dots, K$   
**for**  $t = 1, 2, \dots, T$  **do**  
     For each arm  $i \in [K]$ , sample  $\theta_i(t)$  from the  $Beta(S_i + 1, F_i + 1)$  distribution.  
     Play arm  $i(t) = \arg \max_i \theta_i(t)$ , and observe reward  $r_t$   
     If  $r_t = 1$ , then  $S_{i(t)} = S_{i(t)} + 1$ ; else  $F_{i(t)} = F_{i(t)} + 1$ .  
**end for**

---

The above Thompson sampling algorithm can be easily extended to the general stochastic bandits setting, i.e., when the rewards for arm  $i$  are generated from an arbitrary unknown distribution with support  $[0, 1]$  and mean  $\mu_i$ . It has also been extensively used for contextual bandit settings (Chapelle and Li, 2011; Agrawal and Goyal, 2013).

In theory, for the  $K$ -armed stochastic bandit problem, denoting  $\Delta_i = \mu_1 - \mu_i$ , the Thompson Sampling algorithm has the expected regret given in (Agrawal and Goyal, 2012)

$$\mathbb{E}[R_T] \leq O\left(\left(\sum_{a=2}^K \frac{1}{\Delta_a^2}\right)^2 \ln T\right)$$

*Other Related Work.* In addition to the classic Thompson Sampling algorithm, another notable variant of Bayesian Bandit is called the Bayes-UCB algorithm (Kaufmann et al., 2012), which is a UCB-like algorithm, where the upper confidence bounds are based on the quantiles of Beta posterior distributions, and is able to achieve the lower bound of Lai and Robbins (1985) for Bernoulli rewards. More other extensive and recent studies of Bayesian Bandits can be found in (Scott, 2010; May et al., 2012; Russo and Van Roy, 2016).

### 5.3 Adversarial Bandits

In the previous setting of stochastic bandits, we generally assume that the rewards are i.i.d., which are drawn independently from some unknown but fixed distribution. We now relax such stochastic assumption on the rewards. We assume the reward distribution can be affected by the previous actions taken by the player, which is termed as ‘‘Adversarial Bandits’’ problems Auer et al. (1995). In the following, we review several classes of adversarial bandits, including some fundamentals of adversarial MAB and other active topics such as linear bandits and combinatorial bandits.

### 5.3.1 ADVERSARIAL MULTI-ARMED BANDIT

We consider a  $K$ -armed adversarial bandit problem where  $K > 1$  and the learner receives an arbitrary sequence of loss vectors  $(\ell_1, \dots, \ell_T)$  where  $\ell_t \in [0, 1]^K \forall t \in [T]$ . On each round, the learner plays an action  $I_t \in [K]$  and observes the loss  $\ell_t(I_t)$ . For adversarial bandits, a randomized policy is commonly used. In particular, given some policy  $\pi$ , the conditional distribution over the actions having observed  $\Omega_{t-1} = \{(I_1, \ell_1), \dots, (I_{t-1}, \ell_{t-1})\}$  is  $P_t = \pi(\cdot | \Omega_{t-1}) \in \mathcal{P}_{K-1}$ . The performance of a policy  $\pi$  on the environment can be measured by the expected regret which is the expected loss of the policy relative to the best fixed action in hindsight:

$$\mathbb{E}[R_T] = \mathbb{E} \left[ \sum_{t=1}^T \ell_t(I_t) \right] - \min_{i \in [K]} \sum_{t=1}^T \ell_t(i)$$

**Exp3.** The Exponential-weights for Exploration and Exploitation algorithm (Exp3) (Auer et al., 2002b) is a popular algorithm for adversarial MAB. It follows the similar idea of prediction with expert advice and applies the Hedge (or Weighted-Majority) algorithm to the tradeoff of exploration and exploitation. Specifically, we first define a probability vector  $\mathbf{p}_t \in \mathbb{R}^K$  in which the  $i$ -th element  $p_t(i)$  indicates the probability of drawing arm  $i$  at time  $t$ . This vector is initialized uniformly and updated at each round. On each round, the learner plays an action by drawing  $I_t \sim \mathbf{p}_t$  where  $p_t$  is set as follows

$$p_t(i) = (1 - \gamma) \frac{w_t(i)}{\sum_{j=1}^K w_t(j)} + \frac{\gamma}{K}, \forall i \in [K]$$

where  $w_t(i)$  is the importance weight of each arm  $i$  learned by the Hedge algorithm, and  $\gamma$  is a parameter for weighting the exploration term. Algorithm 16 gives a summary of the Exp3 algorithm. By tuning the optimal parameter of  $\gamma$ , the Exp3 algorithm is able to achieve the regret  $O(\sqrt{TK \ln K})$  in the adversarial setting.

---

**Algorithm 16: Exp3**


---

**INPUT:** parameter  $\gamma \in (0, 1]$   
**INIT:**  $w_1(i) = 1, \forall i \in [K]$   
**for**  $t = 1, 2, \dots, T$  **do**  
     Set  $p_t(i) = (1 - \gamma) \frac{w_t(i)}{\sum_{j=1}^K w_t(j)} + \frac{\gamma}{K}, \forall i \in [K]$   
     Play action by drawing  $I_t \sim \mathbf{p}_t$   
     Receive  $\ell_t(I_t) \in [0, 1]$ ,  
     Update  $w_t(i) = w_t(i) e^{-\gamma \frac{\ell_t(i)}{p_t(i)}}$ , if  $i = I_t$ .  
**end for**

---

*Other Related Works.* This is generally more challenging than the stochastic setting. A variety of algorithms have been explored in literature (Bubeck and Cesa-Bianchi, 2012). For example, the Exp3.P algorithm in (Auer et al., 2002b) improves the loss estimation and probability update strategies to get a high probability bound. The Exp3.M algorithm in (Uchiya et al., 2010) explores the new problem setting of multiple plays.

### 5.3.2 LINEAR BANDIT AND COMBINATORIAL BANDIT

We first introduce the problem setting of the *Linear Bandit* optimization problem (Auer, 2002; Rusmevichientong and Tsitsiklis, 2010; Jun et al., 2017). During each iteration, the player makes its decision by choosing a vector from a finite set  $\mathcal{S} \subseteq \mathbb{R}^d$  of elements  $\mathbf{v}(i)$  for  $i = 1, \dots, k$ . The chosen action at iteration  $t$  is indexed as  $I_t$ . The environment chooses a loss vector  $\ell_t \in \mathbb{R}^d$  and returns the linear loss as  $c_t(I_t) = \ell_t^\top \mathbf{v}(I_t)$ . Note that the player has no access to the full knowledge of loss vector and the only information revealed to the player is the loss of its own decision  $c_t(I_t)$ . Obviously, when setting  $d = k$  and  $\mathbf{v}(i)$  is the standard basis vector, this problem is identical to that in the previous section.

**Combinatorial Bandit.** *Combinatorial bandit* (Cesa-Bianchi and Lugosi, 2012) is a special case of Linear Bandits, where  $\mathcal{S}$  is a subset of binary hypercube  $\{0, 1\}^d$ . The loss vector  $\ell_t$  may be generated from an unknown but fixed distribution, which is termed as stochastic combinatorial bandit, or chosen from some adversarial environment, which is termed as adversarial combinatorial bandit. The goal is to minimize the expected regret

$$\bar{R}_T = \mathbb{E}[\sum_{t=1}^T c_t(I_t)] - \min_{i \in [K]} L_T(i)$$

where  $L_T(i) = \mathbb{E} \sum_{t=1}^T c_t(i)$  is the expected sum of loss for choosing action  $i$  in all  $T$  iterations, not a random variable in stochastic setting. One example algorithm for Combinatorial Bandit is the COMBAND algorithm (Cesa-Bianchi and Lugosi, 2012). It first defines a sampling probability vector  $\mathbf{p}_t \in \mathbb{R}^k$  for sampling  $\mathbf{v}(I_t)$  from  $\mathcal{S}$

$$\mathbf{p}_t = (1 - \gamma)\mathbf{q}_t + \gamma\boldsymbol{\mu}$$

where  $\mathbf{q}_t \in \mathbb{R}^k$  is the exploitation probability vector that is updated during all iterations to follow the best action,  $\boldsymbol{\mu} \in \mathbb{R}^d$  is a fixed exploration probability, and  $\gamma \in [0, 1]$  is the weight to control the exploitation and exploration trade-off. The algorithm draws the action  $I_t$  based on distribution  $\mathbf{p}_t$  and gets the loss  $c_t(I_t)$  from the environment. Second, an estimation of the loss vector  $\ell_t$  is calculated with the new information,

$$\tilde{\ell}_t = c_t(I_t)P_t^+ \mathbf{v}(I_t)$$

where  $P_t^+$  is the pseudo-inverse of the expected correlation matrix  $\mathbb{E}_{\mathbf{p}_t}[\mathbf{v}\mathbf{v}^\top]$ . Finally, the exploitation weights are scaled based on the estimated loss vector,

$$\mathbf{q}_{t+1}(i) \propto \mathbf{q}_t(i) \exp(-\eta \tilde{\ell}_t^\top \mathbf{v}(i))$$

where  $\eta > 0$  is a learning rate parameter and  $\propto$  indicates that this scaling step is followed by a normalization step so that  $\sum_{i=1}^k \mathbf{q}(i) = 1$ . The COMBAND algorithm achieves a regret bound better than  $O(\sqrt{Td \ln |\mathcal{S}|})$  for a variety of concrete choices of  $\mathcal{S}$ .

**Other Related Works.** Recently, many studies also address linear bandits and combinatorial bandits in different settings. The ESCB algorithm (Combes et al., 2015) efficiently exploits the structure of the problem and gets a better regret bound of  $O(\ln(T))$ . The CUCB algorithm (Chen et al., 2016b) addresses the problem where the loss may be nonlinear. (Combes et al., 2015) provided a useful survey for closely related works (Bubeck et al., 2012; Cesa-Bianchi and Lugosi, 2012) and gave a novel algorithm with promising bounds.

## 5.4 Contextual Bandits

Contextual Bandit is a widely used extension of MAB by associating contextual information with each arm (Zeng et al., 2016; Li et al., 2017b). For example, in personalized recommendation problem, the task is to select products that are most likely to be purchased by a user. In this case, each product corresponds to an arm and the features of each product are easy to acquire (Li et al., 2010b).

In a contextual bandits problem, there is a set of policies  $\mathcal{F}$ , which may be finite or infinite. Each  $f \in \mathcal{F}$  maps a context  $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^d$  to an arm  $i \in [k]$ . Different from the previous setting where the regret is defined by competing with the arm with the highest expected reward, the regret here is defined by comparing the decision  $I_t$  with the best policy  $f^* = \arg \inf_{f \in \mathcal{F}} \ell_D(f)$ , where  $D$  is the data distribution.

$$R_T(f) = \sum_{t=1}^T [\ell_{I_t, t} - \ell_t(f^*)]$$

In literature, there are comprehensive surveys on contextual bandit algorithms in both stochastic and adversarial settings (Zhou, 2015; Bubeck and Cesa-Bianchi, 2012). Below we focus on two settings of contextual bandits: multiclass classification and expert advice.

### 5.4.1 THE MULTICLASS SETTING.

In this setting, contextual bandit is regarded as a special case of online multi-class classification tasks in bandit setting. The goal is to learn a mapping from context space  $\mathbb{R}^d$  to label space  $\{1, \dots, k\}$  from a sequence of instances  $\mathbf{x}_t \in \mathbb{R}^d$ . Different from classic online multi-class classification problems where a class label  $y_t \in \{1, \dots, k\}$  is revealed at the end of each iteration, in bandit setting, the learner only gets a partial feedback on whether  $\hat{y}_t$  equals to  $y_t$ . In the following, we briefly review some representative works of contextual bandits for multi-class classification.

*Banditron* is the first bandit algorithm for online multiclass prediction (Kakade et al., 2008), which is a variant of the Perceptron. To efficiently make prediction and update the model, the Banditron algorithm keep a linear model  $W^t$ , which is initialized as  $W^1 = 0 \in \mathbb{R}^{k \times d}$ . At the  $t$ -th iteration, after receiving the instance  $\mathbf{x}_t \in \mathbb{R}^d$ , it will first set

$$\hat{y}_t = \arg \max_{r \in [k]} (W^t \mathbf{x}_t)_r$$

where  $(\mathbf{z})_r$  denotes the  $r$ -th element of  $\mathbf{z}$ . Then the algorithm will define a distribution as

$$\Pr(r) = (1 - \gamma) \mathbb{I}(r = \hat{y}_t) + \gamma/k, \forall r \in [k]$$

which roughly implies that the algorithm exploits with probability  $1 - \gamma$  and explores with the remaining probability by uniformly predicting a random label from  $[k]$ . The parameter  $\gamma$  controls the exploration-exploitation tradeoff. The algorithm then randomly sample  $\tilde{y}_t$  according to the probability  $\Pr$  and predicts it as the label of  $\mathbf{x}_t$ . After the prediction, the algorithm then receives the bandit feedback  $\mathbb{I}(\tilde{y}_t = y_t)$ . Then the algorithm uses this feedback to construct a matrix

$$\tilde{U}_{r,j}^t = x_{t,j} \left( \mathbb{I}(\hat{y}_t = r) - \frac{\mathbb{I}(\tilde{y}_t = y_t) \mathbb{I}(\tilde{y}_t = r)}{\Pr(r)} \right)$$

since its expectation satisfies  $\mathbb{E}\tilde{U}_{r,j}^t = U_{r,j}^t = x_{t,j} (\mathbb{I}(\hat{y}_t = r) - \mathbb{I}(y_t = r))$ , where  $U^t$  is actually a (sub)-gradient of the following hinge loss

$$\ell(W; (\mathbf{x}_t, y_t)) = \max_{r \in [k] / \{y_t\}} [1 - (W\mathbf{x}_t)_{y_t} + (W\mathbf{x}_t)_r]_+$$

where  $[z]_+ = \max(0, z)$ . Then the algorithm will update the model by  $W^{t+1} = W^t - \tilde{U}^t$ . The Banditron algorithm is summarized in Algorithm 17.

---

**Algorithm 17:** Banditron

---

**INIT:**  $\mathbf{w}_{1,1} = 0, \dots, \mathbf{w}_{k,1} = 0$   
**for**  $t = 1, 2, \dots, T$  **do**  
 Receive an incoming instance  $\mathbf{x}_t$   
 $P(r) = (1 - \gamma) \mathbf{1}[r = \arg \max_i \mathbf{w}_{i,t}^\top \mathbf{x}_t] + \frac{\gamma}{k}$ .  
 Sample  $\hat{y}_t$  according to  $P(r), r \in \{1, \dots, k\}$   
 $\mathbf{u}_r = \mathbf{x}_t \left( \frac{\mathbf{1}[y_t = \hat{y}_t = r]}{P(r)} - \mathbf{1}[r = \arg \max_i \mathbf{w}_{i,t}^\top \mathbf{x}_t] \right)$ ;  
 $\mathbf{w}_{r,t+1} = \mathbf{w}_{r,t} + \mathbf{u}_r$   
**end for**

---

This algorithm achieves  $O(\sqrt{T})$  in linear separable case and  $O(T^{\frac{2}{3}})$  in inseparable case.

*Other Related Work.* Following the Banditron, many algorithms have been proposed. For example, Bandit Passive Aggressive (Bandit PA) follows the PA learning principle and adopts the framework of one against all others to make prediction and update the model (Chen et al., 2009). In general, some update principles are based on first-order gradient descent (Wang et al., 2010), while others adopt second order learning (Crammer and Gentile, 2011; Hazan and Kale, 2011; Zhang et al., 2016a; Beygelzimer et al., 2017). Most of these algorithms explore the  $k$  classes uniformly with probability  $\gamma$ , while (Crammer and Gentile, 2011) sample the classes based on the Upper Confidence Bound.

#### 5.4.2 THE EXPERT SETTING.

We now introduce a well-known algorithm called Exp4 (Auer et al., 2002b) for contextual bandits in the expert settings. The Exp4 algorithm assumes that there are  $N$  experts who will give advice on the distribution over arms during all iterations.  $\xi_{i,t}^n$  indicates the probability of picking arm  $i \in [K]$  recommended by expert  $n \in [N]$  during time  $t \in [T]$ . Obviously,  $\sum_{i=1}^k \xi_{i,t}^n = 1$ . During time  $t$ , the true reward vector is denoted by  $\mathbf{r}_t \in [0, 1]^K$ . Thus the expected reward of expert  $n$  is  $\boldsymbol{\xi}_t^n \cdot \mathbf{r}_t$ . The regret is defined by comparing with the expert with the highest expected cumulative reward.

$$R_t = \max_{n \in [N]} \sum_{t=1}^T \boldsymbol{\xi}_t^n \cdot \mathbf{r}_t - \mathbb{E} \sum_{t=1}^T r_{t, I_t}$$

The Exp4 algorithm first defines a weight vector  $\mathbf{w}_t \in \mathbb{R}^N$  that indicates the weights for the  $N$  experts. We set the weight as  $\mathbf{w}_0 = \mathbf{1}$  and update it during each iteration. During iteration  $t$ , we calculate the probability of picking arm  $i$  as the weighted sum of advices



from all  $N$  experts,

$$p_{i,t} = (1 - \gamma) \frac{\sum_{n=1}^N w_{n,t} \xi_{i,t}^n}{\sum_{n=1}^N w_{n,t}} + \frac{\gamma}{K}$$

where  $\gamma \in [0, 1]$  is a parameter to balance exploitation and exploration. We then draw an arm  $I_t$  according to probability  $p_{i,t}$  and calculate an unbiased estimator of  $\hat{r}_{i,t} = \frac{r_{i,t}}{p_{i,t}} \mathbb{I}_{i=I_t}$ , which will be used to calculate the expected reward. Finally the weight  $\mathbf{w}_t$  is updated according to the expected reward of each arm. The Exp4 algorithm is able to achieve the regret bound  $O(\sqrt{TK \ln N})$  as shown in (Auer et al., 2002b) and tighter bounds were also given in (McMahan and Streeter, 2009).

*Other Related Work.* Another general contextual bandit algorithm is the epoch-greedy algorithm in (Langford and Zhang, 2008) that is similar to  $\epsilon$ -greedy with shrinking  $\epsilon$ . This algorithm is computationally efficient given an oracle optimizer but has the weaker regret guarantee of  $O(T^{2/3})$ . LinUCB (Chu et al., 2011) is an extension of UCB to contextual bandit problem, by assuming that there is a feature vector  $\mathbf{x}_{t,i} \in \mathbb{R}^d$  at time  $t$  for each arm  $i$ . Similar to the UCB algorithm, a model is learnt to estimate the upper confidence bound of each arm  $i \in [k]$  given the input of  $\mathbf{x}_{t,i}$ . The algorithm simply chooses the arm with the highest UCB. The LinREL algorithm (Auer, 2002) is similar to LinUCB in that it adopts the same problem setting and same maximizing UCB strategy. While, a different regularization term is used which leads to a different calculation of the UCB.

## 5.5 Other Bandit Variants

In literature, there are many other studies addressing on various types of bandit variants. We refer readers for more comprehensive studies on bandit topics in (Bubeck and Cesa-Bianchi, 2012). Below we briefly introduce a few other major variants.

Other than stochastic bandits and adversarial bandits, another fundamental topic of multi-armed bandits is called “*Markovian bandits*”, which generally assumes the reward processes are neither i.i.d. (like in stochastic MAB) nor adversarial. Specifically, arms are associated with  $K$  Markov processes, each with its own state space. On each round, an arm is chosen in some state, a stochastic reward is drawn from some probability distribution, and the state of the reward process for the arm changes in a Markovian fashion, based on an underlying stochastic transition matrix. Both reward and new state are revealed to the player. The seminal work of Gittins (1979) gives an optimal greedy policy that can be computed efficiently. A special case of Markovian bandits is Bayesian bandits (Scott, 2010; Gittins et al., 2011; Kaufmann et al., 2012), which are parametric stochastic bandits where the parameters of the reward distributions are assumed to be drawn from known priors, and the regret is computed by also averaging over the draw of parameters from the prior.

Another topic is to study *infinitely many-armed bandits* problems where the number of arms can be larger than the possible number of experiments or even infinite (Berry et al., 1997; Wang et al., 2009; Bubeck et al., 2011). Among these studies, one niche sub-topic is *continuum-armed bandits* (Kleinberg, 2005a; Chowdhury and Gopalan, 2017), where the arms lie in some Euclidean (or metric) space and their mean-reward is a deterministic and smooth (e.g., Lipschitz) function of the arms, a.k.a. *Lipschitz Bandit* (Magureanu et al., 2014).

## 6. Online Active Learning

### 6.1 Overview

In a standard online learning task (e.g., online binary classification), the learner receives and makes prediction for a sequence of instances generated from some unknown distribution. At the end of every round, it always assumes the learner will receive the true label (feedback) from the environment. For many real-world applications, obtaining the labels could be very expensive, and sometimes it is not always necessary/informative to query the true labels of every instance, e.g., if an instance is correctly classified with a high confidence. Motivated to address this challenge, online active learning is a special class of online learner that observes a sequence of unlabeled instances each time deciding whether to query the label of the incoming instance; if the label is queried, then the learner can use the labelled instance to update the prediction model; otherwise, the model will be kept unchanged.

In literature, there are two major kinds of settings for online active learning. One is called the “*selective sampling*” setting (Atlas et al., 1990; Freund et al., 1997; Cesa-Bianchi et al., 2003) by adapting classical online learning for active learning. The other is *online active learning with expert advice* by adapting the setting of prediction with expert advice for active learning (Helmbold and Panizza, 1997; Zhao et al., 2013). Both operate in the similar problem settings where true label of an instance is only queried when some condition is satisfied, e.g., predictive confidence is below some threshold.

### 6.2 Selective Sampling Algorithms

In this section we review a family of popular Selective Sampling (SS) algorithms for online active learning tasks. In the following discussions, we use a typical online binary classification task as a running example. For notation, an example is a pair  $(\mathbf{x}, y)$ , where  $\mathbf{x} \in \mathbb{R}^d$  is an instance vector and  $y \in \{-1, +1\}$  is the binary class label. Assume the learning proceeds in a sequence of  $T$  rounds, where  $T$  may not be known in advance. On each round  $t$ , a learner observes an instance  $\mathbf{x}_t$ , then outputs a prediction  $\hat{y}_t \in \{-1, +1\}$  as the label for the instance, and then decides whether or not to query the label  $y_t$ . Whenever  $\hat{y}_t \neq y_t$ , the learner’s prediction outcome is considered as a mistake, no matter if it has decided to query the label or not. For notation, we denote  $M_t = \mathbb{I}(\hat{y}_t \neq y_t) \in \{0, 1\}$  as an indicator whether the learner makes a mistake at round  $t$ . For most cases, we also assume the learner adopts a linear model to predict the class label using  $\hat{y}_t = \text{sign}(\hat{p}_t)$ , where  $\hat{p}_t = \mathbf{w}_t^\top \mathbf{x}_t$ .

#### 6.2.1 FIRST-ORDER SELECTIVE SAMPLING ALGORITHMS

**Selective-sampling Perceptron.** This algorithm (Cesa-Bianchi et al., 2006) decides whether or not to query the label  $y_t$  through a simple randomized rule: drawing a Bernoulli random variable  $Z_t \in \{0, 1\}$  with probability

$$\Pr(Z_t = 1) = \frac{\delta}{\delta + |\hat{p}_t|}$$

where  $\delta > 0$  is a smooth parameter that can be used to control the number of labels queried during the online active learning process. If  $\delta$  increases, the number of queried labels increases. If  $Z_t = 1$ , then the label  $y_t$  of  $\mathbf{x}_t$  will be queried, and the model will be

---

**Algorithm 18:** Selective-sampling Perceptron
 

---

**INPUT:** parameter  $\delta > 0$   
**INIT:**  $\mathbf{w}_0 = (0, \dots, 0)^\top$   
**for**  $t = 1, 2, \dots, T$  **do**  
     Observe an input instance  $\mathbf{x}_t \in \mathbb{R}^d$   
     Predict  $\hat{y}_t = \text{sign}(\hat{p}_t)$ , where  $\hat{p}_t = \mathbf{w}_t^\top \mathbf{x}_t$   
     Draw a Bernoulli random variable  $Z_t \in \{0, 1\}$  of probability  $\frac{\delta}{\delta + |\hat{p}_t|}$   
     **IF**  $Z_t = 1$  **THEN**  
         Query label  $y_t \in \{-1, +1\}$  and Update  $\mathbf{w}_t$  by Perceptron:  $\mathbf{w}_{t+1} = \mathbf{w}_t + M_t y_t \mathbf{x}_t$ .  
     **end for**

---

updated using the Perceptron rule. Algorithm 18 gives a summary of the Selective-sampling Perceptron algorithm.

In theory, assuming  $\|\mathbf{x}_t\| \leq R$ , for any  $\mathbf{w} \in \mathbb{R}^d$ , the expected number of mistakes of the Selective-sampling Perceptron algorithm can be bounded as:

$$\mathbb{E}[\sum_{t=1}^T M_t] \leq (1 + \frac{R^2}{2\delta}) \frac{\bar{L}_{\gamma,T}(\mathbf{w})}{\gamma} + \frac{\|\mathbf{w}\|^2 (2\delta + R^2)^2}{8\delta\gamma^2}.$$

where  $\bar{L}_{\gamma,T}(\mathbf{w}) = \mathbb{E}[\sum_{t=1}^T Z_t M_t \ell_{\gamma,t}(\mathbf{w})]$ , and  $\ell_{\gamma,t}(\mathbf{w}) = \max(0, \gamma - y_t \mathbf{w}^\top \mathbf{x}_t)$ . Furthermore, the expected number of labels queried by the algorithm equals  $\sum_{t=1}^T \mathbb{E}[\frac{\delta}{\delta + |\hat{p}_t|}]$ . This bound depends on the value of the parameter  $\delta$ . By choosing the optimal value of  $\delta$  as

$$\delta = \frac{R^2}{2} \sqrt{1 + \frac{4\gamma^2}{\|\mathbf{w}\|^2 R^2} \frac{\bar{L}_{\gamma,T}(\mathbf{w})}{\gamma}}$$

the expected number of mistakes can be bounded

$$\mathbb{E}[\sum_{t=1}^T M_t] \leq \frac{\bar{L}_{\gamma,T}(\mathbf{w})}{\gamma} + \frac{\|\mathbf{w}\|^2 R^2}{2\gamma^2} + \frac{\|\mathbf{w}\| R}{\gamma} \sqrt{\frac{\bar{L}_{\gamma,T}(\mathbf{w})}{\gamma} + \frac{\|\mathbf{w}\|^2 R^2}{4\gamma^2}}$$

This is an expectation version of the mistake bound for the standard Perceptron Algorithm. Especially, in the special case when the data is linearly separable, the optimal value of  $\delta$  is  $R^2/2$  and this bound becomes the familiar Perceptron bound  $(\|\mathbf{w}\| R)^2 / \gamma^2$ . Instead of using a fixed constant parameter, (Cesa-Bianchi et al., 2006) also proposed an adaptive parameter version of the selective sampling Perceptron algorithm as follows:

$$\Pr(Z_t = 1) = \frac{\delta_t}{\delta_t + |\hat{p}_t|}, \quad \text{s.t.} \quad \delta_t = \beta(R')^2 \sqrt{1 + \sum_{i=1}^{t-1} Z_i M_i},$$

where  $\beta > 0$  is a predefined parameter,  $R' = \max R_{t-1}, \|\mathbf{x}_t\|$ ,  $R_{t-1} = \max\{\|\mathbf{x}_i\| | Z_i M_i = 1\}$ .

**Other first-order approaches.** Instead of using Perceptron, the Passive-Aggressive Active learning algorithms in (Lu et al., 2014, 2016b) are selective sampling algorithms by extending the framework of PA online learning algorithms. They also extended their algorithms for multi-class classification and cost-sensitive classification tasks. (Zhao and Hoi, 2013) proposed a cost-sensitive online active learning approach that directly optimizes cost-sensitive measures using PA-like algorithms for class-imbalanced classification tasks.

### 6.2.2 SECOND-ORDER SELECTIVE SAMPLING ALGORITHMS

**Selective-sampling Second-order Perceptron.** Instead of using the standard Perceptron algorithm, (Cesa-Bianchi et al., 2006) also proposed a selective-sampling algorithm based on the Second-order Perceptron.

---

**Algorithm 19:** Selective-sampling Second-order Perceptron

---

**INPUT:** parameter  $\delta > 0$   
**INIT:**  $A_0 = I$ ,  $\mathbf{w}_0 = (0, \dots, 0)^\top$   
**for**  $t = 1, 2, \dots, T$  **do**  
    Observe an input instance  $\mathbf{x}_t \in \mathbb{R}^d$   
    Computer  $\hat{p}_t = [(A_t + \mathbf{x}_t \mathbf{x}_t^\top)^{-\frac{1}{2}} \mathbf{u}_t]^\top [(A_t + \mathbf{x}_t \mathbf{x}_t^\top)^{-\frac{1}{2}} \mathbf{x}_t] = \mathbf{u}_t^\top (A_t + \mathbf{x}_t \mathbf{x}_t^\top)^{-1} \mathbf{x}_t$   
    Predict  $\hat{y}_t = \text{sign}(\hat{p}_t)$   
    Draw a Bernoulli random variable  $Z_t \in \{0, 1\}$  of probability  $\frac{\delta}{\delta + |\hat{p}_t|}$   
    **IF**  $Z_t = 1$  **THEN**  
        Query label  $y_t \in \{-1, +1\}$  and Update  $\mathbf{w}_t$  by Second-order Perceptron:  
         $\mathbf{u}_{t+1} = \mathbf{u}_t + M_t y_t \mathbf{x}_t$ , and  $A_{t+1} = A_t + M_t \mathbf{x}_t \mathbf{x}_t^\top$   
    **end for**

---

Let  $\mathbf{u}_t$  denote the weight vector computed by standard Perceptron, and  $A_t = I + \sum_{i \leq t-1, Z_i M_i = 1} \mathbf{x}_i \mathbf{x}_i^\top$  denote the correlation matrix over the mistaken trials plus an identity matrix  $I$ , then the second-order Perceptron predicts the label of current instance  $\mathbf{x}_t$  as

$$\hat{y}_t = \text{sign}(\hat{p}_t), \text{ where } \hat{p}_t = [(A_t + \mathbf{x}_t \mathbf{x}_t^\top)^{-\frac{1}{2}} \mathbf{u}_t]^\top [(A_t + \mathbf{x}_t \mathbf{x}_t^\top)^{-\frac{1}{2}} \mathbf{x}_t] = \mathbf{u}_t^\top (A_t + \mathbf{x}_t \mathbf{x}_t^\top)^{-1} \mathbf{x}_t$$

The second-order algorithm differs from standard Perceptron in that, before each prediction, a linear transformation  $(A_t + \mathbf{x}_t \mathbf{x}_t^\top)^{-1/2}$  is applied to both current Perceptron weight  $\mathbf{u}_t$  and current instance  $\mathbf{x}_t$ . After prediction, the query strategy of this algorithm is the same with the previous selective sampling: draw a Bernoulli random variable  $Z_t \in \{0, 1\}$  with

$$\Pr(Z_t = 1) = \frac{\delta}{\delta + |\hat{p}_t|}.$$

Algorithm 19 gives a summary of the Selective-sampling Second-order Perceptron algorithm. In theory, if the algorithm runs on a sequence of  $T$  rounds, for any  $\mathbf{w}$ , the expected number of mistakes made by the algorithm is bounded:

$$\mathbb{E}[\sum_{t=1}^T M_t] \leq \frac{\bar{L}_{\gamma, T}(\mathbf{w})}{\gamma} + \frac{\delta}{2\gamma^2} \mathbf{w}^\top \mathbb{E}[A_T] \mathbf{w} + \frac{1}{2\delta} \sum_{i=1}^d \mathbb{E} \ln(1 + \lambda_i)$$

where  $\bar{L}_{\gamma,T}(\mathbf{w}) = \mathbb{E}[\sum_{t=1}^T Z_t M_t \ell_{\gamma,t}(\mathbf{w})]$  with  $\ell_{\gamma,t}(\mathbf{w}) = \max(0, \gamma - y_t \mathbf{w}^\top \mathbf{x}_t)$ ,  $\lambda_1, \dots, \lambda_d$  are the eigenvalues of the random correlation matrix  $\sum_{t=1}^T Z_t M_t \mathbf{x}_t \mathbf{x}_t^\top$  and  $A_T = I + \sum_{t=1}^T M_t Z_t \mathbf{x}_t \mathbf{x}_t^\top$ . Moreover, the expected number of queries by the algorithm equals  $\sum_{t=1}^T \mathbb{E}[\frac{\delta}{\delta + |\hat{p}_t|}]$ .

Furthermore, by setting  $\delta = \gamma \sqrt{\frac{\sum_{i=1}^d \mathbb{E} \ln(1 + \lambda_i)}{\mathbf{w}^\top \mathbb{E}[A_T] \mathbf{w}}}$ , it leads to the optimal bound

$$\mathbb{E}[\sum_{t=1}^T M_t] \leq \frac{\bar{L}_{\gamma,T}(\mathbf{w})}{\gamma} + \frac{1}{\gamma} \sqrt{(\mathbf{w}^\top \mathbb{E}[A_T] \mathbf{w}) \sum_{i=1}^d \mathbb{E} \ln(1 + \lambda_i)}$$

(Cesa-Bianchi et al., 2006) also proposed an improved selective-sampling algorithm based on second-order Perceptron, which modifies the sampling probability by incorporating the second-order information, i.e., with

$$\Pr(Z_t = 1) = \frac{\delta}{\delta + |\hat{p}_t| + \frac{1}{2} \hat{p}_t^2 (1 + \mathbf{x}_t^\top A_t^{-1} \mathbf{x}_t)},$$

**Other second-order approaches.** Cesa-Bianchi et al. (2003) proposed a margin-based selective sampling algorithm which also exploits second-order information in the model:

$$\hat{y}_t = \text{sign}(p_t), \text{ where } p_t = \mathbf{w}_t^\top \mathbf{x}_t, \mathbf{w}_t = A_t^{-1} \mathbf{u}_t, \mathbf{u}_t = \sum_{i=1}^{t-1} Z_i y_i \mathbf{x}_i^\top, A_t = (I + \sum_{i=1}^{t-1} Z_i \mathbf{x}_i \mathbf{x}_i^\top)$$

But the query strategy is a margin-based sampling approach without explicitly exploiting the second-order information:

$$\Pr(Z_{t+1} = 1) = \mathbb{I}(p_t \leq \frac{4 \ln t}{\sum_{i=1}^{t-1} Z_i})$$

Hao et al. (2016, 2017a) proposed second-order online active learning algorithms by fully exploiting both the first-order and second-order information for online active learning tasks and also gave cost-sensitive extensions for class-imbalanced tasks.

### 6.2.3 OTHER SELECTIVE SAMPLING APPROACHES

There are also a few other selective sampling approaches in which the base classifier is based on the Regularized Least Squares (RLS). In particular, on each round  $t$ , the linear classification model can be updated by the RLS estimate

$$\mathbf{w}_t = (I + S_{t-1} S_{t-1}^\top + \mathbf{x}_t \mathbf{x}_t^\top)^{-1} S_{t-1} Y_{t-1}$$

where matrix  $S_{t-1} = [\mathbf{x}'_1, \dots, \mathbf{x}'_{N_{t-1}}]$  is the collection of  $N_{t-1}$  instances queried up to time  $t - 1$ , and the vector  $Y_{t-1} = (Y'_1, \dots, Y'_{N_{t-1}})$  is the set of queried labels for the instances. The selective sampling algorithms that follow this paradigm include the Bound on Bias Query (BBQ) algorithms (Cesa-Bianchi et al., 2009; Orabona and Cesa-Bianchi, 2011) and their improved variants Dekel et al. (2010); Orabona and Cesa-Bianchi (2011). A major drawback of these methods is that the RLS-based base learner is more like a fashion of batch learner instead of truly online learning, and thus the overall learning scheme might be inefficient or non-scalable if the number of queried labeled examples can be large.

### 6.3 Online Active Learning with Expert Advice

The idea of online active learning with expert advice dates back to classical Query by Committee (QBC) in (Seung et al., 1992; Freund et al., 1997), where the idea is to query the label of an instance based on the principle of *maximal disagreement* among a set of experts, i.e., the confidence criteria in this case is how much the expert hypotheses disagree on their evaluation of instance predictions. QBC bounds from below the average information gain provided by each requested label. Baram et al. (2004) considers the setting of how to online combine an ensemble of active learners, which is executed based on a maximum entropy criterion. Another perhaps more dominating line of studies in (Helmbold and Panizza, 1997; Cesa-Bianchi et al., 2005b; Zhao et al., 2013) explore the exponentiated weighted average forecaster for online active learning tasks, where an instance is stochastically queried based on the available feedback on the importance of each expert in the pool.

Next we describe in detail one of the most recent approaches for online active learning with expert advice in (Zhao et al., 2013). Consider an unknown sequence of instances  $\mathbf{x}_1, \dots, \mathbf{x}_T \in \mathbb{R}^d$ , a “forecaster” aims to predict the class labels of every incoming instance  $\mathbf{x}_t$ . The forecaster sequentially computes its predictions based on the predictions from a set of  $N$  “experts”. Specifically, at the  $t$ -th round, after receiving an instance  $\mathbf{x}_t$ , the forecaster first accesses the predictions of the experts  $\{f_{i,t} : \mathbb{R}^d \rightarrow [0, 1] | i = 1, \dots, N\}$ , and then computes its own prediction  $p_t \in [0, 1]$  based on the predictions of the  $N$  experts. After  $p_t$  is computed, the true outcome  $y_t \in \{0, 1\}$  is disclosed. To solve this problem, the “Exponentially Weighted Average Forecaster” (EWA) makes the following prediction:

$$p_t = \frac{\sum_{i=1}^N \exp(-\eta L_{i,t-1}) f_i(\mathbf{x}_t)}{\sum_{i=1}^N \exp(-\eta L_{i,t-1})}, \quad (13)$$

where  $\eta$  is a learning rate,  $L_{i,t} = \sum_{j=1}^t \ell(f_i(\mathbf{x}_j), y_j)$ ,  $L_t = \sum_{j=1}^t \ell(p_j, y_j)$  with  $\ell(p_t, y_t) = |p_t - y_t|$ . Unlike the above regular learning, in an active learning with expert advice task, the outcome of an incoming instance is *only* revealed whenever the learner requests the label from the environment/oracle. To solve this problem, binary variables  $z_s \in \{0, 1\}$ ,  $s = 1, \dots, t$  are introduced to indicate if an active forecaster has requested the label of an instance at  $s$ -th trial.  $\hat{L}_{i,t}$  is used to denote the loss function experienced by the active learner w.r.t. the  $i^{th}$  expert, i.e.,  $\hat{L}_{i,t} = \sum_{s=1}^t \ell(f_i(\mathbf{x}_s), y_s) z_s$ . For this problem setting, Zhao et al. (2013) proposed a general framework of active forecasters, as shown in Algorithm 20.

---

**Algorithm 20:** Online Active Learning with Expert Advice

---

**INPUT:** a pool of experts  $f_i$ ,  $i = 1, \dots, N$ .  
**INIT:** tolerance threshold  $\delta$  and  $\hat{L}_{i,t} = 0$ ,  $i \in [N]$ .  
**for**  $t = 1, 2, \dots, T$  **do**  
  Receive  $\mathbf{x}_t$  and compute  $f_i(\mathbf{x}_t)$ ,  $i \in [N]$ ;  
  Compute  $\hat{p}_t = \frac{\sum_{i=1}^N \exp(-\eta \hat{L}_{i,t-1}) f_i(\mathbf{x}_t)}{\sum_{i=1}^N \exp(-\eta \hat{L}_{i,t-1})}$ ;  
  If a *confidence condition* is not satisfied  
    request label  $y_t$  and update  $\hat{L}_{i,t} = \hat{L}_{i,t-1} + \ell(f_i(\mathbf{x}_t), y_t)$ ,  $i \in [N]$ ;  
**end for**

---

At each round, after receiving an instance  $\mathbf{x}_t$ , we compute the prediction of class label for the instance by aggregating the prediction of each expert in the pool, i.e.,  $f_i(\mathbf{x}_t)$ . Then, we examine if a confidence condition is satisfied. If so, we will skip the label request; otherwise, the learner will request the class label for this instance from the environment. To decide when to request the class label or not, the key idea is to seek a confidence condition by estimating the difference between  $p_t$  and  $\hat{p}_t$ . Intuitively, the smaller the difference, the more confident we have for the prediction made by the forecaster. More specifically, the work in (Zhao et al., 2013) proved that for a small constant  $\delta > 0$ ,  $\max_{1 \leq i, j \leq N} |f_i(\mathbf{x}_t) - f_j(\mathbf{x}_t)| \leq \delta$  implies  $|p_t - \hat{p}_t| \leq \delta$ . This roughly means that, if any two experts do not disagree with each other too much on the instance, then we can skip requiring its label.

In addition to the above work, there are also a few other active learning strategies for online learning with expert advices, for example the active greedy forecaster (Zhao et al., 2013). Online active learning with expert advice can be applied in some real-world applications, e.g., crowdsourcing tasks (Hao et al., 2015a) where the learner attempts to address both the diverse quality of annotators' performance with expert learning and efficient annotation in seeking informative data using active learning.

## 7. Online Semi-supervised Learning

### 7.1 Overview

Semi-Supervised Learning (SSL) has been an important class of machine learning tasks and techniques, which aims to make use of unlabeled data for learning tasks. It has been extensively studied mostly in the settings of batch learning and some comprehensive surveys can be found in (Zhu, 2006; Zhu and Goldberg, 2009). When online learning meets semi-supervised learning, there are two major branches of research. One major branch of studies is to turn traditional batch semi-supervised learning methods into online algorithms such that they can work from data streams of both labeled and unlabeled data, which we call this setting as "Online Semi-supervised Learning" and we will review a popular framework of "online manifold regularization". The other branch of studies is to study classical online learning tasks in transductive learning settings (e.g., by assuming unlabeled data can be made available before online learning tasks), which we call this setting as "Transductive Online Learning". We note that online active learning as introduced previously can be generally viewed as a special type of online semi-supervised learning where an online learner has to deal with both labeled and unlabeled data.

### 7.2 Online Manifold Regularization

In the area of semi-supervised learning, one major framework for semi-supervised learning is based on manifold regularization (Belkin et al., 2006), where the learner not only minimizes the loss on the labeled data, but also minimizes the difference of predictions on the unlabeled instances which are similar on the manifold. Specifically, consider instances  $(\mathbf{x}_t, y_t)$ ,  $t \in \{1, \dots, T\}$ , the idea is to minimize the following objective function

$$J(f) = \frac{1}{l} \sum_{t=1}^T \delta(y_t) \ell(f(\mathbf{x}_t), y_t) + \frac{\lambda_1}{2} \|f\|^2 + \frac{\lambda_2}{2T} \sum_{s,t=1}^T (f(\mathbf{x}_s) - f(\mathbf{x}_t)) w_{st}$$

where the first term is the loss on labeled instances where  $\delta(y_t) = 1$  if and only if  $y_t$  exists and  $l$  is the number of labeled data, the second term is a classic regularization term for supervised learning, and the last term is the manifold regularization on unlabeled data.

In literature, online manifold regularization has been explored (Goldberg et al., 2008), which attempts to turn batch manifold regularization algorithms into online/incremental algorithms. Specifically, the above objective can be returned online for each instance:

$$J_t(f) = \frac{T}{l} \delta(y_t) \ell(f(\mathbf{x}_t), y_t) + \frac{\lambda_1}{2} \|f\|^2 + \lambda_2 \sum_{i=1}^{t-1} (f(\mathbf{x}_i) - f(\mathbf{x}_t)) w_{it}$$

It can be solved using Online Gradient Descent in  $O(T^2)$  time. Unfortunately, such straightforward solution is expensive in both time and space, since the calculation of the last term requires to store all instances and measure the similarity  $w_{it}$  between the incoming instances and all existing ones.

To address this problem, the authors offer two sparse approximations of the objective function. The first solution is not to keep all instances but to keep only the newest  $\tau$  ones, where  $\tau$  is the buffer size. This strategy is simple but not very efficient since the discarded old instances may contain important information. The second solution adopts a random projection tree to find  $s$  cluster centers during online learning. Finally, instead of calculating the similarity between  $\mathbf{x}_t$  and all existing instances, the algorithm only consider the  $s$  cluster centers as the most representative instances.

In addition to the above work, (Valko et al., 2010) proposed a fast approximate algorithm for online semi-supervised learning. which leverages the incremental k-center quantization method to group neighboring points so as to yield a set of reduced representative points, and as a result an approximate similarity graph can be constructed to find the harmonic solution in semi-supervised learning (Zhu et al., 2003).

Finally, there were some related efforts on online active semi-supervised learning (Goldberg et al., 2011), which extends active learning in the online semi-supervised learning settings. For example, following such kind of setting, (Goldberg et al., 2011) developed the OASIS algorithm by using a general online Bayesian learning framework.

### 7.3 Transductive Online Learning

Transductive online learning (Cesa-Bianchi and Shamir, 2013; Ben-David et al., 1997) is a niche class of online learning tasks, where we want to learn from an arbitrary sequence of labeled examples  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$  by making the assumption that the set of unlabeled instances  $(\mathbf{x}_1, \dots, \mathbf{x}_T)$  can be given in advance to the learner before an online learning task begins. In particular, the work (Cesa-Bianchi and Shamir, 2013) proposed an efficient algorithm based on the principle of prediction with expert advice by combining “random payout” and “randomized rounding” of loss subgradients. We note that this niche topic has received very few attention, possibly because of their assumption of obtaining unlabeled data in advance, which may be unrealistic in many applications. However, the studies in this niche family of studies may provide some theory insights about the linkage between online learning and batch learning as demonstrated in (Cesa-Bianchi and Shamir, 2013).



## 8. Online Unsupervised Learning

### 8.1 Overview

In this section we briefly review some key work in the literature of online unsupervised learning, where models are learned from unlabeled data streams where no explicit feedback is available. Broadly, we categorize the existing work into four major groups: Online Clustering, Online Dimension Reduction, Online Anomaly Detection, and Online Density Estimation. Due to the vast number of ways in which unsupervised learning in online settings have been explored in literature, and numerous applications for which algorithms are designed, it is almost impossible to make a comprehensive treatment on this topic in this survey. Instead, we try to focus on the key areas and give a general overview of the main ideas in each area which are closely related to online learning.

### 8.2 Online Clustering

Clustering is an unsupervised learning process of grouping unlabeled data instances such that instances in the same group are similar, and instances between groups are dissimilar. It gives an effective mechanism to summarize the data, and does not require labels of the instances in order to perform the clustering. For batch learning settings, clustering is usually classified into following categories: partition based clustering, hierarchical clustering, density-based clustering, and grid-based clustering (Berkhin, 2006). For online settings (Aggarwal, 2013), partition-based and density-based clustering have been studied more extensively. In the following we briefly review some of online learning approaches for clustering on streaming data especially for these two categories.

*Partitioning Based* clustering methods split the instances into partitions where each partition represents a cluster. The partitions are designed on the basis of some distance measures (e.g. Euclidean distance). The number of clusters is usually pre-defined by the user. The most popular algorithms in this category are those based on *k-MEANS* and *k-MEDOIDS* algorithm. The *k-MEANS* algorithm is one of the oldest and most popular clustering methods, where the idea is to identify  $k$  centroids, where each centroid corresponds to one out of  $k$  clusters by minimizing the sum of square errors between each instance to their corresponding centroids. Sequential algorithms performing *k-MEDOID* or *k-MEDIAN* clustering usually try to break the stream of instances into chunks where the size of each chunk is set based on some pre-specified memory budget. Given a data stream  $D$ , it is broken into several chunks denoted by  $D_1, D_2, \dots, D_t, \dots$  where each chunk contains at most  $m$  instances, where  $m$  is the budget of the chunks. In such a case, *k-MEDIAN*s can be directly applied to each chunk. This framework is called the *STREAM* framework. (Guha et al., 2000; O’callaghan et al., 2002; Guha et al., 2003). There are also sampling approaches designed for clustering when the data streams are extremely large (Kaufman and Rousseeuw, 2008). Another method is the *StreamKM++* (Ackermann et al., 2012). In this approach, first, an adaptive non-uniform sampling approach is used to obtain small coresets from the data streams. The coreset construction is done by the utilization of coreset tree proposed in this paper which helps in significant speed up.

*Density-based clustering.* Most clustering techniques suffer from several drawbacks. First, many of them (e.g. *k-MEANS*) are designed for only spherical clusters and can

not adapt to arbitrary cluster shapes. In addition, the value of  $k$ , or the number of clusters has to be known a priori. Lastly, these methods are susceptible to outliers. Density based clustering algorithms (most popularly DBSCAN and its variants) are able to address all these challenges. Density based approaches cluster dense regions which are separated by sparse regions. A cluster based on density can take on arbitrary shapes, does not require prior knowledge of the number of clusters, and is robust to outliers in the data. However, performing density based clustering on streaming data in an online manner is plagued with several challenges including dynamic evolution of the clusters, limited memory space, etc. Following (Aggarwal, 2013; Amini et al., 2014), we categorize the online density-based clustering algorithms into *Micro-clustering Algorithms* and *Grid-based clustering Algorithms*. The micro-clustering algorithms aim to summarize a data in an online manner, and the clustering is performed using these summaries (Cao et al., 2006; Tasoulis et al., 2007; Ruiz et al., 2009; Li-xiong et al., 2009; Ren and Ma, 2009; Ntoutsi et al., 2012). Grid-based methods, divide the entire instance space into grids, and each instance upon arrival is assigned a grid, and then the clustering is then done based on the density of the grids (Gao et al., 2005; Chen and Tu, 2007; Jia et al., 2008; Tu and Chen, 2009; Wan et al., 2009; Ren et al., 2011; Amini and Ying, 2012; Bhatnagar et al., 2014).

*Other Clustering Methods.* *Hierarchical Clustering* is a paradigm in which either a bottom-up approach or a top-down approach is used to gradually agglomerate the data points together. This results in a tree of clusters, which is also called a dendrogram. Among the earliest approaches to incremental hierarchical clustering was CobWeb (Fisher, 1987), which determines how to insert a new data point into the tree structure based on a category utility criteria. Recent hierarchical clustering algorithms include the ClusTree (Kranen et al., 2011), which offers a compact self adapting index structure for storing stream summaries in addition to giving more importance to recent data, and Perch (Kobren et al., 2017), which allows the clustering to scale to a large number of data points and clusters. (Tu et al., 2012) propose an incremental approach to do Hierarchical clustering of the data, in addition to accounting for variance and density of the data. Some techniques have been developed for online clustering for very high dimensional data where the data sparsity makes it very hard to perform clustering as many instances tend to be equidistant from one another. HPStream (Aggarwal et al., 2004) introduces a concept of projected clustering to data streams. There are online clustering algorithms for other specific scenarios such as clustering of discrete and categorical streams, text streams, uncertain data streams, graph streams as well as distributed clustering (Aggarwal, 2013).

### 8.3 Online Dimension Reduction

When the feature dimensions are very large, Dimension Reduction (DR) techniques can be used to improve learning efficiency, compress original data, visualize data better, and improve its applicability to real-world applications. Consider instance  $\mathbf{x}_t \in \mathbb{R}^d$ , the goal of dimension reduction is to learn a new instance  $\hat{\mathbf{x}}_t \in \mathbb{R}^k$  where  $k \leq d$  by following some principle of unsupervised learning. There have been several approaches to unsupervised dimensional reduction. We broadly categorize them into two major groups of studies: subspace learning and manifold learning. More comprehensive surveys of classic dimension reduction techniques can be found in (Burgess et al., 2010).

**Subspace Learning.** This class of DR methods aims to find an optimal linear mapping of input data in high-dimensional space to a lower-dimensional space. In general, there are two major types of approaches: linear methods and nonlinear methods. Popular linear subspace methods include Principal Component Analysis (PCA) and Independent Component Analysis (ICA), etc. Nonlinear methods often extend the linear subspace learning methods using kernel tricks. Examples include Kernel PCA, Kernel ICA, etc.

For online dimension reduction tasks, more popular efforts have been focused on addressing online PCA for unsupervised learning on streaming data settings in literature (Warmuth and Kuzmin, 2008; Arora et al., 2013, 2012a; Mitliagkas et al., 2013; Feng et al., 2013), while there are also a few studies for online ICA (Li et al., 2016b; Wang and Lu, 2017). For nonlinear space learning methods, online Kernel-PCA has also received some research interests (Kuzmin and Warmuth, 2007; Honeine, 2012).

**Manifold learning.** This class of DR methods generally belongs to nonlinear DR techniques. Manifold learning assumes that input data lie on an embedded non-linear manifold within the high-dimensional space. DR by manifold learning aims to find a low-dimensional representation by preserving some properties of the manifold. For example, some methods preserving global properties include Multi-dimensional scaling (MDS), and IsoMap (Tenenbaum et al., 2000), while some preserving local properties including Locally Linear Embedding (LLE) (Roweis and Saul, 2000) and Laplacian Eigenmaps.

For online manifold learning settings, there are some efforts for achieving the incremental approaches of manifold learning in literature. For example, Law and Jain (2006) proposed an incremental learning algorithm for ISOMAP and Schuon et al. (2008) presented an online approach for LLE.

#### 8.4 Online Density Estimation

Online density estimation refers to constructing an estimate of an underlying unobservable probability density function based on observed data streams (Silverman, 2018). In literature, there are many different approaches to perform density estimation, e.g., histograms, naive estimator, nearest neighbour methods, Parzen windows, etc. Among various approaches, Kernel Density Estimation (KDE) is probably one of the most extensively explored topics in density estimation, which is a non-parametric way to estimate the probability density function of a target random variable (Scott, 2015). Here, we briefly review and categorize some of the commonly used approaches for kernel density estimation in online-learning settings. Given a sequence of instances  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ , where  $\mathbf{x}_t \in \mathbb{R}^d$ , KDE estimates the density at a point  $\mathbf{x}$  as

$$\mathbf{f}(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T \kappa(\mathbf{x}, \mathbf{x}_t) = \frac{1}{Th} \sum_{t=1}^T \kappa\left(\frac{\mathbf{x} - \mathbf{x}_t}{h}\right)$$

where the kernel  $\kappa(\mathbf{x}, \mathbf{x}_t)$  is a radially symmetric unimodal function that integrates to 1 and  $h$  is a smoothing parameter called the bandwidth. Like in the case of Online Learning with Kernels (Kivinen et al., 2004; Lu et al., 2015), this problem suffers from the curse of kernelization, which means to estimate the density at any point  $\mathbf{x}$ , it requires computing the kernel function with respect to all the data points observed in the data stream so far.

There have been several attempts to overcome this curse of kernelization, and can be grouped into *Merging* and *Sampling* approaches. Merging approaches require a pre-specified budget on how many instances or kernels can be stored in memory. A newly arriving sample will (typically) be stored in memory as a kernel, unless the budget is exceeded. If the budget is exceeded, two or more similar kernels get merged. The merging criteria depends on some objective function. Some efforts in this direction include (Zhou et al., 2003; Boedihardjo et al., 2008; Kristan et al., 2011), which usually differ in how they select the bandwidth values. Another approach in (Cao et al., 2012) performs clustering using self-organizing maps, and leverages this to perform kernel merging. Sampling approaches randomly select points to be kept in memory, but attempts to maintain a certain level of accuracy (Zheng et al., 2013). Recent approaches (Qahtan et al., 2017) try to perform online density estimation with efficient methods for bandwidth selection and also to capture changes in the data distribution. Finally, online KDE techniques can be applied and integrated with real-world applications, such as real-time visual tracking (Han et al., 2008).

## 8.5 Online Anomaly Detection

Anomaly Detection (AD), also known as “outlier detection” or “novelty detection”, is the process of detecting abnormal behavior in the data. The definition of abnormal behaviors can be very subjective, and the notion of “anomaly” varies from domain to domain. Anomaly detection research is abundant in literature due to its wide applications. Example applications include but not limited to intrusion detection, fraud detection, medical anomaly detection, industrial damage detection, amongst others. Anomaly detection has been extensively studied by many communities in a wide range of diverse settings, ranging from supervised to unsupervised and semi-supervised learning, and batch learning to online learning settings. More comprehensive surveys of classic anomaly detection studies can be found in (Chandola et al., 2009; Gupta et al., 2014). In this survey, we will focus online anomaly detection in unsupervised learning settings, which we believe it is one of the most popular and dominating scenarios in many real-world applications.

According to the literature surveys, unsupervised anomaly detection can be grouped into several major categories, including Distance based, Density based, Clustering based, Statistical methods, and others (such as subspace and one-class learning, etc). In the following, we briefly review some of popular work by focusing on online learning settings.

In literature, distance based online AD algorithms have been extensively studied in the context of unsupervised learning over data streams using distance-based methods (Angiulli and Fassetti, 2007; Yang et al., 2009a; Bu et al., 2009; Kloft and Laskov, 2012). Some typical strategies of distance-based online AD approaches is to apply the sliding window model where distance-based anomalies/outliers can be detected in the current window. In addition to distance-based methods, there are also some other studies that explore different methods for online AD in data streams, such as using one-class anomaly detector (Tan et al., 2011) or online clustering based approaches (Spinosa et al., 2009). We note that, despite extensive and diverse studies in the field of online anomaly detection, from a machine learning perspective, many of these approaches (e.g., sliding windows based) not purely learn in an online-learning fashion and many are not designed in machine learning based manners. We therefore keep the review of this part brief and concise.

## 9. Related Areas and Other Terminologies

### 9.1 Overview

In this section, we discuss the relationship of online learning with other related areas and terminologies which sometimes may be confused. We note that some of the following remarks may be somewhat subjective, and their meanings may vary in diverse contexts whereas some terms and notions may be used interchangeably.

### 9.2 Incremental Learning

Incremental learning, or decremental learning, represents a family of machine learning techniques (Michalski et al., 1986; Poggio, 2001; Read et al., 2012), which are particularly suitable for learning from data streams. There are various definitions of incremental learning/decremental learning. The basic idea of incremental learning is to learn some models from a stream of training instances with limited space and computational costs, often attempting to approximate a traditional batch machine learning counterpart as much as possible. For example, incremental SVM (Poggio, 2001) aims to train an SVM classifier the same as a batch SVM in an incremental manner where one training instance is added for updating the model each time (and similarly a training instance can be removed by updating the model decrementally).

Incremental learning can work either in online learning or batch learning manners (Read et al., 2012). For the incremental online learning (Poggio, 2001), only one example is presented for updating the model at one time, while for the incremental batch learning (Wang et al., 2003), a batch of multiple training examples are used for updating the model each time. Incremental learning (or decremental learning) methods are often natural extensions of existing supervised learning or unsupervised learning techniques for addressing efficiency and scalability when dealing with real-world data particularly arriving in stream-based settings. Generally speaking, incremental learning can be viewed as a branch of online learning and extensions for adapting traditional offline learning counterparts in data-stream settings.

### 9.3 Sequential Learning

Sequential learning is mainly concerned with learning from sequential training data (Dietterich, 2002), formulated as follows: a learner trains a model from a collection of  $N$  training data pairs  $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}), i = 1, \dots, N\}$  where  $\mathbf{x}^{(i)} = (x_1^i, x_2^i, \dots, x_{N_i}^i)$  is an  $N_i$ -dimensional instance vector and  $\mathbf{y}^{(i)} = (y_1^i, y_2^i, \dots, y_{N_i}^i)$  is an  $N_i$ -dimensional label vector. It can be viewed as a special type of supervised learning, known as structured prediction or structured (output) learning (BakIr, 2007), where the goal is to predict structured objects (e.g., sequence or graphs), rather than simple scalar discrete (“classification”) or real values (“regression”). Unlike traditional supervised learning that often assume data is independently and identically distributed, sequential learning attempts to exploit significant sequential correlation of sequential data when training the predictive models. Some classical methods of sequential learning include sliding window methods, recurrent sliding windows, hidden Markov models, conditional random fields, and graph transformer networks, etc. There are also many recent studies for structured prediction with application to sequential learning (BakIr, 2007; Roth et al., 2009). In general, sequential learning can be solved by either batch or online learning

algorithms. Finally, it is worth mentioning another closely related learning, i.e., “sequence classification”, whose goal is to predict a single class output for a whole input “sequence” instance. Sequence classification is a special case of sequential learning with the target class vector reduced to a single variable. It is generally simpler than regular sequential learning, and can be solved by either batch or online learning algorithms.

#### 9.4 Stochastic Learning

Stochastic learning refers to a family of machine learning algorithms by following the theory and principles of stochastic optimization (Bottou, 2004; Zhang, 2004; Bottou, 2010), which have achieved great successes for solving large-scale machine learning tasks in practice (Bousquet and Bottou, 2008). Stochastic learning is closely related to online learning. Typically, stochastic learning algorithms are motivated to accelerate the training speed of some existing batch machine learning methods for large-scale machine learning tasks, which may be often solved by batch gradient descent algorithms. Stochastic learning algorithms, e.g., Stochastic Gradient Descent (SGD) or a.k.a Online Gradient Descent (OGD) in online learning terminology, often operate sequentially by processing one training instance (randomly chosen) each time in an online learning manner, which thus are computationally more efficient and scalable than the batch GD algorithms for large-scale applications. Rather than processing a single training instance each time, a more commonly used stochastic learning technique in practice is the mini-batch SGD algorithm (Bousquet and Bottou, 2008; Shalev-Shwartz et al., 2011), which processes a small batch of training instances each time. Thus, stochastic learning can be viewed as a special family of online learning algorithms and extensions, while online learning may explore more other topics and challenges beyond stochastic learning/optimizations.

#### 9.5 Adaptive Learning

This term is occasionally used in the machine learning and neural networks fields. There is no a very formal definition about what exactly is adaptive learning in literature. In literature, there are quite a lot of different studies more or less concerned with adaptive learning (Carpenter et al., 1991, 1992), which attempt to adapt a learning model/system (e.g., neural networks) for dynamically changing environments over time. In general, these existing works are similar to online learning in that the environment is often changing and evolving dynamically. But they are different in that they are not necessarily purely based on online learning theory and algorithms. Some of these works are based on heuristic adaptation/modification of existing batch learning algorithms for updating the models with respect to the environment changes. Last but not least, most of these existing works are motivated by different kinds of heuristics, generally lack solid theoretical analysis and thus can seldom give performance guarantee in theory.

#### 9.6 Interactive Learning

Traditional machine learning mostly works in a fully automated process where training data are collected and prepared typically with the aid of domain experts. By contrast, interactive (machine) learning aims to make the machine learning procedure interactive by

engaging human (users or domain experts) in the loop (Ware et al., 2001; Johnson et al., 2003). The advantages of interactive learning include the natural integration of domain knowledge in the learning process, effective communication and continuous improvements for learning efficacy through the interaction between learning systems and users/experts. Online learning often plays an important role in an interactive learning system, in which active (online) learning can be used in finding the most informative instances to save labeling costs, incremental (online) learning algorithms could be applied for updating the models sequentially, and/or bandit online learning algorithms may be used for decision-making to trade off exploration and exploitation in some scenarios.

## 9.7 Reinforcement Learning

Reinforcement Learning (RL) (Kaelbling et al., 1996; Sutton and Barto, 1998) is a branch of machine learning inspired by behaviorist psychology, which is often concerned with how software agents should take actions in an environment to maximize cumulative rewards. The goal of an agent in RL is to find a good policy and state-update function by attempting to maximize the the expected sum of discounted rewards. RL is different from supervised learning (BARTO and DIETTERICH, 2004) in that the goal of supervised learning is to reconstruct an unknown function  $f$  that can assign the desired output values  $y$  to input data  $x$ ; while the goal of RL is to find the input (policy/action)  $x$  that gives the maximum reward  $R(x)$ . In general, RL can work either in batch or online learning manners. In practice, RL methods are commonly applied to problems involving sequential dynamics and optimization of some objectives, typically with online exploration of the effects of actions. RL is closely related to bandit online learning with the similar goal of finding a good policy that has to balance the tradeoff between exploration (of uncertainty) and exploitation (of known knowledge). Many RL solutions follow the same ideas of multi-armed bandits, and some bandit algorithms were also inspired by the field of RL studies too. However, RL can be more general when learning to interact with more complex scenarios and environments.

## 9.8 Continual Learning

*Continual Learning*, also called “Lifelong Learning” (Ruvolo and Eaton, 2013; Silver et al., 2013; Parisi et al., 2018) is a field of machine learning inspired by human ability to learn new tasks throughout their lifespan. When new tasks arrive, humans are able to leverage existing knowledge, and more effectively learn the new tasks, and at the same time, they do not forget how to perform the old tasks. In formal settings, the tasks arrive sequentially, but instances for each task arrive as a batch, and thus each task is still learned in batch settings. While older methods used linear models for lifelong learning (Ruvolo and Eaton, 2013), recent efforts have been focused on continual learning with neural networks (Parisi et al., 2018), in which one of key challenges is to address the *catastrophic forgetting*, a problem which traditional machine learning including neural networks is often susceptible to, but humans are immune to. When new tasks are learned, traditional machine learning tends to forget how to perform older tasks, and a major research direction in continual learning is to develop algorithms that can address this catastrophic forgetting. Although continual learning is closely related to online learning, most existing studies still follow the paradigm of batch training, which are not considered as online learning algorithms.

## 10. Conclusions

### 10.1 Concluding Remarks

This paper gave a comprehensive survey of existing online learning works and reviewed ongoing trends of online learning research. In theory, online learning methodologies are founded primarily based on learning theory, optimization theory, and game theory. According to the type of feedback to the learner, the existing online learning methods can be roughly grouped into the following three major categories:

- **Supervised online learning** is concerned with the online learning tasks where full feedback information is always revealed to the learner, which can be further divided into three groups: (i) “linear online learning” that aims to learn a linear predictive model, (ii) “nonlinear online learning” that aims to learn a nonlinear predictive model, and (iii) non-traditional online learning that addresses a variety of supervised online learning tasks which are different from traditional supervised prediction models for classification and regression.
- **Online learning with limited feedback** is concerned with the online learning tasks where the online learner receives partial feedback information from the environment during the learning process. The learner often has to make online predictions or decisions by achieving a tradeoff between the exploitation of disclosed knowledge and the exploration of unknown information.
- **Unsupervised online learning** is concerned with the online learning tasks where the online learner only receives the sequence of data instances without any additional feedback (e.g., true class label) during the online learning tasks. Examples of unsupervised online learning include online clustering, online representation learning, and online anomaly detection tasks, etc.

In this survey, we have focused more on the first category of work because it has received more research attention than the other two categories in literature. This is mainly because supervised online learning is a natural extension of traditional batch supervised learning, and thus an online supervised learning technique could be directly applied to a wide range of real-world applications especially for real-time machine learning from data streams where conventional batch supervised learning techniques may suffer critical limitations. However, we do note that in contrast to supervised online learning, the problems of online learning with limited feedback or unsupervised online learning are generally much more challenging, and thus should attract more research attentions and efforts in the future.

### 10.2 Future Directions

Despite the extensive studies in literature, there are still many open issues and challenges, which have not been fully solved by the existing works and need to be further explored by the community efforts in the future work. In the following, we highlight a few important and emerging research directions for researchers who are interested in online learning.

First of all, although supervised online learning has been extensively studied, learning from non-stationary data streams remain an open challenge. In particular, one critical



challenge with supervised online learning is to address “concept drifting” where the target concepts to be predicted may change over time in unforeseeable ways. Although many online learning studies have attempted to address concept drifting by a variety of approaches, they are fairly limited in that they often make some restricted assumptions for addressing certain types of concept drifting patterns. In general, there still lacks of formal theoretical frameworks or principled ways for resolving all types of concept drifting issues, particularly for non-stationary settings where target concepts may drift over time in arbitrary ways.

Second, an important growing trend of online learning research is to explore large-scale online learning for real-time big data analytics. Although online learning has huge advantages over batch learning in efficiency and scalability, it remains a non-trivial task when dealing with real-world big data analytics with extremely high volume and high velocity. Despite extensive research in large-scale batch machine learning, more future research efforts should address parallel online learning and distributed online learning by exploiting various computational resources, such as high-performance computing machines, cloud computing infrastructures, and perhaps low-cost IoT computing environments.

Third, another challenge of online learning is to address the “variety” in online data analytics tasks. Most existing online learning studies are often focused on handling single-source structured data typically by vector space representations. In many real-world data analytics applications, data may come from multiple diverse sources and could contain different types of data (such as structured, semi-structured, and unstructured data). Some existing studies, such as the series of online multiple kernel learning works, have attempted to address some of these issues, but certainly have not yet fully resolved all the challenges of variety. In the future, more research efforts should address the “variety” challenges, such as multi-source online learning, multi-modal online learning, etc.

Fourth, existing online learning works seldom address the data “veracity” issue, that is, the quality of data, which can considerably affect the efficacy of online learning. Conventional online learning studies often implicitly assume data and feedback are given in perfect quality, which is not always true for many real-world applications, particularly for real-time data analytics tasks where data arriving on-the-fly may be contaminated with noise or may have missing values or incomplete data without applying advanced pre-processing. More future research efforts should address the data veracity issue by improving the robustness of online learning algorithms particularly when dealing with real data of poor quality.

Fifth, due to the remarkable successes and impact of deep learning techniques for various applications in recent years, another emerging and increasingly important topic is “online deep learning” (Sahoo et al., 2018), i.e., learning deep neural networks from data streams on the fly in an online fashion. Despite some preliminary research, we note there are still many research challenges in this field, e.g., how to balance the tradeoff between learning accuracy, computational efficiency, learning scalability and model complexity.

Last but not least, we believe it can be valuable to explore “online continual learning” by extending traditional continual learning methods for pure online-learning settings, which is more natural for many real-world applications where data for either existing and novel tasks are often arriving in a streaming and continuous fashion. Some research progress in online deep learning might be applied here, but the key challenge of continual online learning is to resolve the catastrophic forgetting problem across tasks during the online learning process.

## 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.
- Jacob Abernethy, Kevin Canini, John Langford, and Alex Simma. Online collaborative filtering. *University of California at Berkeley, Tech. Rep*, 2007.
- Jacob Abernethy, Elad Hazan, and Alexander Rakhlin. Competing in the dark: An efficient algorithm for bandit linear optimization. In *COLT*, pages 263–274, 2008.
- Marcel R Ackermann, Marcus Mörtens, Christoph Raupach, Kamil Swierkot, Christiane Lammersen, and Christian Sohler. Streamkm++: A clustering algorithm for data streams. *Journal of Experimental Algorithmics (JEA)*, 17:2–4, 2012.
- Alekh Agarwal, Martin J Wainwright, and John C Duchi. Distributed dual averaging in networks. In *Advances in Neural Information Processing Systems*, pages 550–558, 2010.
- Amit Agarwal, Elad Hazan, Satyen Kale, and Robert E Schapire. Algorithms for portfolio management based on the newton method. In *ICML*, pages 9–16. ACM, 2006.
- Charu C Aggarwal. A survey of stream clustering algorithms., 2013.
- Charu C Aggarwal, Jiawei Han, Jianyong Wang, and Philip S Yu. A framework for projected clustering of high dimensional data streams. In *VLDB*, 2004.
- Shmuel Agmon. The relaxation method for linear inequalities. *Canadian Journal of Mathematics*, 6(3):382–392, 1954.
- Shipra Agrawal and Navin Goyal. Analysis of thompson sampling for the multi-armed bandit problem. In *Conference on Learning Theory*, pages 39–1, 2012.
- Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear payoffs. In *International Conference on Machine Learning*, pages 127–135, 2013.
- Karhan Akcoglu, Petros Drineas, and Ming-Yang Kao. Fast universalization of investment strategies. *SIAM Journal on Computing*, 34(1):1–22, 2004.
- Susanne Albers. Online algorithms: a survey. *Mathematical Programming*, 2003.
- Muqet Ali, Christopher C Johnson, and Alex K Tang. Parallel collaborative filtering for streaming data. *University of Texas Austin, Tech. Rep*, 2011.
- Amineh Amini and W Ying. Dengris-stream: A density-grid based clustering algorithm for evolving data streams over sliding window. In *Proc. International Conference on Data Mining and Computer Engineering*, pages 206–210, 2012.
- Amineh Amini, Teh Ying Wah, and Hadi Saboohi. On density-based data streams clustering algorithms: A survey. *Journal of Computer Science and Technology*, 29(1):116–141, 2014.

- Oren Anava, Elad Hazan, Shie Mannor, and Ohad Shamir. Online learning for time series prediction. In *Conference on Learning Theory*, pages 172–184, 2013.
- Fabrizio Angiulli and Fabio Fasseti. Detecting distance-based outliers in streams of data. In *Proceedings of the sixteenth ACM conference on Conference on information and knowledge management*, pages 811–820. ACM, 2007.
- Raman Arora, Andrew Cotter, Karen Livescu, and Nathan Srebro. Stochastic optimization for pca and pls. In *Allerton Conference*, pages 861–868. Citeseer, 2012a.
- Raman Arora, Andy Cotter, and Nati Srebro. Stochastic optimization of pca with capped msg. In *Advances in Neural Information Processing Systems*, pages 1815–1823, 2013.
- Sanjeev Arora, Elad Hazan, and Satyen Kale. The multiplicative weights update method: a meta-algorithm and applications. *Theory of Computing*, 8(1):121–164, 2012b.
- Les E. Atlas, David A. Cohn, and Richard E. Ladner. Training connectionist networks with queries and selective sampling. In D. S. Touretzky, editor, *Advances in Neural Information Processing Systems 2*, pages 566–573. Morgan-Kaufmann, 1990.
- Jean-Yves Audibert, Rémi Munos, and Csaba Szepesvári. Exploration–exploitation tradeoff using variance estimates in multi-armed bandits. *Theoretical Computer Science*, 410(19):1876–1902, 2009.
- Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422, 2002.
- Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. Gambling in a rigged casino: The adversarial multi-armed bandit problem. In *focs*, page 322. IEEE, 1995.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002a.
- Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002b.
- Gökhan Bakır. *Predicting structured data*. MIT press, 2007.
- Yoram Baram, Ran El Yaniv, and Kobi Luz. Online choice of active learning algorithms. *Journal of Machine Learning Research*, 5(Mar):255–291, 2004.
- ANDREW G. BARTO and THOMAS G. DIETTERICH. Reinforcement learning and its relationship to supervised learning. *Handbook of learning and approximate dynamic programming*, 2:47, 2004.
- Mikhail Belkin, Partha Niyogi, and Vikas Sindhwani. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *Journal of machine learning research*, 7(Nov):2399–2434, 2006.

- Shai Ben-David, Eyal Kushilevitz, and Yishay Mansour. Online learning versus offline learning. *Machine Learning*, 29(1):45–63, 1997.
- Pavel Berkhin. A survey of clustering data mining techniques. In *Grouping multidimensional data*, pages 25–71. Springer, 2006.
- Donald A Berry, Robert W Chen, Alan Zame, David C Heath, and Larry A Shepp. Bandit problems with infinitely many arms. *The Annals of Statistics*, pages 2103–2116, 1997.
- Alina Beygelzimer, Francesco Orabona, and Chicheng Zhang. Efficient online bandit multiclass learning with  $\sqrt{T}$  regret. In *International Conference on Machine Learning*, 2017.
- Vasudha Bhatnagar, Sharanjit Kaur, and Sharma Chakravarthy. Clustering data streams using grid-based synopsis. *Knowledge and Information Systems*, 41(1):127–152, 2014.
- H Bhatt, Richa Singh, Mayank Vatsa, and N Ratha. Improving cross-resolution face matching using ensemble based co-transfer learning. 2014.
- Himanshu S Bhatt, Richa Singh, Mayank Vatsa, and Nalini Ratha. Matching cross-resolution face images using co-transfer learning. In *Image Processing (ICIP), 2012 19th IEEE International Conference on*, pages 1453–1456. IEEE, 2012.
- Avrim Blum. *On-line algorithms in machine learning*. Springer, 1998.
- Arnold P Boediardjo, Chang-Tien Lu, and Feng Chen. A framework for estimating complex probability density structures in data streams. In *Proceedings of the 17th ACM conference on Information and knowledge management*, pages 619–628. ACM, 2008.
- Allan Borodin, Ran El-Yaniv, and Vincent Gogan. Can we learn to beat the best stock. In *Advances in Neural Information Processing Systems*, pages 345–352, 2004.
- Léon Bottou. Online algorithms and stochastic approximations. In David Saad, editor, *Online Learning and Neural Networks*. Cambridge University Press, Cambridge, UK, 1998a. revised, oct 2012.
- Léon Bottou. Online learning and stochastic approximations. *On-line learning in neural networks*, 17(9):142, 1998b.
- Léon Bottou. Stochastic learning. In *Advanced lectures on machine learning*, pages 146–168. Springer, 2004.
- Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT’2010*, pages 177–186. Springer, 2010.
- Olivier Bousquet and Léon Bottou. The tradeoffs of large scale learning. In *Advances in neural information processing systems*, pages 161–168, 2008.
- Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning*, 3(1):1–122, 2011.

- Yingyi Bu, Lei Chen, Ada Wai-Chee Fu, and Dawei Liu. Efficient anomaly monitoring over moving object trajectory streams. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 159–168. ACM, 2009.
- Sébastien Bubeck and Nicolo Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1): 1–122, 2012.
- Sébastien Bubeck, Rémi Munos, Gilles Stoltz, and Csaba Szepesvári. X-armed bandits. *Journal of Machine Learning Research*, 12(May):1655–1695, 2011.
- Sébastien Bubeck, Nicolo Cesa-Bianchi, Sham M Kakade, et al. Towards minimax policies for online linear optimization with bandit feedback. In *COLT*, volume 23, 2012.
- Christopher JC Burges et al. Dimension reduction: A guided tour. *Foundations and Trends® in Machine Learning*, 2(4):275–365, 2010.
- Feng Cao, Martin Ester, Weining Qian, and Aoying Zhou. Density-based clustering over an evolving data stream with noise. In *SDM*, volume 6, pages 328–339. SIAM, 2006.
- Yuan Cao, Haibo He, and Hong Man. Somke: Kernel density estimation over data streams by sequences of self-organizing maps. *IEEE transactions on neural networks and learning systems*, 23(8):1254–1268, 2012.
- Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. Learning to rank: from pairwise approach to listwise approach. In *Proceedings of the 24th international conference on Machine learning*, pages 129–136. ACM, 2007.
- Gail A Carpenter, Stephen Grossberg, and John H Reynolds. Artmap: Supervised real-time learning and classification of nonstationary data by a self-organizing neural network. *Neural networks*, 4(5):565–588, 1991.
- Gail A Carpenter, Stephen Grossberg, Natalya Markuzon, John H Reynolds, and David B Rosen. Fuzzy artmap: A neural network architecture for incremental supervised learning of analog multidimensional maps. *Neural Networks, IEEE Transactions on*, 3(5):698–713, 1992.
- Rich Caruana. Multitask learning. In *Learning to learn*, pages 95–133. Springer, 1998.
- Giovanni Cavallanti, Nicolò Cesa-Bianchi, and Claudio Gentile. Tracking the best hyper-plane with a simple budget perceptron. *Machine Learning*, 69(2-3):143–167, 2007.
- Giovanni Cavallanti, Nicolo Cesa-Bianchi, and Claudio Gentile. Linear algorithms for online multitask classification. *Journal of Machine Learning Research*, 11(Oct):2901–2934, 2010.
- Nicolò Cesa-Bianchi and Claudio Gentile. Improved risk tail bounds for on-line algorithms. *IEEE Transactions on Information Theory*, 54(1):386–390, 2008.
- Nicolò Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge University Press, New York, NY, USA, 2006. ISBN 978-0-521-84108-5.

- Nicolo Cesa-Bianchi and Gábor Lugosi. Combinatorial bandits. *Journal of Computer and System Sciences*, 78(5):1404–1422, 2012.
- Nicolò Cesa-Bianchi and Ohad Shamir. Efficient transductive online learning via randomized rounding. In *Empirical Inference*, pages 177–194. Springer, 2013.
- Nicolò Cesa-Bianchi, Alex Conconi, and Claudio Gentile. Learning probabilistic linear-threshold classifiers via selective sampling. In *Computational Learning Theory and Kernel Machines, 16th Annual Conference on Computational Learning Theory and 7th Kernel Workshop, COLT/Kernel 2003*, pages 373–387, 2003.
- Nicolò Cesa-Bianchi, Alex Conconi, and Claudio Gentile. On the generalization ability of on-line learning algorithms. *IEEE Transactions on Information Theory*, 50(9):2050–2057, 2004.
- Nicolò Cesa-Bianchi, Alex Conconi, and Claudio Gentile. A second-order perceptron algorithm. *SIAM Journal on Computing*, 34(3):640–668, 2005a.
- Nicolo Cesa-Bianchi, Gábor Lugosi, and Gilles Stoltz. Minimizing regret with label efficient prediction. *IEEE Transactions on Information Theory*, 51(6):2152–2162, 2005b.
- Nicolò Cesa-Bianchi, Claudio Gentile, and Luca Zaniboni. Worst-case analysis of selective sampling for linear classification. *Journal of Machine Learning Research*, 7:1205–1230, 2006.
- Nicolò Cesa-Bianchi, Claudio Gentile, and Francesco Orabona. Robust bounds for classification via selective sampling. In *Proceedings of the 26th Annual International Conference on Machine Learning, ICML2009*, pages 121–128, 2009.
- Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM computing surveys (CSUR)*, 41(3):15, 2009.
- Yin-Wen Chang, Cho-Jui Hsieh, Kai-Wei Chang, Michael Ringgaard, and Chih-Jen Lin. Training and testing low-degree polynomial data mappings via linear svm. *The Journal of Machine Learning Research*, 11:1471–1490, 2010.
- Olivier Chapelle and S Sathiya Keerthi. Efficient algorithms for ranking with svms. *Information Retrieval*, 13(3):201–215, 2010.
- Olivier Chapelle and Lihong Li. An empirical evaluation of thompson sampling. In *Advances in neural information processing systems*, pages 2249–2257, 2011.
- Chris Chatfield. *Time-series forecasting*. CRC Press, 2000.
- Kamalika Chaudhuri, Yoav Freund, and Daniel J Hsu. A parameter-free hedging algorithm. In *Advances in neural information processing systems*, pages 297–305, 2009.
- Guangyun Chen, Gang Chen, Jianwen Zhang, Shuo Chen, and Changshui Zhang. Beyond banditron: A conservative and efficient reduction for online multiclass prediction with bandit setting model. In *9th IEEE International Conference on Data Mining (ICDM2009)*, pages 71–80, 2009.

- Ning Chen, Steven CH Hoi, Shaohua Li, and Xiaokui Xiao. Simapp: A framework for detecting similar mobile applications by online kernel learning. In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining*, pages 305–314. ACM, 2015.
- Ning Chen, Steven CH Hoi, Shaohua Li, and Xiaokui Xiao. Mobile app tagging. In *Proceedings of the Ninth ACM International Conference on Web Search and Data Mining*, pages 63–72. ACM, 2016a.
- Wei Chen, Yajun Wang, Yang Yuan, and Qinshi Wang. Combinatorial multi-armed bandit and its extension to probabilistically triggered arms. *Journal of Machine Learning Research*, 17(50):1–33, 2016b.
- Yixin Chen and Li Tu. Density-based clustering for real-time stream data. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 133–142. ACM, 2007.
- Zhong Chen, Zhide Fang, Wei Fan, Andrea Edwards, and Kun Zhang. Cstg: An effective framework for cost-sensitive sparse online learning. In *Proceedings of the 2017 SIAM International Conference on Data Mining*, pages 759–767. SIAM, 2017.
- Alexey Chernov and Vladimir Vovk. Prediction with advice of unknown number of experts. In *Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*, pages 117–125. AUAI Press, 2010.
- Sayak Ray Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *International Conference on Machine Learning*, 2017.
- Wei Chu, Lihong Li, Lev Reyzin, and Robert E Schapire. Contextual bandits with linear payoff functions. In *AISTATS*, volume 15, pages 208–214, 2011.
- Michael Clements and David Hendry. *Forecasting economic time series*. Cambridge University Press, 1998.
- Richard Combes, Mohammad Sadegh Talebi Mazraeh Shahi, Alexandre Proutiere, et al. Combinatorial bandits revisited. In *Advances in Neural Information Processing Systems*, pages 2116–2124, 2015.
- Thomas M Cover. Universal portfolios. In *The Kelly Capital Growth Investment Criterion: Theory and Practice*, pages 181–209. World Scientific, 2011.
- Koby Crammer and Claudio Gentile. Multiclass classification with bandit feedback using adaptive regularization. In *Proceedings of 28th International Conference on Machine Learning (ICML2011)*, pages 273–280, 2011.
- Koby Crammer and Daniel D Lee. Learning via gaussian herding. In *Advances in neural information processing systems*, pages 451–459, 2010.
- Koby Crammer and Yoram Singer. Online ranking by projecting. *Neural Computation*, 17(1):145–175, 2005.

- Koby Crammer, Yoram Singer, et al. Pranking with ranking. In *Nips*, volume 1, pages 641–647, 2001.
- Koby Crammer, Jaz S Kandola, and Yoram Singer. Online classification on a budget. In *NIPS*, volume 2, page 5, 2003.
- Koby Crammer, Ofer Dekel, Joseph Keshet, Shai Shalev-Shwartz, and Yoram Singer. Online passive-aggressive algorithms. *The Journal of Machine Learning Research*, 7:551–585, 2006.
- Koby Crammer, Mark Dredze, and Alex Kulesza. Multi-class confidence weighted algorithms. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 496–504, 2009a.
- Koby Crammer, Alex Kulesza, and Mark Dredze. Adaptive regularization of weight vectors. *Machine Learning*, pages 1–33, 2009b.
- Ashok Cutkosky and Kwabena A Boahen. Online convex optimization with unconstrained domains and losses. In *Advances In Neural Information Processing Systems*, pages 748–756, 2016.
- Abhinandan S Das, Mayur Datar, Ashutosh Garg, and Shyam Rajaram. Google news personalization: scalable online collaborative filtering. In *Proceedings of the 16th international conference on World Wide Web*, pages 271–280. ACM, 2007.
- Jason V Davis, Brian Kulis, Prateek Jain, Suvrit Sra, and Inderjit S Dhillon. Information-theoretic metric learning. In *Proceedings of the 24th international conference on Machine learning*, pages 209–216. ACM, 2007.
- Ofer Dekel, Shai Shalev-Shwartz, and Yoram Singer. The forgetron: A kernel-based perceptron on a fixed budget. In *NIPS*, 2005.
- Ofer Dekel, Philip M Long, and Yoram Singer. Online multitask learning. In *International Conference on Computational Learning Theory*, pages 453–467. Springer, 2006.
- Ofer Dekel, Claudio Gentile, and Karthik Sridharan. Robust selective sampling from single and multiple teachers. In *COLT 2010 - The 23rd Conference on Learning Theory, Haifa, Israel, June 27-29, 2010*, pages 346–358, 2010.
- Ofer Dekel, Ran Gilad-Bachrach, Ohad Shamir, and Lin Xiao. Optimal distributed online prediction using mini-batches. *The Journal of Machine Learning Research*, 13(1):165–202, 2012.
- Thomas G Dietterich. Machine learning for sequential data: A review. In *Structural, syntactic, and statistical pattern recognition*, pages 15–30. Springer, 2002.
- Yi Ding, Peilin Zhao, Steven CH Hoi, and Yew-Soon Ong. An adaptive gradient method for online auc maximization. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.



- Mark Dredze and Koby Crammer. Active learning with confidence. In *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics on Human Language Technologies: Short Papers*, pages 233–236, 2008.
- Mark Dredze, Koby Crammer, and Fernando Pereira. Confidence-weighted linear classification. In *Proceedings of the 25th international conference on Machine learning*, pages 264–271. ACM, 2008.
- John Duchi and Yoram Singer. Efficient online and batch learning using forward backward splitting. *The Journal of Machine Learning Research*, 10:2899–2934, 2009.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *The Journal of Machine Learning Research*, 12:2121–2159, 2011a.
- John C. Duchi, Shai Shalev-Shwartz, Yoram Singer, and Ambuj Tewari. Composite objective mirror descent. In *COLT*, pages 14–26, 2010.
- John C. Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011b.
- Ryan Elwell and Robi Polikar. Incremental learning of concept drift in nonstationary environments. *Neural Networks, IEEE Transactions on*, 22(10):1517–1531, 2011.
- Theodoros Evgeniou and Massimiliano Pontil. Regularized multi-task learning. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 109–117. ACM, 2004.
- Jason Farquhar, David Hardoon, Hongying Meng, John S Shawe-taylor, and Sandor Szepesvari. Two view learning: Svm-2k, theory and practice. In *Advances in neural information processing systems*, pages 355–362, 2006.
- Jiashi Feng, Huan Xu, Shie Mannor, and Shuicheng Yan. Online pca for contaminated data. In *Advances in Neural Information Processing Systems*, pages 764–772, 2013.
- Amos Fiat and Gerhard Woeginger. *Online algorithms: The state of the art*. Springer Heidelberg, 1998.
- Douglas H Fisher. Knowledge acquisition via incremental conceptual clustering. *Machine learning*, 2(2):139–172, 1987.
- Yoav Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *J. Comput. Syst. Sci.*, 55(1):119–139, 1997.
- Yoav Freund and Robert E. Schapire. Large margin classification using the perceptron algorithm. *Mach. Learn.*, 37(3):277–296, 1999a.
- Yoav Freund and Robert E Schapire. Adaptive game playing using multiplicative weights. *Games and Economic Behavior*, 29(1-2):79–103, 1999b.

- Yoav Freund, H Sebastian Seung, Eli Shamir, and Naftali Tishby. Selective sampling using the query by committee algorithm. *Machine learning*, 28(2-3):133–168, 1997.
- Daniel Gabay and Bertrand Mercier. A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Computers & Mathematics with Applications*, 2(1):17–40, 1976.
- Alexei A Gaivoronski and Fabio Stella. Stochastic nonstationary optimization for finding universal portfolios. *Annals of Operations Research*, 100(1):165–188, 2000.
- Jing Gao, Jianzhong Li, Zhaogong Zhang, and Pang-Ning Tan. An incremental data stream clustering algorithm based on dense units detection. In *Advances in Knowledge Discovery and Data Mining*, pages 420–425. Springer, 2005.
- Wei Gao, Rong Jin, Shenghuo Zhu, and Zhi-Hua Zhou. One-pass auc optimization. In *ICML*, 2013.
- Xingyu Gao, Steven CH Hoi, Yongdong Zhang, Ji Wan, and Jintao Li. Soml: Sparse online metric learning with application to image retrieval. *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence*, 2014.
- Xingyu Gao, Steven CH Hoi, Yongdong Zhang, Jianshe Zhou, Ji Wan, Zhenyu Chen, Jintao Li, and Jianke Zhu. Sparse online learning of image similarity. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 8(5):64, 2017.
- Liang Ge, Jing Gao, and Aidong Zhang. Oms-tl: a framework of online multiple source transfer learning. In *Proceedings of the 22nd ACM international conference on Conference on information & knowledge management*, pages 2423–2428. ACM, 2013.
- Liang Ge, Jing Gao, Hung Ngo, Kang Li, and Aidong Zhang. On handling negative transfer and imbalanced distributions in multiple source transfer learning. *Statistical Analysis and Data Mining: The ASA Data Science Journal*, 7(4):254–271, 2014.
- Claudio Gentile. A new approximate maximal margin classification algorithm. *Journal of Machine Learning Research*, 2:213–242, 2001.
- Box George. *Time Series Analysis: Forecasting & Control*, 3/e. Pearson Education India, 1994.
- John Gittins, Kevin Glazebrook, and Richard Weber. *Multi-armed bandit allocation indices*. John Wiley & Sons, 2011.
- John C Gittins. Bandit processes and dynamic allocation indices. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 148–177, 1979.
- Andrew B Goldberg, Ming Li, and Xiaojin Zhu. Online manifold regularization: A new learning setting and empirical study. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 393–407. Springer, 2008.
- Andrew B Goldberg, Xiaojin Zhu, Alex Furger, and Jun-Ming Xu. Oasis: Online active semi-supervised learning. In *AAAI*, 2011.

- Sudipto Guha, Nina Mishra, Rajeev Motwani, and Liadan O’Callaghan. Clustering data streams. In *Foundations of computer science, 2000. proceedings. 41st annual symposium on*, pages 359–366. IEEE, 2000.
- Sudipto Guha, Adam Meyerson, Nina Mishra, Rajeev Motwani, and Liadan O’Callaghan. Clustering data streams: Theory and practice. *Knowledge and Data Engineering, IEEE Transactions on*, 15(3):515–528, 2003.
- Manish Gupta, Jing Gao, Charu C Aggarwal, and Jiawei Han. Outlier detection for temporal data: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 26(9):2250–2267, 2014.
- L Györfi and D Schafer. Nonparametric prediction. *NATO SCIENCE SERIES SUB SERIES III COMPUTER AND SYSTEMS SCIENCES*, 190:341–356, 2003.
- László Györfi, Frederic Udina, and Harro Walk. Nonparametric nearest neighbor based empirical portfolio selection strategies. *Statistics & Decisions International mathematical journal for stochastic methods and models*, 26(2):145–157, 2008.
- Bohyung Han, Dorin Comaniciu, Ying Zhu, and Larry S Davis. Sequential kernel density approximation and its application to real-time visual tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(7):1186–1197, 2008.
- LI Hang. A short introduction to learning to rank. *IEICE TRANSACTIONS on Information and Systems*, 94(10):1854–1862, 2011.
- James Hannan. Approximation to bayes risk in repeated play. *Contributions to the Theory of Games*, 3(97-139):2, 1957.
- Shuji Hao, Steven CH Hoi, Chunyan Miao, and Peilin Zhao. Active crowdsourcing for annotation. In *Web Intelligence and Intelligent Agent Technology (WI-IAT), 2015 IEEE/WIC/ACM International Conference on*, volume 2, pages 1–8. IEEE, 2015a.
- Shuji Hao, Peilin Zhao, Steven CH Hoi, and Chunyan Miao. Learning relative similarity from data streams: Active online learning approaches. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, pages 1181–1190. ACM, 2015b.
- Shuji Hao, Peilin Zhao, Jing Lu, Steven CH Hoi, Chunyan Miao, and Chi Zhang. Soal: Second-order online active learning. In *Data Mining (ICDM), 2016 IEEE 16th International Conference on*, pages 931–936. IEEE, 2016.
- Shuji Hao, Jing Lu, Peilin Zhao, Chi Zhang, Steven CH Hoi, and Chunyan Miao. Second-order online active learning and its applications. *IEEE Transactions on Knowledge and Data Engineering*, 2017a.
- Shuji Hao, Peilin Zhao, Yong Liu, Steven CH Hoi, and Chunyan Miao. Online multitask relative similarity learning. *International Joint Conference on Artificial Intelligence*, 2017b.

- Edward F Harrington. Online ranking/collaborative filtering using the perceptron algorithm. In *ICML*, volume 20, pages 250–257, 2003.
- Elad Hazan and Satyen Kale. Newtron: an efficient bandit algorithm for online multiclass prediction. In *Advances in Neural Information Processing Systems*, pages 891–899, 2011.
- Elad Hazan and Comandur Seshadhri. Efficient learning algorithms for changing environments. In *Proceedings of the 26th annual international conference on machine learning*, pages 393–400. ACM, 2009.
- Elad Hazan, Amit Agarwal, and Satyen Kale. Logarithmic regret algorithms for online convex optimization. *Machine Learning*, 69(2-3):169–192, 2007a.
- Elad Hazan, Alexander Rakhlin, and Peter L Bartlett. Adaptive online gradient descent. In *Advances in Neural Information Processing Systems*, pages 65–72, 2007b.
- Elad Hazan et al. Introduction to online convex optimization. *Foundations and Trends® in Optimization*, 2(3-4):157–325, 2016.
- Reinhard Heckel and Kannan Ramchandran. The sample complexity of online one-class collaborative filtering. In *International Conference on Machine Learning*, 2017.
- David Helmbold and Sandra Panizza. Some label efficient learning results. In *Proceedings of the tenth annual conference on Computational learning theory*, pages 218–230. ACM, 1997.
- David P Helmbold, Robert E Schapire, Yoram Singer, and Manfred K Warmuth. On-line portfolio selection using multiplicative updates. *Mathematical Finance*, 8(4):325–347, 1998.
- Ralf Herbrich, Thore Graepel, and Klaus Obermayer. Support vector learning for ordinal regression. 1999.
- Steven C. Hoi, Rong Jin, Peilin Zhao, and Tianbao Yang. Online multiple kernel classification. *Machine Learning*, 90(2):289–316, 2013.
- Steven CH Hoi, Jiale Wang, Peilin Zhao, and Rong Jin. Online feature selection for mining big data. In *Proceedings of the 1st International Workshop on Big Data, Streams and Heterogeneous Source Mining: Algorithms, Systems, Programming Models and Applications*, pages 93–100. ACM, 2012.
- Steven CH Hoi, Jiale Wang, and Peilin Zhao. Libol: a library for online learning algorithms. *The Journal of Machine Learning Research*, 15(1):495–499, 2014.
- Paul Honeine. Online kernel principal component analysis: A reduced-order model. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (9):1814–1826, 2012.
- Junjie Hu, Haiqin Yang, Irwin King, Michael R Lyu, and Anthony Man-Cho So. Kernelized online imbalanced learning with fixed budgets. In *AAAI*, pages 2666–2672, 2015.

- Dingjiang Huang, Junlong Zhou, Bin Li, Steven CH Hoi, and Shuigeng Zhou. Robust median reversion strategy for on-line portfolio selection. In *Proceedings of the Twenty-Third international joint conference on Artificial Intelligence*, pages 2006–2012. AAAI Press, 2013.
- Dingjiang Huang, Yan Zhu, Bin Li, Shuigeng Zhou, and Steven CH Hoi. Semi-universal portfolios with transaction costs. In *Twenty-Fourth International Joint Conference on Artificial Intelligence*, 2015.
- Prateek Jain, Brian Kulis, Inderjit S Dhillon, and Kristen Grauman. Online metric learning and fast similarity search. In *Advances in neural information processing systems*, pages 761–768, 2009.
- Rodolphe Jenatton, Jim Huang, and Cédric Archambeau. Adaptive algorithms for online convex optimization with long-term constraints. *NIPS*, 2016.
- Chen Jia, ChengYu Tan, and Ai Yong. A grid and density-based clustering algorithm for processing data stream. In *Genetic and Evolutionary Computing, 2008. WGEC’08. Second International Conference on*, pages 517–521. IEEE, 2008.
- Luo Jie, Francesco Orabona, Marco Fornoni, Barbara Caputo, and Nicolo Cesa-bianchi. Om-2: An online multi-class multi-kernel learning algorithm. In *Proc. of the 4th IEEE Online Learning for Computer Vision Workshop*, 2010.
- Rong Jin, Shijun Wang, and Yang Zhou. Regularized distance metric learning: Theory and algorithm. In *Advances in neural information processing systems*, pages 862–870, 2009.
- Rong Jin, Steven C. H. Hoi, and Tianbao Yang. Online multiple kernel learning: Algorithms and mistake bounds. In *Algorithmic Learning Theory, 21st International Conference, ALT 2010, Canberra, Australia, October 6-8, 2010. Proceedings*, pages 390–404, 2010.
- David Johnson, Sylvie Levesque, and Tong Zhang. Interactive machine learning system for automated annotation of information in text, July 31 2003. US Patent App. 10/630,854.
- Kwang-Sung Jun, Aniruddha Bhargava, Robert Nowak, and Rebecca Willett. Scalable generalized linear bandits: Online computation and hashing. In *Advances in Neural Information Processing Systems*, 2017.
- Leslie Pack Kaelbling, Michael L Littman, and Andrew W Moore. Reinforcement learning: A survey. *Journal of artificial intelligence research*, pages 237–285, 1996.
- Sham M Kakade and Ambuj Tewari. On the generalization ability of online strongly convex programming algorithms. In *Advances in Neural Information Processing Systems*, pages 801–808, 2009.
- Sham M. Kakade, Shai Shalev-Shwartz, and Ambuj Tewari. Efficient bandit algorithms for online multiclass prediction. In *ICML*, pages 440–447, 2008.
- Adam Tauman Kalai and Santosh Vempala. Efficient algorithms for online decision problems. *J. Comput. Syst. Sci.*, 71(3):291–307, 2005.

- Satyen Kale, Zohar Karnin, Tengyuan Liang, and Dávid Pál. Adaptive feature selection: Computationally efficient online sparse linear regression under rip. In *International Conference on Machine Learning*, 2017.
- Purushottam Kar, Bharath K Sriperumbudur, Prateek Jain, and Harish C Karnick. On the generalization ability of online learning algorithms for pairwise loss functions. In *ICML*, 2013.
- Michael N Katehakis and Arthur F Veinott Jr. The multi-armed bandit problem: decomposition and computation. *Mathematics of Operations Research*, 12(2):262–268, 1987.
- L Kaufman and Pr J Rousseeuw. Clustering large applications (program clara). *Finding groups in data: an introduction to cluster analysis*, pages 126–146, 2008.
- Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On bayesian upper confidence bounds for bandit problems. In *Artificial Intelligence and Statistics*, pages 592–600, 2012.
- John L Kelly Jr. A new interpretation of information rate. In *The Kelly Capital Growth Investment Criterion: Theory and Practice*, pages 25–34. World Scientific, 2011.
- Jyrki Kivinen and Manfred K. Warmuth. Additive versus exponentiated gradient updates for linear prediction. In *Proceedings of the Twenty-Seventh Annual ACM Symposium on Theory of Computing (STOC’95)*, pages 209–218, 1995.
- Jyrki Kivinen, Alexander J Smola, and Robert C Williamson. Online learning with kernels. *Signal Processing, IEEE Transactions on*, 52(8):2165–2176, 2004.
- Robert D Kleinberg. Nearly tight bounds for the continuum-armed bandit problem. In *Advances in Neural Information Processing Systems*, pages 697–704, 2005a.
- Robert David Kleinberg. *Online decision problems with large strategy sets*. PhD thesis, Massachusetts Institute of Technology, 2005b.
- Marius Kloft and Pavel Laskov. Security analysis of online centroid anomaly detection. *Journal of Machine Learning Research*, 13(Dec):3681–3724, 2012.
- Ari Kobren, Nicholas Monath, Akshay Krishnamurthy, and Andrew McCallum. An online hierarchical algorithm for extreme clustering. *arXiv preprint arXiv:1704.01858*, 2017.
- Wouter M Koolen and Tim Van Erven. Second-order quantile methods for experts and combinatorial games. In *Conference on Learning Theory*, pages 1155–1175, 2015.
- Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *Computer*, (8):30–37, 2009.
- Philipp Kranen, Ira Assent, Corinna Baldauf, and Thomas Seidl. The clustree: indexing micro-clusters for anytime stream mining. *Knowledge and information systems*, 29(2): 249–272, 2011.

- Werner Krauth and Marc Mézard. Learning algorithms with optimal stability in neural networks. *Journal of Physics A: Mathematical and General*, 20(11):L745, 1987.
- Matej Kristan, Ales Leonardis, and Danijel Skocaj. Multivariate online kernel density estimation with gaussian kernels. *Pattern Recognition*, 44(10-11):2630–2642, 2011.
- Abhishek Kumar and Hal Daumé III. Learning task grouping and overlap in multi-task learning. In *Proceedings of the 29th International Conference on Machine Learning*, pages 1723–1730. Omnipress, 2012.
- Dima Kuzmin and Manfred K Warmuth. Online kernel pca with entropic matrix updates. In *Proceedings of the 24th international conference on Machine learning*, pages 465–472. ACM, 2007.
- Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22, 1985.
- John Langford and Tong Zhang. The epoch-greedy algorithm for multi-armed bandits with side information. In *NIPS*, pages 817–824, 2008.
- John Langford, Lihong Li, and Alex Strehl. Vowpal wabbit online learning project, 2007.
- John Langford, Lihong Li, and Tong Zhang. Sparse online learning via truncated gradient. *The Journal of Machine Learning Research*, 10:777–801, 2009.
- Martin HC Law and Anil K Jain. Incremental nonlinear dimensionality reduction by manifold learning. *IEEE transactions on pattern analysis and machine intelligence*, 28(3):377–391, 2006.
- Trung Le, Tu Nguyen, Vu Nguyen, and Dinh Phung. Dual space gradient descent for online learning. In *Advances In Neural Information Processing Systems*, pages 4583–4591, 2016.
- Yann A LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient back-prop. In *Neural networks: Tricks of the trade*, pages 9–48. Springer, 1998.
- Kfir Y Levy. Online to offline conversions and adaptive minibatch sizes. In *Advances in Neural Information Processing Systems*, 2017.
- Bin Li. *Online Portfolio Selection*. PhD thesis, Nanyang Technological University, 2013.
- Bin Li and Steven CH Hoi. On-line portfolio selection with moving average reversion. *arXiv preprint arXiv:1206.4626*, 2012.
- Bin Li and Steven CH Hoi. Online portfolio selection: A survey. *ACM Computing Surveys (CSUR)*, 46(3):35, 2014.
- Bin Li and Steven Chu Hong Hoi. *Online Portfolio Selection: Principles and Algorithms*. Crc Press, 2015.
- Bin Li, Steven CH Hoi, and Vivekanand Gopalkrishnan. Corn: Correlation-driven non-parametric learning approach for portfolio selection. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2(3):21, 2011a.

- Bin Li, Steven CH Hoi, Peilin Zhao, and Vivekanand Gopalkrishnan. Confidence weighted mean reversion strategy for on-line portfolio selection. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages 434–442, 2011b.
- Bin Li, Peilin Zhao, Steven CH Hoi, and Vivekanand Gopalkrishnan. Pamr: Passive aggressive mean reversion strategy for portfolio selection. *Machine learning*, 87(2):221–258, 2012.
- Bin Li, Steven CH Hoi, Peilin Zhao, and Vivekanand Gopalkrishnan. Confidence weighted mean reversion strategy for online portfolio selection. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 7(1):4, 2013.
- Bin Li, Steven CH Hoi, Doyen Sahoo, and Zhi-Yong Liu. Moving average reversion strategy for on-line portfolio selection. *Artificial Intelligence*, 222:104–123, 2015.
- Bin Li, Doyen Sahoo, and Steven CH Hoi. Olps: a toolbox for on-line portfolio selection. *Journal of Machine Learning Research*, 17(35):1–5, 2016a.
- Bin Li, Jialei Wang, Dingjiang Huang, and Steven CH Hoi. Transaction cost optimization for online portfolio selection. *Quantitative Finance*, pages 1–14, 2017a.
- Chris Junchi Li, Zhaoran Wang, and Han Liu. Online ica: Understanding global dynamics of nonconvex optimization via diffusion processes. In *Advances in Neural Information Processing Systems*, pages 4967–4975, 2016b.
- Guangxia Li, Steven CH Hoi, Kuiyu Chang, and Ramesh Jain. Micro-blogging sentiment detection by collaborative online learning. In *IEEE Intl. Conference on Data Mining*, pages 893–898. IEEE, 2010a.
- Guangxia Li, Steven CH Hoi, Kuiyu Chang, Wenting Liu, and Ramesh Jain. Collaborative online multitask learning. *IEEE Transactions on Knowledge and Data Engineering*, 26(8):1866–1876, 2014.
- Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670. ACM, 2010b.
- Lihong Li, Yu Lu, and Dengyong Zhou. Provable optimal algorithms for generalized linear contextual bandits. In *International Conference on Machine Learning*, 2017b.
- Yaoyong Li, Hugo Zaragoza, Ralf Herbrich, John Shawe-Taylor, and Jaz Kandola. The perceptron algorithm with uneven margins. In *ICML*, volume 2, pages 379–386, 2002.
- Yi Li and Philip M. Long. The relaxed online maximum margin algorithm. *Machine Learning*, 46(1-3):361–387, 2002.
- Yingming Li, Ming Yang, and Zhongfei Zhang. Multi-view representation learning: A survey from shallow methods to deep methods. *arXiv preprint arXiv:1610.01206*, 2016c.



- Liu Li-xiong, Kang Jing, Guo Yun-fei, and Huang Hai. A three-step clustering algorithm over an evolving data stream. In *Intelligent Computing and Intelligent Systems, 2009. ICIS 2009. IEEE International Conference on*, volume 1, pages 160–164. IEEE, 2009.
- Nan-Ying Liang, Guang-Bin Huang, Paramasivan Saratchandran, and Narasimhan Sundararajan. A fast and accurate online sequential learning algorithm for feedforward networks. *Neural Networks, IEEE Transactions on*, 17(6):1411–1423, 2006.
- Keng-Pei Lin and Ming-Syan Chen. Efficient kernel approximation for large-scale support vector machine classification. In *Proceedings of the Eleventh SIAM International Conference on Data Mining*, pages 211–222. SIAM, 2011.
- Guang Ling, Haiqin Yang, Irwin King, and Michael R Lyu. Online learning for collaborative filtering. In *Neural Networks (IJCNN), The 2012 International Joint Conference on*, pages 1–8. IEEE, 2012.
- Nick Littlestone. Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm. *Machine learning*, 2(4):285–318, 1988.
- Nick Littlestone. From on-line to batch learning. In *Proceedings of the Second Annual Workshop on Computational Learning Theory, COLT 1989, Santa Cruz, CA, USA, July 31 - August 2, 1989.*, pages 269–284, 1989.
- Nick Littlestone and Manfred K Warmuth. The weighted majority algorithm. In *30th Annual Symposium on Foundations of Computer Science*, pages 256–261, 1989.
- Nick Littlestone and Manfred K. Warmuth. The weighted majority algorithm. *Inf. Comput.*, 108(2):212–261, 1994.
- Chenghao Liu, Steven CH Hoi, Peilin Zhao, and Jianling Sun. Online arima algorithms for time series prediction. 2016a.
- Chenghao Liu, Steven CH Hoi, Peilin Zhao, Jianling Sun, and Ee-Peng Lim. Online adaptive passive-aggressive methods for non-negative matrix factorization and its applications. In *Proceedings of the 25th ACM International Conference on Information and Knowledge Management*, pages 1161–1170. ACM, 2016b.
- Chenghao Liu, Tao Jin, Steven CH Hoi, Peilin Zhao, and Jianling Sun. Collaborative topic regression for online recommender systems: an online and bayesian approach. *Machine Learning*, 106(5):651–670, 2017.
- Nathan N Liu, Min Zhao, Evan Xiang, and Qiang Yang. Online evolutionary collaborative filtering. In *Proceedings of 4th ACM conference on Recommender systems*, pages 95–102, 2010.
- Jing Lu, Steven Hoi, and Jiale Wang. Second order online collaborative filtering. In *Asian Conference on Machine Learning*, pages 325–340, 2013.
- Jing Lu, Peilin Zhao, and Steven C.H. Hoi. Online passive aggressive active learning and its applications. *The 6th Asian Conference on Machine Learning (ACML2014)*, 2014.

- Jing Lu, Steven C.H. Hoi, Jialei Wang, Peilin Zhao, and Zhi-Yong Liu. Large scale online kernel learning. *The Journal of Machine Learning Research*, 2015.
- Jing Lu, Peilin Zhao, and Steven CH Hoi. Online sparse passive aggressive learning with kernels. In *Proceedings of the 2016 SIAM International Conference on Data Mining*, pages 675–683. SIAM, 2016a.
- Jing Lu, Peilin Zhao, and Steven CH Hoi. Online passive-aggressive active learning. *Machine Learning*, 103(2):141–183, 2016b.
- Jing Lu, Doyen Sahoo, Peilin Zhao, and Steven CH Hoi. Sparse passive-aggressive learning for bounded online kernel methods. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 9(4):45, 2018.
- Haipeng Luo and Robert E Schapire. Achieving all with no parameters: Adanormalhedge. In *Conference on Learning Theory*, pages 1286–1304, 2015.
- Haipeng Luo, Alekh Agarwal, Nicolò Cesa-Bianchi, and John Langford. Efficient second order online learning by sketching. In *Advances in Neural Information Processing Systems*, pages 902–910, 2016.
- Stefan Magureanu, Richard Combes, and Alexandre Proutière. Lipschitz bandits: Regret lower bound and optimal algorithms. In *Proceedings of The 27th Conference on Learning Theory, COLT 2014, Barcelona, Spain, June 13-15, 2014*, pages 975–999, 2014.
- André Filipe Torres Martins, Noah A. Smith, Eric P. Xing, Pedro M. Q. Aguiar, and Mário A. T. Figueiredo. Online learning of structured predictors with multiple kernels. *Journal of Machine Learning Research - Proceedings Track*, 15:507–515, 2011.
- Benedict C May, Nathan Korda, Anthony Lee, and David S Leslie. Optimistic bayesian sampling in contextual-bandit problems. *Journal of Machine Learning Research*, 13(Jun): 2069–2106, 2012.
- H Brendan McMahan and Matthew J Streeter. Tighter bounds for multi-armed bandits with expert advice. In *COLT*, 2009.
- Avihai Mejer and Koby Crammer. Confidence in structured-prediction using confidence-weighted models. In *Proceedings of the 2010 conference on empirical methods in natural language processing*, pages 971–981. Association for Computational Linguistics, 2010.
- Ryszard S Michalski, Igor Mozetic, Jiarong Hong, and Nada Lavrac. The multi-purpose incremental learning system aq15 and its testing application to three medical domains. *Proc. AAAI 1986*, pages 1–041, 1986.
- Ioannis Mitliagkas, Constantine Caramanis, and Prateek Jain. Memory limited, streaming pca. In *Advances in Neural Information Processing Systems*, pages 2886–2894, 2013.
- João FC Mota, João MF Xavier, Pedro MQ Aguiar, and Markus Püschel. D-admm: A communication-efficient distributed algorithm for separable optimization. *IEEE Transactions on Signal Processing*, 61(10):2718–2723, 2013.

- Keerthiram Murugesan, Hanxiao Liu, Jaime Carbonell, and Yiming Yang. Adaptive smoothed online multi-task learning. In *Advances in Neural Information Processing Systems*, pages 4296–4304, 2016.
- Yurii Nesterov. Primal-dual subgradient methods for convex problems. *Mathematical programming*, 120(1):221–259, 2009.
- Tam T Nguyen, Kuiyu Chang, and Siu Cheung Hui. Two-view online learning. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 74–85. Springer, 2012.
- Tu Dinh Nguyen, Trung Le, Hung Bui, and Dinh Phung. Large-scale online kernel learning with random feature reparameterization. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence (IJCAI-17)*, pages 2543–2549, 2017.
- Noam Nisan, Tim Roughgarden, Eva Tardos, and Vijay V Vazirani. *Algorithmic game theory*, volume 1. Cambridge University Press Cambridge, 2007.
- Albert B Novikoff. On convergence proofs for perceptrons. Technical report, STANFORD RESEARCH INST MENLO PARK CA, 1963.
- Irene Ntoutsis, Arthur Zimek, Themis Palpanas, Peer Kröger, and Hans-Peter Kriegel. Density-based projected clustering over high dimensional data streams. In *SDM*, pages 987–998. SIAM, 2012.
- Liadan O’callaghan, Adam Meyerson, Rajeev Motwani, Nina Mishra, and Sudipto Guha. Streaming-data algorithms for high-quality clustering. In *IEEE 29th International Conference on Data Engineering (ICDE)*, pages 0685–0685, 2002.
- Francesco Orabona and Nicolò Cesa-Bianchi. Better algorithms for selective sampling. In *Proc. 28th International Conference on Machine Learning (ICML2011)*, pages 433–440, 2011.
- Francesco Orabona and Koby Crammer. New adaptive algorithms for online classification. In *Advances in neural information processing systems*, pages 1840–1848, 2010.
- Francesco Orabona, Joseph Keshet, and Barbara Caputo. Bounded kernel-based online learning. *The Journal of Machine Learning Research*, 10:2643–2666, 2009.
- Mihály Ormos and András Urbán. Performance analysis of log-optimal portfolio strategies with transaction costs. *Quantitative Finance*, 13(10):1587–1597, 2013.
- Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *Knowledge and Data Engineering, IEEE Transactions on*, 22(10):1345–1359, 2010.
- German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *arXiv preprint arXiv:1802.07569*, 2018.
- John Platt. A resource-allocating network for function interpolation. *Neural computation*, 3(2):213–225, 1991.

- John Platt et al. Fast training of support vector machines using sequential minimal optimization. *Advances in kernel methodssupport vector learning*, 3, 1999.
- Gert CauwenberghsTomaso Poggio. Incremental and decremental support vector machine learning. In *Advances in Neural Information Processing Systems 13: Proceedings of the 2000 Conference*, volume 13, page 409. MIT Press, 2001.
- Robi Polikar, Lalita Upda, Satish S Upda, and Vasant Honavar. Learn++: An incremental learning algorithm for supervised neural networks. *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, 31(4):497–508, 2001.
- Abdulhakim Qahtan, Suojin Wang, and Xiangliang Zhang. Kde-track: An efficient dynamic density estimator for data streams. *IEEE Transactions on Knowledge and Data Engineering*, 29(3):642–655, 2017.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pages 1177–1184, 2007.
- Alexander Rakhlin. Lecture notes on online learning. *Notes appeared in the Statistical Learning Theory course at UC Berkeley*, 2008.
- Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. Online learning: Random averages, combinatorial parameters, and learnability. In *Advances in Neural Information Processing Systems*, pages 1984–1992, 2010.
- Alain Rakotomamonjy, Francis R. Bach, Stephane Canu, and Yves Grandvalet. Simplemkl. *J. Mach. Learn. Res. (JMLR)*, 11:2491–2521, 2008.
- Jesse Read, Albert Bifet, Bernhard Pfahringer, and Geoff Holmes. Batch-incremental versus instance-incremental learning in dynamic and evolving data. In *Advances in Intelligent Data Analysis XI*, pages 313–323. Springer, 2012.
- Jiadong Ren and Ruiqing Ma. Density-based data streams clustering over sliding windows. In *Fuzzy Systems and Knowledge Discovery, 2009. FSKD’09. Sixth International Conference on*, volume 5, pages 248–252. IEEE, 2009.
- Jiadong Ren, Binlei Cai, and Changzhen Hu. Clustering over data streams based on grid density and index tree. *Journal of Convergence Information Technology*, 6(1), 2011.
- Herbert Robbins. Some aspects of the sequential design of experiments. In *Herbert Robbins Selected Papers*, pages 169–177. Springer, 1985.
- Frank Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386, 1958.
- Dan Roth, Kevin Small, and Ivan Titov. Sequential learning of classifiers for structured prediction problems. In *International Conference on Artificial Intelligence and Statistics*, pages 440–447, 2009.
- Tim Roughgarden and Okke Schrijvers. Online prediction with selfish experts. In *Advances In Neural Information Processing Systems*, 2017.

- Sam T Roweis and Lawrence K Saul. Nonlinear dimensionality reduction by locally linear embedding. *science*, 290(5500):2323–2326, 2000.
- Carlos Ruiz, Ernestina Menasalvas, and Myra Spiliopoulou. C-denstream: Using domain knowledge on a data stream. In *Discovery Science*, pages 287–301. Springer, 2009.
- Paat Rusmevichientong and John N Tsitsiklis. Linearly parameterized bandits. *Mathematics of Operations Research*, 35(2):395–411, 2010.
- Daniel Russo and Benjamin Van Roy. An information-theoretic analysis of thompson sampling. *The Journal of Machine Learning Research*, 17(1):2442–2471, 2016.
- Paul Ruvolo and Eric Eaton. Ella: An efficient lifelong learning algorithm. In *International Conference on Machine Learning*, pages 507–515, 2013.
- Avishek Saha, Piyush Rai, Hal Daum<sup>Ã</sup>, Suresh Venkatasubramanian, et al. Online learning of multiple tasks and their relationships. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages 643–651, 2011.
- Doyen Sahoo, Steven CH Hoi, and Bin Li. Online multiple kernel regression. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 293–302. ACM, 2014.
- Doyen Sahoo, Abhishek Sharma, Steven CH Hoi, and Peilin Zhao. Temporal kernel descriptors for learning with time-sensitive patterns. In *Proceedings of the 2016 SIAM International Conference on Data Mining*, pages 540–548. SIAM, 2016a.
- Doyen Sahoo, Peilin Zhao, and Steven CH Hoi. Cost-sensitive online multiple kernel classification. In *Proceedings of The 8th Asian Conference on Machine Learning*, pages 65–80, 2016b.
- Doyen Sahoo, Quang Pham, Jing Lu, and Steven C. H. Hoi. Online deep learning: Learning deep neural networks on the fly. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence IJCAI-18*, pages 2660–2666, 7 2018.
- Nicholas I Sapankevych and Ravi Sankar. Time series prediction using support vector machines: a survey. *Computational Intelligence Magazine, IEEE*, 4(2):24–38, 2009.
- Bernhard Schölkopf, Ralf Herbrich, and Alex J. Smola. A generalized representer theorem. In *COLT/EuroCOLT*, pages 416–426, 2001.
- Sebastian Schuon, Marko Durković, Klaus Diepold, Jürgen Scheuerle, and Stefan Markward. Truly incremental locally linear embedding. In *CoTeSys 1st International Workshop on Cognition for Technical Systems*, 2008.
- David W Scott. *Multivariate density estimation: theory, practice, and visualization*. John Wiley & Sons, 2015.
- Steven L Scott. A modern bayesian look at the multi-armed bandit. *Applied Stochastic Models in Business and Industry*, 26(6):639–658, 2010.

- H Sebastian Seung, Manfred Opper, and Haim Sompolsky. Query by committee. In *Proc. 5th annual workshop on Computational learning theory*, pages 287–294, 1992.
- Parikshit Shah, Akshay Soni, and Troy Chevalier. Online ranking with constraints: A primal-dual algorithm and applications to web traffic-shaping. In *KDD*, 2017.
- Shai Shalev-Shwartz. *Online learning: Theory, algorithms, and applications*. PhD thesis, The Hebrew University of Jerusalem, 2007.
- Shai Shalev-Shwartz. Online learning and online convex optimization. *Foundations and Trends in Machine Learning*, 4(2):107–194, 2011.
- Shai Shalev-Shwartz and Yoram Singer. A primal-dual perspective of online learning algorithms. *Machine Learning*, 69(2-3):115–142, 2007.
- Shai Shalev-Shwartz and Ambuj Tewari. Stochastic methods for  $l_1$ -regularized loss minimization. *The Journal of Machine Learning Research*, 999999:1865–1892, 2011.
- Shai Shalev-Shwartz, Yoram Singer, and Andrew Y Ng. Online and batch learning of pseudo-metrics. In *Proceedings of the twenty-first international conference on Machine learning*, page 94. ACM, 2004.
- Shai Shalev-Shwartz, Yoram Singer, Nathan Srebro, and Andrew Cotter. Pegasos: Primal estimated sub-gradient solver for svm. *Mathematical programming*, 127(1):3–30, 2011.
- Yue Shi, Martha Larson, and Alan Hanjalic. Collaborative filtering beyond the user-item matrix: A survey of the state of the art and future challenges. *ACM Computing Surveys (CSUR)*, 47(1):3, 2014.
- Daniel L Silver, Qiang Yang, and Lianghao Li. Lifelong machine learning systems: Beyond learning algorithms. In *AAAI Spring Symposium: Lifelong Machine Learning*, volume 13, page 05, 2013.
- Bernard W Silverman. *Density estimation for statistics and data analysis*. Routledge, 2018.
- Padhraic Smyth, Max Welling, and Arthur U Asuncion. Asynchronous distributed learning of topic models. In *NIPS*, pages 81–88, 2009.
- Soeren Sonnenburg and Vojtvech Franc. Coffin: a computational framework for linear svms. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, pages 999–1006. Omnipress, 2010.
- Sören Sonnenburg, Gunnar Rätsch, Christin Schäfer, and Bernhard Schölkopf. Large scale multiple kernel learning. *J. Mach. Learn. Res. (JMLR)*, 7:1531–1565, 2006.
- Ricardo Sousa, Luis M Silva, Luis A Alexandre, Jorge Santos, and Joaquim Marques de Sá. Transfer learning: Current status, trends and challenges.
- Eduardo J Spinoso, F de Leon, André Ponce, and João Gama. Novelty detection with application to data streams. *Intelligent Data Analysis*, 13(3):405–422, 2009.

- Xiaoyuan Su and Taghi M Khoshgoftaar. A survey of collaborative filtering techniques. *Advances in artificial intelligence*, 2009:4, 2009.
- Shiliang Sun. A survey of multi-view machine learning. *Neural Computing and Applications*, 23(7-8):2031–2038, 2013.
- Richard S. Sutton and Andrew G. Barto. *Reinforcement learning: An introduction*. MIT press, 1998.
- Swee Chuan Tan, Kai Ming Ting, and Tony Fei Liu. Fast anomaly detection for streaming data. In *IJCAI Proceedings-International Joint Conference on Artificial Intelligence*, volume 22, page 1511, 2011.
- Dimitris K Tasoulis, Gordon Ross, and Niall M Adams. Visualising the cluster structure of data streams. In *Advances in Intelligent Data Analysis VII*, pages 81–92. Springer, 2007.
- Joshua B Tenenbaum, Vin De Silva, and John C Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.
- 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, 1933.
- Tatiana Tommasi, Francesco Orabona, Mohsen Kaboli, Barbara Caputo, and CH Martigny. Leveraging over prior knowledge for online learning of visual categories. In *BMVC*, 2012.
- Andrew Trotman. Learning to rank. *Information Retrieval*, 8(3):359–381, 2005.
- P. Tseng. On accelerated proximal gradient methods for Convex-Concave optimization. *SIAM Journal on Optimization*, 2008.
- Li Tu and Yixin Chen. Stream data clustering based on grid density and attraction. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 3(3):12, 2009.
- Q Tu, JF Lu, B Yuan, JB Tang, and Jing-Yu Yang. Density-based hierarchical clustering for streaming data. *Pattern Recognition Letters*, 33(5):641–645, 2012.
- Taishi Uchiya, Atsuyoshi Nakamura, and Mineichi Kudo. Algorithms for adversarial bandit problems with multiple plays. In *International Conference on Algorithmic Learning Theory*, pages 375–389. Springer, 2010.
- Michal Valko, Branislav Kveton, Huang Ling, and Ting Daniel. Online semi-supervised learning on quantized graphs. In *Uncertainty in Artificial Intelligence*, 2010.
- Tim van Erven and Wouter M Koolen. Metagrad: Multiple learning rates in online learning. In *Advances in Neural Information Processing Systems*, pages 3666–3674, 2016.
- Vladimir N Vapnik. An overview of statistical learning theory. *IEEE transactions on neural networks*, 10(5):988–999, 1999.
- Vladimir Naumovich Vapnik and Vlamimir Vapnik. *Statistical learning theory*, volume 1. Wiley New York, 1998.

- Joannes Vermorel and Mehryar Mohri. Multi-armed bandit algorithms and empirical evaluation. In *Machine Learning: ECML 2005*, pages 437–448. Springer, 2005.
- Volodya Vovk and Chris Watkins. Universal portfolio selection. In *Proceedings of the eleventh annual conference on Computational learning theory*, pages 12–23. ACM, 1998.
- Ji Wan, Pengcheng Wu, Steven CH Hoi, Peilin Zhao, Xingyu Gao, Dayong Wang, Yongdong Zhang, and Jintao Li. Online learning to rank for content-based image retrieval. In *IJCAI*, pages 2284–2290, 2015.
- Li Wan, Wee Keong Ng, Xuan Hong Dang, Philip S Yu, and Kuan Zhang. Density-based clustering of data streams at multiple resolutions. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 3(3):14, 2009.
- Chuang Wang and Yue Lu. The scaling limit of high-dimensional online independent component analysis. In *Advances in Neural Information Processing Systems*, pages 6638–6647, 2017.
- Dayong Wang, Pengcheng Wu, Peilin Zhao, Yue Wu, Chunyan Miao, and Steven CH Hoi. High-dimensional data stream classification via sparse online learning. In *Data Mining (ICDM), 2014 IEEE International Conference on*, pages 1007–1012. IEEE, 2014a.
- Dayong Wang, Pengcheng Wu, Peilin Zhao, and Steven CH Hoi. A framework of sparse online learning and its applications. *arXiv preprint arXiv:1507.07146*, 2015a.
- Haixun Wang, Wei Fan, Philip S Yu, and Jiawei Han. Mining concept-drifting data streams using ensemble classifiers. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 226–235. ACM, 2003.
- Huahua Wang and Arindam Banerjee. Online alternating direction method. In *Proceedings of the 29th International Conference on Machine Learning, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012*, 2012.
- Jialei Wang, Peilin Zhao, and Steven C. H. Hoi. Cost-sensitive online classification. In *12th IEEE International Conference on Data Mining (ICDM2012)*, pages 1140–1145, 2012a.
- Jialei Wang, Peilin Zhao, and Steven CH Hoi. Exact soft confidence-weighted learning. In *Proceedings of the 29th International Conference on Machine Learning*, pages 107–114. Omnipress, 2012b.
- Jialei Wang, Steven CH Hoi, Peilin Zhao, and Zhi-Yong Liu. Online multi-task collaborative filtering for on-the-fly recommender systems. In *Proceedings of the 7th ACM conference on Recommender systems*, pages 237–244. ACM, 2013a.
- Jialei Wang, Steven CH Hoi, Peilin Zhao, Jinfeng Zhuang, and Zhi-yong Liu. Large scale online kernel classification. In *Proceedings of the Twenty-Third international joint conference on Artificial Intelligence*, pages 1750–1756. AAAI Press, 2013b.
- Jialei Wang, Peilin Zhao, and Steven CH Hoi. Cost-sensitive online classification. *Knowledge and Data Engineering, IEEE Transactions on*, 26(10):2425–2438, 2014b.



- Jialei Wang, Peilin Zhao, Steven CH Hoi, and Rong Jin. Online feature selection and its applications. *IEEE Transactions on Knowledge and Data Engineering*, 26(3):698–710, 2014c.
- Jialei Wang, Ji Wan, Yongdong Zhang, and Steven CH Hoi. Solar: Scalable online learning algorithms for ranking. *ACL*, 2015b.
- Jialei Wang, Peilin Zhao, and Steven CH Hoi. Soft confidence-weighted learning. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 8(1):15, 2016a.
- S Wang, Leandro Lei Minku, and X Yao. Dealing with multiple classes in online class imbalance learning. In *International Joint Conferences on Artificial Intelligence*, 2016b.
- Shijun Wang, Rong Jin, and Hamed Valizadegan. A potential-based framework for online multi-class learning with partial feedback. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 900–907, 2010.
- Yizao Wang, Jean-Yves Audibert, and Rémi Munos. Algorithms for infinitely many-armed bandits. In *Advances in Neural Information Processing Systems*, pages 1729–1736, 2009.
- Yuyang Wang, Roni Khardon, Dmitry Pechyony, and Rosie Jones. Generalization bounds for online learning algorithms with pairwise loss functions. In *COLT*, volume 23, pages 13–1, 2012c.
- Zhuang Wang and Slobodan Vucetic. Tighter perceptron with improved dual use of cached data for model representation and validation. In *Neural Networks, 2009. IJCNN 2009. International Joint Conference on*, pages 3297–3302. IEEE, 2009.
- Zhuang Wang and Slobodan Vucetic. Online passive-aggressive algorithms on a budget. *Journal of Machine Learning Research - Proceedings Track*, 9:908–915, 2010.
- Zhuang Wang, Koby Crammer, and Slobodan Vucetic. Breaking the curse of kernelization: Budgeted stochastic gradient descent for large-scale svm training. *The Journal of Machine Learning Research*, 13(1):3103–3131, 2012d.
- Malcolm Ware, Eibe Frank, Geoffrey Holmes, Mark Hall, and Ian H Witten. Interactive machine learning: letting users build classifiers. *International Journal of Human-Computer Studies*, 55(3):281–292, 2001.
- Manfred K Warmuth and Dima Kuzmin. Randomized online pca algorithms with regret bounds that are logarithmic in the dimension. *Journal of Machine Learning Research*, 9(10), 2008.
- Jason Weston, Antoine Bordes, Léon Bottou, et al. Online (and offline) on an even tighter budget. In *Proceedings of the 10th International Workshop on Artificial Intelligence and Statistics*, pages 413–420, 2005.
- Christopher K. I. Williams and Matthias Seeger. Using the nyström method to speed up kernel machines. In T.K. Leen, T.G. Dietterich, and V. Tresp, editors, *Advances in Neural Information Processing Systems 13*, pages 682–688. MIT Press, 2001.

- Ronald J Williams and David Zipser. A learning algorithm for continually running fully recurrent neural networks. *Neural computation*, 1(2):270–280, 1989.
- Pengcheng Wu, Steven CH Hoi, Hao Xia, Peilin Zhao, Dayong Wang, and Chunyan Miao. Online multimodal deep similarity learning with application to image retrieval. In *Proceedings of 21st ACM international conference on Multimedia*, pages 153–162, 2013.
- Pengcheng Wu, Steven CH Hoi, Peilin Zhao, Chunyan Miao, and Zhi-Yong Liu. Online multi-modal distance metric learning with application to image retrieval. *IEEE Transactions on Knowledge and Data Engineering*, 28(2):454–467, 2016.
- Yue Wu, Steven CH Hoi, and Tao Mei. Massive-scale online feature selection for sparse ultra-high dimensional data. *arXiv preprint arXiv:1409.7794*, 2014.
- Yue Wu, Steven CH Hoi, Chenghao Liu, Jing Lu, Doyen Sahoo, and Nenghai Yu. Sol: A library for scalable online learning algorithms. *Neurocomputing*, 260:9–12, 2017a.
- Yue Wu, Steven CH Hoi, Tao Mei, and Nenghai Yu. Large-scale online feature selection for ultra-high dimensional sparse data. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 11(4):48, 2017b.
- Hao Xia, Pengcheng Wu, and Steven CH Hoi. Online multi-modal distance learning for scalable multimedia retrieval. In *Proceedings of the sixth ACM international conference on Web search and data mining*, pages 455–464. ACM, 2013.
- Hao Xia, Steven CH Hoi, Rong Jin, and Peilin Zhao. Online multiple kernel similarity learning for visual search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(3):536–549, 2014.
- Lin Xiao. Dual averaging method for regularized stochastic learning and online optimization. In *Advances in Neural Information Processing Systems*, pages 2116–2124, 2009.
- Chang Xu, Dacheng Tao, and Chao Xu. A survey on multi-view learning. *arXiv preprint arXiv:1304.5634*, 2013.
- Zenglin Xu, Rong Jin, Irwin King, and Michael R. Lyu. An extended level method for efficient multiple kernel learning. In *NIPS*, 2008.
- Di Yang, Elke A Rundensteiner, and Matthew O Ward. Neighbor-based pattern detection for windows over streaming data. In *Proceedings of the 12th International Conference on Extending Database Technology (EDBT)*, pages 529–540, 2009a.
- Haiqin Yang, Irwin King, and Michael R Lyu. Online learning for multi-task feature selection. In *Proceedings of the 19th ACM international conference on Information and knowledge management*, pages 1693–1696. ACM, 2010.
- Liu Yang and Rong Jin. Distance metric learning: A comprehensive survey. *Michigan State University*, 2(2):4, 2006.
- Liu Yang, Rong Jin, and Jieping Ye. Online learning by ellipsoid method. In *Proceedings of 26th International Conference on Machine Learning*, pages 1153–1160, 2009b.

- Tianbao Yang, Mehrdad Mahdavi, Rong Jin, Jinfeng Yi, and Steven CH Hoi. Online kernel selection: Algorithms and evaluations. In *AAAI*, 2012.
- Peilin Zhao Steven Hoi Yi Ding, Chenghao Liu. Large scale kernel methods for online auc maximization. In *The IEEE International Conference on Data Mining (ICDM)*, 2017.
- Yiming Ying, Longyin Wen, and Siwei Lyu. Stochastic online auc maximization. In *Advances in Neural Information Processing Systems*, pages 451–459, 2016.
- Li Yuan-Xiang, Li Zhi-Jie, Wang Feng, and Kuang Li. Accelerated online learning for collaborative filtering and recommender systems. In *Data Mining Workshop (ICDMW), 2014 IEEE International Conference on*, pages 879–885. IEEE, 2014.
- Chunqiu Zeng, Qing Wang, Shekoofeh Mokhtari, and Tao Li. Online context-aware recommendation with time varying multi-arm bandit. *KDD*, 2016.
- Lijun Zhang, Rong Jin, Chun Chen, Jiajun Bu, and Xiaofei He. Efficient online learning for large-scale sparse kernel logistic regression. In *AAAI*, 2012.
- Lijun Zhang, Tianbao Yang, Rong Jin, Yichi Xiao, and Zhi-Hua Zhou. Online stochastic linear optimization under one-bit feedback. In *International Conference on Machine Learning*, pages 392–401, 2016a.
- Tong Zhang. Solving large scale linear prediction problems using stochastic gradient descent algorithms. In *Proc. 21th International Conference on Machine Learning (ICML’04)*, 2004.
- Tong Zhang. Data dependent concentration bounds for sequential prediction algorithms. In *18th Annual Conference on Learning Theory(COLT’05)*, pages 173–187, 2005.
- Wenpeng Zhang, Peilin Zhao, Wenwu Zhu, Steven CH Hoi, and Tong Zhang. Projection-free distributed online learning in networks. In *International Conference on Machine Learning*, pages 4054–4062, 2017.
- Xiaoxuan Zhang, Tianbao Yang, and Padmini Srinivasan. Online asymmetric active learning with imbalanced data. *KDD*, 2016b.
- Peilin Zhao. *Kernel Based Online Learning*. PhD thesis, Nanyang Technological University, 2013.
- Peilin Zhao and Steven C Hoi. Otl: A framework of online transfer learning. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 1231–1238, 2010.
- Peilin Zhao and Steven CH Hoi. Bduol: double updating online learning on a fixed budget. In *Machine Learning and Knowledge Discovery in Databases*, pages 810–826. 2012.
- Peilin Zhao and Steven CH Hoi. Cost-sensitive online active learning with application to malicious url detection. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 919–927. ACM, 2013.

- Peilin Zhao, Steven C Hoi, and Rong Jin. Duol: A double updating approach for online learning. In *Advances in Neural Information Processing Systems*, pages 2259–2267, 2009.
- Peilin Zhao, Steven CH Hoi, and Rong Jin. Double updating online learning. *The Journal of Machine Learning Research*, 12:1587–1615, 2011a.
- Peilin Zhao, Rong Jin, Tianbao Yang, and Steven C Hoi. Online auc maximization. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 233–240, 2011b.
- Peilin Zhao, Jiale Wang, Pengcheng Wu, Rong Jin, and Steven C. H. Hoi. Fast bounded online gradient descent algorithms for scalable kernel-based online learning. In *ICML*, 2012.
- Peilin Zhao, Steven C. H. Hoi, and Jinfeng Zhuang. Active learning with expert advice. In *Proceedings of the 29th Conference on Uncertainty in Artificial Intelligence*, 2013.
- Peilin Zhao, Steven CH Hoi, Jiale Wang, and Bin Li. Online transfer learning. *Artificial Intelligence*, 216:76–102, 2014.
- Peilin Zhao, Furen Zhuang, Min Wu, Xiao-Li Li, and Steven CH Hoi. Cost-sensitive online classification with adaptive regularization and its applications. In *Data Mining (ICDM), 2015 IEEE International Conference on*, pages 649–658. IEEE, 2015.
- Yan Zheng, Jeffrey Jests, Jeff M Phillips, and Feifei Li. Quality and efficiency for kernel density estimates in large data. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, pages 433–444. ACM, 2013.
- Aoying Zhou, Zhiyuan Cai, Li Wei, and Weining Qian. M-kernel merging: Towards density estimation over data streams. In *DASFAA*, pages 285–292. IEEE, 2003.
- Li Zhou. A survey on contextual multi-armed bandits. *CoRR*, abs/1508.03326, 2015.
- Xiaojin Zhu. Semi-supervised learning literature survey. *Computer Science, University of Wisconsin-Madison*, 2(3):4, 2006.
- Xiaojin Zhu and Andrew B Goldberg. Introduction to semi-supervised learning. *Synthesis lectures on artificial intelligence and machine learning*, 3(1):1–130, 2009.
- Xiaojin Zhu, Zoubin Ghahramani, and John D Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *Proceedings of the 20th International conference on Machine learning (ICML-03)*, pages 912–919, 2003.
- Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the Twentieth International Conference on Machine Learning (ICML 2003)*, pages 928–936, 2003.
- Masrour Zoghi, Tomas Tunys, Mohammad Ghavamzadeh, Branislav Kveton, Csaba Szepesvari, and Zheng Wen. Online learning to rank in stochastic click models. In *International Conference on Machine Learning*, pages 4199–4208, 2017.