

Risk Bounds for Low Cost Bipartite Ranking

San Gultekin and John Paisley

November 24, 2019

Abstract

Bipartite ranking is an important supervised learning problem; however, unlike regression or classification, it has a quadratic dependence on the number of samples. To circumvent the prohibitive sample cost, many recent work focus on stochastic gradient-based methods. In this paper we consider an alternative approach, which leverages the structure of the widely-adopted pairwise squared loss, to obtain a stochastic and low cost algorithm that does not require stochastic gradients or learning rates. Using a novel uniform risk bound—based on matrix and vector concentration inequalities—we show that the sample size required for competitive performance against the all-pairs batch algorithm does not have a quadratic dependence. Generalization bounds for both the batch and low cost stochastic algorithms are presented. Experimental results show significant speed gain against the batch algorithm, as well as competitive performance against state-of-the-art bipartite ranking algorithms on real datasets.

1 Introduction

Binary classification is among the most widely studied machine learning problems, with many applications. Given a binary labeled dataset, the aim is to learn a mapping from the features to the labels. The performance of a learning algorithm is typically gauged in terms of classification error. However, in situations such as cost-sensitive learning [1] and imbalanced learning [2, 3] this choice may not be appropriate. For example, in online advertising [1] one is typically concerned with separating the interesting ads from the rest. This problem is also known as bipartite ranking, where the aim is to rank the “positive” inputs higher than the “negative” ones.

1.1 Problem Setup

Let \mathcal{X} be a D -dimensional input domain and \mathcal{Y} the label domain. The bipartite ranking problem is defined for binary labeled samples, $\mathcal{Y} = \{0, 1\}$. These sample pairs are generated i.i.d. from

an unknown distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$. A ranking function $f : \mathcal{X} \rightarrow \mathbb{R}$ is a mapping from the inputs to a scalar-valued score. In this paper we are interested in linear ranking functions of the form $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$ with parameter vector \mathbf{w} . The goal in bipartite ranking is to assign higher scores to the inputs with label 1. Let \mathbf{x}^1 and \mathbf{x}^0 denote two samples with corresponding labels of 1 and 0. In many ranking settings we are given a set of features that encodes the preference of one sample over another [1], which can be represented by $\Psi(\mathbf{x}^1, \mathbf{x}^0)$. It is also possible that the individual features \mathbf{x}_0 and \mathbf{x}_1 are not defined explicitly; for example when a single user is shown two ads and prefers one over another. Given this, a widely used performance metric is the Wilcoxon-Mann-Whitney statistic

$$\ell_{\text{WMW}}(\mathbf{w}, \Psi(\mathbf{x}^1, \mathbf{x}^0)) = \mathbb{I}\{\mathbf{w}^\top \Psi(\mathbf{x}^1, \mathbf{x}^0) > 0\} + \frac{1}{2} \mathbb{I}\{\mathbf{w}^\top \Psi(\mathbf{x}^1, \mathbf{x}^0) = 0\} . \quad (1)$$

Based on this we define the risks

$$R^{\text{AUC}}(\mathbf{w}) = \mathbb{E}_{\substack{\mathbf{x}^1 \sim D^1 \\ \mathbf{x}^0 \sim D^0}} [\ell_{\text{WMW}}(\mathbf{w}, \Psi(\mathbf{x}^1, \mathbf{x}^0))] , \quad R_{\mathcal{N}}^{\text{UC}}(\mathbf{w}) = \frac{1}{N_1 N_0} \sum_{i=1}^{N_1} \sum_{j=1}^{N_0} \ell_{\text{WMW}}(\mathbf{w}, \Psi(\mathbf{x}_i^1, \mathbf{x}_j^0)) , \quad (2)$$

where $R^{\text{AUC}}(\mathbf{w})$ corresponds to the actual AUC risk and is obtained by taking expectation over the class-conditional distributions $D^i = D(\mathbf{x}|y=i)$. Let the sample set $\mathcal{N} = \{\mathbf{x}_1^0, \dots, \mathbf{x}_{N_0}^0, \mathbf{x}_1^1, \dots, \mathbf{x}_{N_1}^1\}$, containing N_0 and N_1 samples with labels 0 and 1, respectively, and define $N := N_0 + N_1$. This gives the empirical risk $R_{\mathcal{N}}^{\text{UC}}(\mathbf{w})$ to be minimized in practice. The main drawback of this approach is that the objective function is now a sum of indicator functions, an NP-hard problem. A widely-used approach to handle this problem is to replace the intrinsic loss function ℓ_{WMW} with a convex one, often written as ℓ_ϕ [4, 5].

In this paper we are interested in the pairwise squared loss as it is a consistent estimator of AUC [6], and widely preferred by recent work [7, 8, 9]:

$$\ell_\phi(\mathbf{w}, (\mathbf{x}^1, \mathbf{x}^0)) = \frac{1}{2} [1 - \mathbf{w}^\top \Psi(\mathbf{x}_i^1, \mathbf{x}_j^0)]^2 , \quad (3)$$

which produces its own corresponding actual and empirical ϕ -risks that parallels Eq. (2). Let $\Sigma_N = 1/(N_1 N_0) \sum_{i=1}^{N_1} \sum_{j=1}^{N_0} \Psi(\mathbf{x}_i^1, \mathbf{x}_j^0) \Psi(\mathbf{x}_i^1, \mathbf{x}_j^0)^\top$ and $\mu_N = 1/(N_1 N_0) \sum_{i=1}^{N_1} \sum_{j=1}^{N_0} \Psi(\mathbf{x}_i^1, \mathbf{x}_j^0)$. The empirical risk for this problem can then also written as the optimization of $R_{\mathcal{N}}^\phi(\mathbf{w}) = (1/2) \mathbf{w}^\top \Sigma_N \mathbf{w} - \mu_N^\top \mathbf{w}$. For what follows we will drop the ϕ symbol and refer to ℓ_ϕ , R^ϕ , and $R_{\mathcal{N}}^\phi$ simply as ℓ , R , and $R_{\mathcal{N}}$.

The focus of this paper is on two algorithms: The first one is a Batch Bipartite Ranking (BBR) algorithm (Algorithm 1), which aims at minimizing $R_{\mathcal{N}}$ based on all pairs available in \mathcal{N} . We provide a new theoretical analysis of BBR in Section 2. Due to the quadratic growth of sample size, $O(N_1 N_0)$, the sample cost of BBR quickly becomes prohibitive. Therefore we also propose a new Low Cost Bipartite Ranking (LCBR) algorithm (Algorithm 2), which, given the same sample set \mathcal{N} , subsamples S pairs uniformly at random with replacement. The main goal of this paper is to analyze the subsample size S required for LCBR to be competitive with BBR. As we show in Section 3, S does not have such quadratic dependence. Section 4 discusses related work. We show experiments in Section 5 and conclude in Section 6.

Algorithm 1 BBR	Algorithm 2 LCBR
Input: Sample set \mathcal{N} Regularization parameter (W_*)	Input: Sample set \mathcal{N} , subsample size (S) Regularization parameter (W_*)
Output: Linear ranker's weight \mathbf{w}_N	Output: Linear ranker's weight \mathbf{w}_S
1. Initialize $\boldsymbol{\mu}_N \leftarrow \mathbf{0}$, $\boldsymbol{\Sigma}_N \leftarrow \mathbf{0}$	1. Initialize $\boldsymbol{\mu}_S \leftarrow \mathbf{0}$, $\boldsymbol{\Sigma}_S \leftarrow \mathbf{0}$
2. //Accummulation	2. //Accummulation
3. for $i = 1, \dots, N_1$	3. for $s = 1, \dots, S$
4. for $j = 1, \dots, N_0$	4. Sample (i_s, j_s) uniformly with replacement
5. $\boldsymbol{\mu}_N \leftarrow \boldsymbol{\mu}_N + \frac{1}{N_1 N_0}(\mathbf{x}_i^1 - \mathbf{x}_j^0)$	5. $\boldsymbol{\mu}_S \leftarrow \boldsymbol{\mu}_S + \frac{1}{S}(\mathbf{x}_{i_s}^1 - \mathbf{x}_{j_s}^0)$
6. $\boldsymbol{\Sigma}_N \leftarrow \boldsymbol{\Sigma}_N + \frac{1}{N_1 N_0}(\mathbf{x}_i^1 - \mathbf{x}_j^0)(\mathbf{x}_i^1 - \mathbf{x}_j^0)^\top$	6. $\boldsymbol{\Sigma}_S \leftarrow \boldsymbol{\Sigma}_S + \frac{1}{S}(\mathbf{x}_{i_s}^1 - \mathbf{x}_{j_s}^0)(\mathbf{x}_{i_s}^1 - \mathbf{x}_{j_s}^0)^\top$
7. end for	7. end for
8. end for	8.
9. //Empirical Risk Minimization	9. //Empirical Risk Minimization
10. $\mathbf{w}_N \leftarrow \arg \min_{\mathbf{w} \in \mathcal{B}_2(W_*)} \frac{1}{2} \mathbf{w}^\top \boldsymbol{\Sigma}_N \mathbf{w} - \boldsymbol{\mu}_N^\top \mathbf{w}$	10. $\mathbf{w}_S \leftarrow \arg \min_{\mathbf{w} \in \mathcal{B}_2(W_*)} \frac{1}{2} \mathbf{w}^\top \boldsymbol{\Sigma}_S \mathbf{w} - \boldsymbol{\mu}_S^\top \mathbf{w}$

Remark: At this point, for the ease of presentation we focus on the case $\Psi(\mathbf{x}_i^1, \mathbf{x}_j^0) = \mathbf{x}_i^1 - \mathbf{x}_j^0$ (which is reflected in Algorithms 1 and 2). All the theory developed in this paper is valid for the most general case of $\Psi(\mathbf{x}_i^1, \mathbf{x}_j^0)$, however. In particular, imposing a norm bound on the individual features $\mathbf{x}_i^1, \mathbf{x}_j^0$ is equivalent to imposing a scaled norm bound on $\Psi(\mathbf{x}_i^1, \mathbf{x}_j^0)$.

Additional assumptions: We assume the input domain is compact, $\mathcal{X} \subseteq \mathcal{B}_2(X_*)$ which implies $\|\mathbf{x}\|_2 \leq X_*$, and that the domain of ranking functions is compact, taking $\mathcal{W} = \mathcal{B}_2(W_*)$ such that $\|\mathbf{w}\|_2 \leq W_*$. Here $\mathcal{B}_p(r) := \{\mathbf{x} \in \mathcal{X} : \|\mathbf{x}\|_p \leq r\}$ is the ℓ_p -ball of radius r . The first assumption is related to data preprocessing where the features are scaled appropriately; the second assumption corresponds to regularization. For both cases we chose the ℓ_2 -norm as it provides a dimension-independent upper bound; however it is still possible to derive bounds using other norms. Let \mathbf{w}_N and \mathbf{w}_S be the minimizer of the empirical risk objectives (Algorithms 1 and 2, Line 10). Also for clarity we focus on the case where $\boldsymbol{\Sigma}, \boldsymbol{\Sigma}_N, \boldsymbol{\Sigma}_S \succ \mathbf{0}$ in the sequel.¹ This means \mathbf{w}_* , \mathbf{w}_N , and \mathbf{w}_S have unique values. Note that, this last assumption is only for presentation purposes and the results derived in this paper apply to the most general case where $\boldsymbol{\Sigma}, \boldsymbol{\Sigma}_N, \boldsymbol{\Sigma}_S \succeq \mathbf{0}$.

2 Batch Bipartite Ranking (BBR)

In order to analyze the more efficient LCBR, we first derive a risk bound for the corresponding BBR. The risk bounds derived in this paper are with respect to the best-in-class ranking function. For the sample set \mathcal{N} we define $\rho := N_1/(N_0 + N_1)$ as the label skew. Also, for the unit sphere \mathcal{S}^{D-1} let $C_{\mathcal{S}}(\epsilon)$ denote the covering number based on ℓ_2 -balls of radius ϵ . The main result of this section is the following risk bound based on the metric entropy of \mathcal{W} .

¹Where we defined $\boldsymbol{\Sigma} = \mathbb{E}_{\mathbf{x}_1 \sim P_1, \mathbf{x}_0 \sim P_0}[(\mathbf{x}_1 - \mathbf{x}_0)(\mathbf{x}_1 - \mathbf{x}_0)^\top]$.

Theorem 1. For a given sample set \mathcal{N} , define $\mathbf{w}_N := \arg \min_{\mathbf{w} \in \mathcal{W}} R_{\mathcal{N}}(\mathbf{w})$ and $\mathbf{w}_* = \arg \min_{\mathbf{w} \in \mathcal{W}} R(\mathbf{w})$. Also define the constants $C_1 := 8X_*^2 W_* + 4X_*$ and $C_2 := 3X_*^2 W_*^2 + 2X_* W_*$. Then

$$P_{\mathcal{N} \sim \mathcal{D}^N} \left(R(\mathbf{w}_N) - R(\mathbf{w}_*) \geq \epsilon \right) \leq 2 C_S \left(\frac{\epsilon}{4C_1} \right) \exp \left\{ -\frac{\epsilon^2}{8C_2^2} \rho(1-\rho)N \right\}. \quad (4)$$

Here we used the shorthand \mathcal{D}^N to denote the product measure over the samples, $[\mathcal{D}^1]^{N_1} \otimes [\mathcal{D}^0]^{N_0}$ where $[\mathcal{D}^i]^{N_i} = \otimes_{i=1}^{N_i} \mathcal{D}^i$. This result is comparable to the bound given for linear regression based on covering numbers [10]. One distinction, however, is the dependence on skew $\rho(1-\rho)$; when ρ is close to 0 or 1, the learner requires significantly more samples to achieve the same bound. The exponential term in Eq. (4) is comparable to the one obtained for AUC loss in Theorem 5 of [11]. On the other hand, when $\rho(1-\rho) = O(1)$ and N is large enough to bound the covering number by the exponential term in Eq. (4) we get the typical rate $O(\sqrt{\log(1/\delta)/N})$.

As the rate in Eq. (4) depends on $N_1 + N_0$ and not $N_1 N_0$, it is natural to only consider pairs of independent samples instead of their Cartesian product, which reduces the analysis to that of linear regression. However, this is not done in practice as it would discard information [7]. Indeed, for this reason, a number of works consider the all-pair problem, e.g. [7, 4, 5]; we do so similarly. Below, we prove Theorem 1 using the following Lemma.

Lemma 1. For a given sample set \mathcal{N} and constants $C_1 := 8X_*^2 W_* + 4X_*$, $C_2 := 3X_*^2 W_*^2 + 2X_* W_*$,

$$P_{\mathcal{N} \sim \mathcal{D}^N} \left(\sup_{\mathbf{w} \in \mathcal{W}} |R(\mathbf{w}) - R_{\mathcal{N}}(\mathbf{w})| \geq \epsilon \right) \leq 2 C_S \left(\frac{\epsilon}{2C_1} \right) \exp \left\{ -\frac{\epsilon^2}{2C_2^2} \rho(1-\rho)N \right\}. \quad (5)$$

Proof. Recall that pairwise squared loss is defined as $\ell(\mathbf{w}, \mathbf{x}^1, \mathbf{x}^0) = (1/2)(1 - h_{\mathbf{w}}(\mathbf{x}^1, \mathbf{x}^0))^2$ for $h_{\mathbf{w}}(\mathbf{x}^1, \mathbf{x}^0) = \mathbf{w}^\top (\mathbf{x}^1 - \mathbf{x}^0)$. Also define $\Phi_{\mathcal{N}}(\mathbf{w}) := R(\mathbf{w}) - R_{\mathcal{N}}(\mathbf{w})$. For $\mathbf{w}_1, \mathbf{w}_2 \in \mathcal{W}$,

$$\begin{aligned} \left| \ell(\mathbf{w}_1, \mathbf{x}^1, \mathbf{x}^0) - \ell(\mathbf{w}_2, \mathbf{x}^1, \mathbf{x}^0) \right| &= \left| \frac{1}{2}(1 - h_{\mathbf{w}_1}(\mathbf{x}^1, \mathbf{x}^0))^2 - \frac{1}{2}(1 - h_{\mathbf{w}_2}(\mathbf{x}^1, \mathbf{x}^0))^2 \right| \\ &\leq \frac{1}{2} |2 - h_{\mathbf{w}_1}(\mathbf{x}^1, \mathbf{x}^0) - h_{\mathbf{w}_2}(\mathbf{x}^1, \mathbf{x}^0)| |h_{\mathbf{w}_2}(\mathbf{x}^1, \mathbf{x}^0) - h_{\mathbf{w}_1}(\mathbf{x}^1, \mathbf{x}^0)| \\ &\leq (2X_* W_* + 1) |(\mathbf{w}_2 - \mathbf{w}_1)^\top (\mathbf{x}^1 - \mathbf{x}^0)| \\ &\leq (4X_*^2 W_* + 2X_*) \|\mathbf{w}_1 - \mathbf{w}_2\|_2. \end{aligned} \quad (6)$$

Using this bound we have

$$\begin{aligned} |\Phi_{\mathcal{N}}(\mathbf{w}_1) - \Phi_{\mathcal{N}}(\mathbf{w}_2)| &= |R(\mathbf{w}_1) - R_{\mathcal{N}}(\mathbf{w}_1) - R(\mathbf{w}_2) + R_{\mathcal{N}}(\mathbf{w}_2)| \\ &\leq \left| \mathbb{E}_{\substack{\mathbf{x}^1 \sim D^1 \\ \mathbf{x}^0 \sim D^0}} [\ell(\mathbf{w}_1, \mathbf{x}^1, \mathbf{x}^0) - \ell(\mathbf{w}_2, \mathbf{x}^1, \mathbf{x}^0)] \right| + \frac{1}{N_1 N_0} \sum_{i=1}^{N_1} \sum_{j=1}^{N_0} \left| \ell(\mathbf{w}_1, \mathbf{x}_i^1, \mathbf{x}_j^0) - \ell(\mathbf{w}_2, \mathbf{x}_i^1, \mathbf{x}_j^0) \right| \\ &\leq (8X_*^2 W_* + 4X_*) \|\mathbf{w}_1 - \mathbf{w}_2\|_2. \end{aligned} \quad (7)$$

We have shown that $|\Phi_{\mathcal{N}}(\mathbf{w}_1) - \Phi_{\mathcal{N}}(\mathbf{w}_2)| \leq C_1 \|\mathbf{w}_1 - \mathbf{w}_2\|_2$. Let $\{\mathcal{B}_i\}_{i=1}^I$ be a set of ℓ_2 -balls of

radius ϵ covering \mathcal{W} . Then

$$P_{\mathcal{N}} \left(\sup_{\mathbf{w} \in \mathcal{W}} |\Phi_{\mathcal{N}}(\mathbf{w})| \geq \epsilon \right) \leq \sum_{i=1}^I P_{\mathcal{N}} \left(\sup_{\mathbf{w} \in \mathcal{B}_i} |\Phi_{\mathcal{N}}(\mathbf{w})| \geq \epsilon \right) \leq \sum_{i=1}^I P_{\mathcal{N}} (|\Phi_{\mathcal{N}}(\mathbf{w}_i)| \geq \epsilon/2). \quad (8)$$

Now that the weight vector and samples are decoupled, we can bound the deviation of $\Phi_{\mathcal{N}}(\mathbf{w})$ for a fixed \mathbf{w} . First note that $\mathbb{E}[\Phi_{\mathcal{N}}(\mathbf{w})] = 0$ by definition of $R(\mathbf{w})$ and $R_{\mathcal{N}}(\mathbf{w})$. Now consider perturbation of a single variable—define $\mathcal{N}' = (\mathbf{x}^1, \dots, \mathbf{x}', \dots, \mathbf{x}^N)$ which matches \mathcal{N} everywhere except \mathbf{x}' . We have two cases:

(i) When \mathbf{x}' has corresponding label 0, we have

$$\begin{aligned} |\Phi_{\mathcal{N}}(\mathbf{w}) - \Phi_{\mathcal{N}'}(\mathbf{w})| &= \left| \frac{1}{2N_0N_1} \sum_{i=1}^{N_1} [1 - \mathbf{w}^\top (\mathbf{x}_i^1 - \mathbf{x}_j^0)]^2 - \frac{1}{2N_0N_1} \sum_{i=1}^{N_1} [1 - \mathbf{w}^\top (\mathbf{x}_i^1 - \mathbf{x}')]^2 \right| \\ &\leq \frac{1}{2N_0N_1} N_1 6X_*^2 W_*^2 + \frac{1}{2N_0N_1} N_1 4X_* W_* \end{aligned} \quad (9)$$

From this last line it follows that $|\Phi_{\mathcal{N}}(\mathbf{w}) - \Phi_{\mathcal{N}'}(\mathbf{w})| \leq C_2/N_0$.

(ii) Similarly when \mathbf{x}' has corresponding label 1: $|\Phi_{\mathcal{N}}(\mathbf{w}) - \Phi_{\mathcal{N}'}(\mathbf{w})| \leq C_2/N_1$.

As the differences are bounded and the inputs are independent, applying McDiarmid's inequality yields $P(|\Phi_{\mathcal{N}}(\mathbf{w})| \geq \epsilon/2) \leq 2 \exp\{-(\epsilon^2/2C_2^2)\rho(1-\rho)N\}$. This, along with $I = C_S(\epsilon/(2C_1))$, implies the bound in Lemma 1. \square

Proof of Theorem 1. By definition of $\mathbf{w}_{\mathcal{N}}$ and \mathbf{w}_* we have $R(\mathbf{w}_{\mathcal{N}}) - R(\mathbf{w}_*) \geq 0$. Also,

$$\begin{aligned} R(\mathbf{w}_{\mathcal{N}}) - R(\mathbf{w}_*) &= [R(\mathbf{w}_{\mathcal{N}}) - R_{\mathcal{N}}(\mathbf{w}_{\mathcal{N}}) + R_{\mathcal{N}}(\mathbf{w}_*) - R(\mathbf{w}_*)] + [R_{\mathcal{N}}(\mathbf{w}_{\mathcal{N}}) - R_{\mathcal{N}}(\mathbf{w}_*)] \\ &\leq |R(\mathbf{w}_{\mathcal{N}}) - R_{\mathcal{N}}(\mathbf{w}_{\mathcal{N}})| + |R(\mathbf{w}_*) - R_{\mathcal{N}}(\mathbf{w}_*)|. \end{aligned} \quad (10)$$

where the second line follows from $R_{\mathcal{N}}(\mathbf{w}_{\mathcal{N}}) - R_{\mathcal{N}}(\mathbf{w}_*) \leq 0$ and the triangle inequality. We next bound both terms by $\epsilon/2$, and according to Lemma 1, the bounds hold simultaneously with probability at least $1 - C_S(\epsilon/4C_1)2 \exp\{-(\epsilon^2/(8C_2^2))\rho(1-\rho)N\}$. \square

3 Low Cost Bipartite Ranking (LCBR)

In this section we derive risk bounds for LCBR, a subsampling strategy for approximately optimizing BBR. Our main goal is to obtain a bound similar to that in Eq. (4). We start with a brief comparison of LCBR and BBR: In Section 1 we noted that the empirical risk objective of BBR can be written as $R_{\mathcal{N}}(\mathbf{w}) = (1/2)\mathbf{w}^\top \mathbf{\Sigma}_N \mathbf{w} - \mathbf{\mu}_N^\top \mathbf{w}$ where $\mathbf{\Sigma}_N = 1/(N_1N_0) \sum_{i=1}^{N_1} \sum_{j=1}^{N_0} (\mathbf{x}_i^1 - \mathbf{x}_j^0)(\mathbf{x}_i^1 - \mathbf{x}_j^0)^\top$ and $\mathbf{\mu}_N = 1/(N_1N_0) \sum_{i=1}^{N_1} \sum_{j=1}^{N_0} (\mathbf{x}_i^1 - \mathbf{x}_j^0)$. Here $\mathbf{\Sigma}_N$ and $\mathbf{\mu}_N$ are constructed using all N_1N_0 pairs available. On the other hand, LCBR subsamples S pairs from the fixed

$N_1 N_0$ total pairs uniformly at random with replacement. The pairs obtained this way are denoted by the set $\mathcal{S} = \{(x_1^1, x_1^0), \dots, (x_S^1, x_S^0)\}$. It can be seen that the elements of \mathcal{S} are random variables sampled from a uniform distribution conditional on \mathcal{N} . Thus, while the elements of \mathcal{N} are sampled from the class-conditionals, i.e. \mathcal{D}^i , the elements of \mathcal{S} are sampled from the uniform distribution $\mathcal{D}(\mathcal{N})$. With a slight abuse of notation, we will denote this by $(\mathbf{x}_i^1, \mathbf{x}_i^0) \sim \mathcal{N}$. The corresponding objective of LCBR is $R_{\mathcal{S}}(\mathbf{w}) = (1/2)\mathbf{w}^\top \mathbf{\Sigma}_S \mathbf{w} - \boldsymbol{\mu}_S^\top \mathbf{w}$ with $\boldsymbol{\mu}_S$ and $\mathbf{\Sigma}_S$ the first and second moments computed on the subsample.

In order to derive a risk bound for $|R(\mathbf{w}_S) - R(\mathbf{w}_*)|$ we first need to bound $|R_{\mathcal{N}}(\mathbf{w}_N) - R_{\mathcal{N}}(\mathbf{w}_S)|$. However, in the latter expression, the weight vectors and the samples are not independent. For this reason, we again use a uniform convergence argument. Here we use matrix and vector concentration to obtain the bounds necessary. In particular, the following lemma provides concentration inequalities for two key variables.

Lemma 2. For $\mathbf{w}_1, \mathbf{w}_2 \in \mathcal{W}$, define $\Delta_{\mathbf{\Sigma}} := \mathbf{w}_1^\top (\mathbf{\Sigma}_S - \mathbf{\Sigma}_N) \mathbf{w}_1$ and $\Delta_{\sigma} := (\mathbf{w}_1 - \mathbf{w}_2)^\top (\boldsymbol{\mu}_N - \boldsymbol{\mu}_S)$. The following hold:

$$(i) \ P(\sup_{\mathbf{w}_1 \in \mathcal{W}} |\Delta_{\mathbf{\Sigma}}| \geq \epsilon) \leq 2D \exp\{-S\epsilon^2/(8\|\mathbf{\Sigma}_N\|_2 X_*^2 W_*^4 + (16/3)\epsilon X_*^2 W_*^2)\}$$

$$(ii) \ P(\sup_{\mathbf{w}_1, \mathbf{w}_2 \in \mathcal{W}} |\Delta_{\sigma}| \geq \epsilon) \leq 2 \exp\{-1/2(\epsilon\sqrt{S}/(2X_* W_*) - 1)^2\}$$

The proof is given in the appendix. We can now bound the difference of the empirical risks with high probability.

Theorem 2. For a given sample set \mathcal{N} and its subsample set \mathcal{S} let $\mathbf{w}_N = \arg \min_{\mathbf{w} \in \mathcal{W}} R_{\mathcal{N}}(\mathbf{w})$ and $\mathbf{w}_S = \arg \min_{\mathbf{w} \in \mathcal{W}} R_{\mathcal{S}}(\mathbf{w})$. If the subsample size satisfies

$$S \geq \max\left\{\log(4D/\delta) \frac{\|\mathbf{\Sigma}_N\|_2 X_*^2 W_*^4 + (1/3)\epsilon X_*^2 W_*^2}{\epsilon^2/32}, \frac{X_*^2 W_*^2}{\epsilon^2/4} \left[\sqrt{2\log(4/\delta)} + 1\right]^2\right\} \quad (11)$$

then

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{N}} \left(|R_{\mathcal{N}}(\mathbf{w}_S) - R_{\mathcal{N}}(\mathbf{w}_N)| \geq \epsilon \right) \leq 1 - \delta. \quad (12)$$

Proof. We define the following pointwise difference

$$\Delta(\mathbf{w}) := R_{\mathcal{S}}(\mathbf{w}) - R_{\mathcal{N}}(\mathbf{w}) = \frac{1}{2} \mathbf{w}^\top (\mathbf{\Sigma}_S - \mathbf{\Sigma}_N) \mathbf{w} + \mathbf{w}^\top (\boldsymbol{\mu}_N - \boldsymbol{\mu}_S). \quad (13)$$

Our first task is to show that the following inequality holds,

$$0 \leq R_{\mathcal{N}}(\mathbf{w}_S) - R_{\mathcal{N}}(\mathbf{w}_N) \leq \Delta(\mathbf{w}_N) - \Delta(\mathbf{w}_S). \quad (14)$$

The LHS of this inequality holds by the definition of \mathbf{w}_N and \mathbf{w}_S . For the RHS we have

$$R_{\mathcal{N}}(\mathbf{w}_S) - R_{\mathcal{N}}(\mathbf{w}_N) = \Delta(\mathbf{w}_N) - \Delta(\mathbf{w}_S) - [R_{\mathcal{S}}(\mathbf{w}_N) - R_{\mathcal{S}}(\mathbf{w}_S)] \leq \Delta(\mathbf{w}_N) - \Delta(\mathbf{w}_S) \quad (15)$$

since $R_S(\mathbf{w}_N) - R_S(\mathbf{w}_S) \geq 0$. It is therefore sufficient prove a high probability bound for $|\Delta(\mathbf{w}_N) - \Delta(\mathbf{w}_S)|$. Next,

$$\begin{aligned} |\Delta(\mathbf{w}_N) - \Delta(\mathbf{w}_S)| &= \left| \frac{1}{2} \mathbf{w}_N^\top (\boldsymbol{\Sigma}_S - \boldsymbol{\Sigma}_N) \mathbf{w}_N - \frac{1}{2} \mathbf{w}_S^\top (\boldsymbol{\Sigma}_S - \boldsymbol{\Sigma}_N) \mathbf{w}_S + (\mathbf{w}_S - \mathbf{w}_N)^\top (\boldsymbol{\mu}_N - \boldsymbol{\mu}_S) \right| \\ &\leq \left| \frac{1}{2} \mathbf{w}_N^\top (\boldsymbol{\Sigma}_S - \boldsymbol{\Sigma}_N) \mathbf{w}_N - \frac{1}{2} \mathbf{w}_S^\top (\boldsymbol{\Sigma}_S - \boldsymbol{\Sigma}_N) \mathbf{w}_S \right| + \left| (\mathbf{w}_S - \mathbf{w}_N)^\top (\boldsymbol{\mu}_N - \boldsymbol{\mu}_S) \right|. \end{aligned} \quad (16)$$

We bound each of these terms by $\epsilon/2$ with probability at least $1 - \delta/2$. Note that both terms have weight vectors coupled with samples so we cannot apply concentration inequalities directly. However, they can be upper bounded by the expressions in Lemma 2. In particular, for the quadratic term $\left| \frac{1}{2} \mathbf{w}_N^\top (\boldsymbol{\Sigma}_S - \boldsymbol{\Sigma}_N) \mathbf{w}_N - \frac{1}{2} \mathbf{w}_S^\top (\boldsymbol{\Sigma}_S - \boldsymbol{\Sigma}_N) \mathbf{w}_S \right| \leq \sup_{\mathbf{w}_1, \mathbf{w}_2 \in \mathcal{W}} |\Delta_{\boldsymbol{\Sigma}}|$. We then apply Lemma 5(i) with threshold $\epsilon/2$ and probability $\delta/2$ which yields the first term in Eq. (11). For the linear term, $\left| (\mathbf{w}_S - \mathbf{w}_N)^\top (\boldsymbol{\mu}_N - \boldsymbol{\mu}_S) \right| \leq \sup_{\mathbf{w}_1, \mathbf{w}_2 \in \mathcal{W}} |\Delta_{\boldsymbol{\sigma}}|$. Applying Lemma 5(ii) with threshold $\epsilon/2$ and probability $\delta/2$ yields the second term. \square

Theorem 2 shows that the subsample size S required to decrease the difference $|R_{\mathcal{N}}(\mathbf{w}_S) - R_{\mathcal{N}}(\mathbf{w}_N)|$ does not depend on the number of total pairs $N_1 N_0$. Instead, the dependence is on the operator norm of the empirical second moment matrix $\|\boldsymbol{\Sigma}_N\|_2$, and polynomial in X_* , W_* , $\log(1/\delta)$, and $1/\epsilon$. This is favorable, as in many settings the total number of pairs can be prohibitively large. The following is the main result of this section, regarding the actual risk of LCBR solution.

Theorem 3. For a probability target p^* let the sample size be chosen such that

$$S \geq \max \left\{ \log(4D/p^*) \frac{\|\boldsymbol{\Sigma}_N\|_2 X_*^2 W_*^4 + (1/15)\epsilon X_*^2 W_*^2}{\epsilon^2/800}, \frac{X_*^2 W_*^2}{\epsilon^2/100} \left[\sqrt{2 \log(4/p^*)} + 1 \right]^2 \right\}. \quad (17)$$

Then the solution \mathbf{w}_S returned by LCBR satisfies

$$\mathbb{P}_{\mathcal{N} \sim \mathcal{D}^N} \left(R(\mathbf{w}_S) - R(\mathbf{w}_*) \geq \epsilon \right) \leq 2 C_S \left(\frac{\epsilon}{10C_1} \right) \exp \left\{ -\frac{\epsilon^2}{50C_2^2} \rho(1 - \rho)N \right\} + p^*. \quad (18)$$

Proof. We start with the inequality

$$R(\mathbf{w}_S) - R(\mathbf{w}_*) = |R(\mathbf{w}_S) - R(\mathbf{w}_*)| \leq |R(\mathbf{w}_S) - R(\mathbf{w}_N)| + |R(\mathbf{w}_N) - R(\mathbf{w}_*)|. \quad (19)$$

The first term can be bounded as

$$\begin{aligned} |R(\mathbf{w}_S) - R(\mathbf{w}_N)| &= |R(\mathbf{w}_S) - R_{\mathcal{N}}(\mathbf{w}_S) + R_{\mathcal{N}}(\mathbf{w}_S) - R_{\mathcal{N}}(\mathbf{w}_N) + R_{\mathcal{N}}(\mathbf{w}_N) - R(\mathbf{w}_N)| \\ &\leq |R(\mathbf{w}_S) - R_{\mathcal{N}}(\mathbf{w}_S)| + |R(\mathbf{w}_N) - R_{\mathcal{N}}(\mathbf{w}_N)| + |R_{\mathcal{N}}(\mathbf{w}_S) - R_{\mathcal{N}}(\mathbf{w}_N)|. \end{aligned} \quad (20)$$

Note that the empirical risk trick in the proof Theorem 1 no longer applies as it not known which one of $R(\mathbf{w}_N)$ or $R(\mathbf{w}_S)$ is smaller, and all three terms have to be retained. On the other hand, for the second term, from the proof of Theorem 1 we know that

$$|R(\mathbf{w}_N) - R(\mathbf{w}_*)| \leq |R(\mathbf{w}_N) - R_{\mathcal{N}}(\mathbf{w}_N)| + |R(\mathbf{w}_*) - R_{\mathcal{N}}(\mathbf{w}_*)|. \quad (21)$$

Combining Eqs. (20) and (21) we get

$$R(\mathbf{w}_S) - R(\mathbf{w}_*) \leq \left[\sum_{\mathbf{w} \in \mathcal{W}_4} |R(\mathbf{w}) - R_{\mathcal{N}}(\mathbf{w})| \right] + |R_{\mathcal{N}}(\mathbf{w}_S) - R_{\mathcal{N}}(\mathbf{w}_N)| . \quad (22)$$

where we defined the sequence $\mathcal{W}_4 = [\mathbf{w}_*, \mathbf{w}_N, \mathbf{w}_N, \mathbf{w}_S]$. We now consider bounding each term by $\epsilon/5$. For the summation on the right hand side this yields

$$\begin{aligned} P \left(\sum_{\mathbf{w} \in \mathcal{W}_4} |R(\mathbf{w}) - R_{\mathcal{N}}(\mathbf{w})| \leq 4\epsilon/5 \right) &\leq P \left(\sup_{\mathbf{w} \in \mathcal{W}} |R(\mathbf{w}) - R_{\mathcal{N}}(\mathbf{w})| \leq \epsilon/5 \right) \\ &\leq 2 C_S \left(\frac{\epsilon}{10C_1} \right) \exp \left\{ -\frac{\epsilon^2}{50C_2^2} \rho(1-\rho)N \right\} \end{aligned} \quad (23)$$

where the last inequality follows from Lemma 1. For the last term we would like to bound the deviation by $\epsilon/5$ with probability p^* . Plugging these terms into the sample complexity bound of Theorem 2 yields the bound of Eq. (17). \square

Theorem 3 shows that the number of samples S required to make LCBR competitive with BBR does not depend on $N_1 N_0$. Firstly, for a fixed probability target p^* , S depends on \mathcal{N} only through Σ_N . In practice we set p^* such that it is relatively smaller than the exponential term in Eq. (18). In this case S will have an implicit dependence on $N_1 + N_0$ and ρ . The important difference is, while BBR requires $N_1 N_0$ samples to achieve the bound in Eq. (4), LCBR can achieve the comparable bound in Eq. (18) with $S \ll N_1 N_0$. We demonstrate this with experiments in Section 5. Finally, note that we focused on the linear case due to space limitations. Clearly, for any *fixed* nonlinear feature map, our results hold where D is replaced with the dimensionality of the new feature space. Another direction is to consider random feature transforms and provide bounds based to the best ranker in the corresponding function space [12]; we leave this for a longer version of the paper.

4 Related Work

Our proof of the uniform risk bound in Lemma 1 is based on covering numbers, which is also used for analyzing linear regression problems [10]. The covering number-based argument is later extended to online learning with pairwise loss functions [13]; however, their analysis is for sequential updates. Bounds based on Rademacher complexity and U-processes are considered in [9, 14, 15, 16]. These bounds can be tighter, however they are based on the assumption that all samples are drawn i.i.d. from \mathcal{D} , which is different from our setup. Bounds based on algorithmic stability [17, 18] and VC dimension [11] have also been considered. Replacing the discrete AUC loss with a convex one has been investigated in a number of work. In particular [4, 5, 9, 8, 7] consider online algorithms. Online learning based on stochastic saddle point problems is considered in [19, 20]. It is also worthwhile to note that the majority of the aforementioned papers are based on the pairwise squared loss. In [6], the consistency of surrogate loss functions with respect to AUC loss has been considered, extending the results of [21]. On the other hand

[22, 23] provide bounds for the AUC loss in terms of the surrogate loss of the learner. Kernel based methods for bipartite ranking was considered in [3, 24, 25]. (For more related work see also [26, 27, 28].) In terms of low sample complexity, the stochastic gradient descent (SGD) based online learning algorithms also provide bounds in terms of the samples used; however these algorithms require a step size to tune, whereas LCBR in Algorithm 2 does not need this parameter, which can be an important practical advantage. As we will show in experiments, LCBR can achieve better performance with fewer samples, compared to the SGD-based methods.

5 Experiments

5.1 Gaussian Mixture Distribution

We first consider experiments where the generating distribution is a K -component Gaussian Mixture. So for $i \in \{0, 1\}$ we can write

$$\mathcal{D}^i(\mathbf{x}) = \sum_{k=1}^K \frac{c_k}{\sqrt{2\pi\sigma^2}} \exp\{-\|\mathbf{x} - \boldsymbol{\mu}_k^i\|_2^2 / (2\sigma^2)\} \quad (24)$$

where c_k 's are mixture weights. We consider the case where the covariance matrix is isotropic and controlled with a single scale parameter σ . We consider three cases where $K = 1, 2, 3$ and $\sigma = 2, 3, 4$ where the increasing value of K and σ makes the problem gradually more difficult. The mean values for class one are sampled from the unit cube in positive orthant and class zero from unit cube in negative orthant. The curves are obtained by averaging 50 experiments. For any given pair of conditional distributions the optimal ranking function minimizing the ϕ -risk is found by calculating $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ and solving the ϕ -risk minimization problem. On the other hand, the optimal ranking function maximizing the AUC is the likelihood ratio of class-conditionals, which is a consequence of the Neyman-Pearson lemma [29]. Figure 1 shows the results of our experiments. For panels (a)-(f) we have $\rho = 0.5$ with 10^3 positive and negative samples. So BBR uses 10^6 pairs whereas LCBR uses the number shown on the x-axis. In panels (a)-(c) we show the ϕ -risk and in (d)-(f) we show the AUC. As the number of components increase the problem becomes more difficult, resulting in higher ϕ -risk and lower AUC for optimal ranking functions. In all cases BBR is very close to the optimal ranker; but it has high sample cost. LCBR, on the other hand, catches on with $S = 3000$ at most. This is also reflected in wall clock times; as shown in panels (g)-(i) LCBR is roughly 100 times faster than BBR when $S = 5000$ subsamples are used. Finally we illustrate the performance of BBR and LCBR as a function of label skew ρ . In the proof of Lemma 1, we saw that applying McDiarmid's inequality gives the term $\rho(1 - \rho)$ in the exponent. This suggests, as ρ goes to 0 or 1 the generalization performance should decrease. Panel (j) shows that this is indeed the case: here BBR uses 2000ρ positive and $2000(1 - \rho)$ negative samples, whereas LCBR uses 5,000 subsamples. The performance of BBR and LCBR are once again close; but as ρ deviates from 0.5 the generalization of both algorithms get worse.

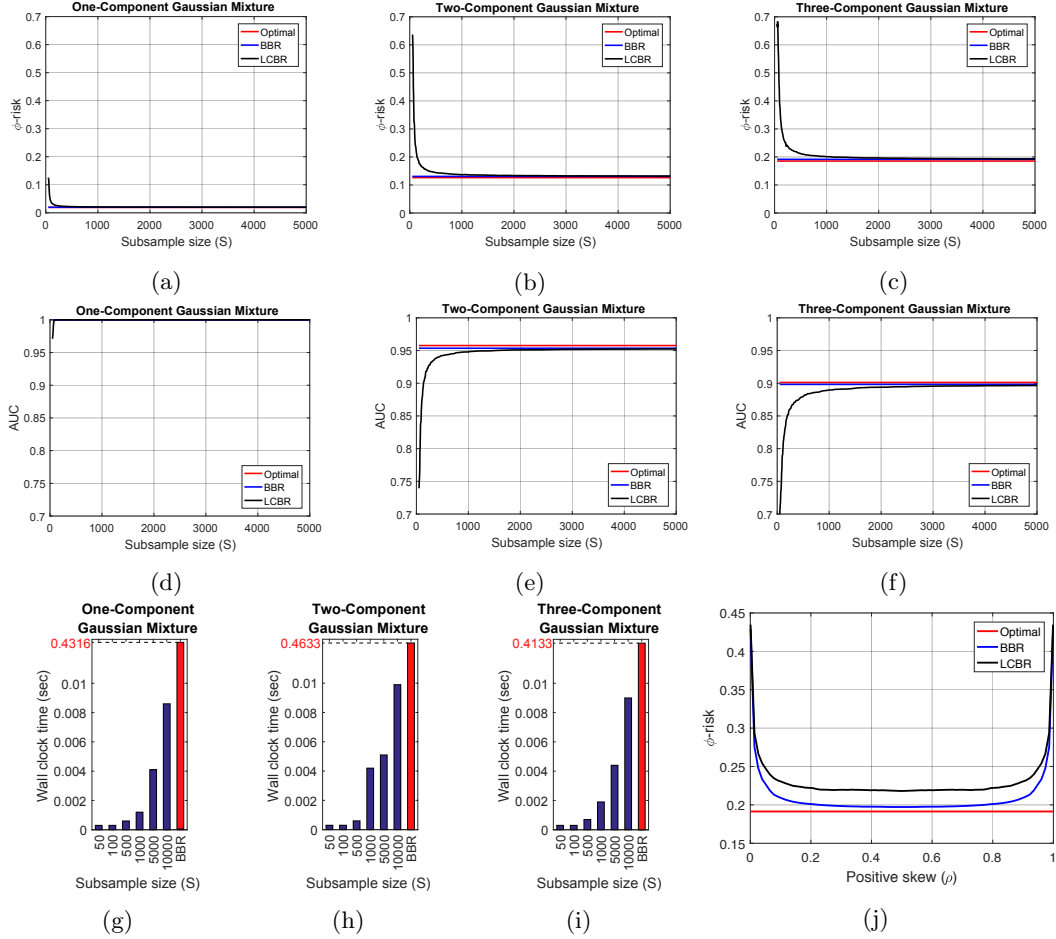


Figure 1: Results of experiments with Gaussian Mixture distributions.

5.2 LIBSVM Datasets

We now compare LCBR against state-of-the-art algorithms on three datasets from LIBSVM. The algorithms we implement, in addition to LCBR, are the following: AdaOAM [8] uses adaptive stochastic gradient descent with pairwise squared loss, whereas PGD is SGD based approach [7]. SPAM is one of the most recent works where AUC optimization is formulated as a stochastic saddle point problem [19]. OAM is a relatively older algorithm, but we include it as it is based in a different objective, the pairwise hinge loss [4]. As widely done in the literature, we use AUC on test set as the performance measure. We use the train and test splits provided by LIBSVM. Regularization parameters are determined by cross-validation and step sizes are chosen based on the references. The experiments are averaged over 50 runs.

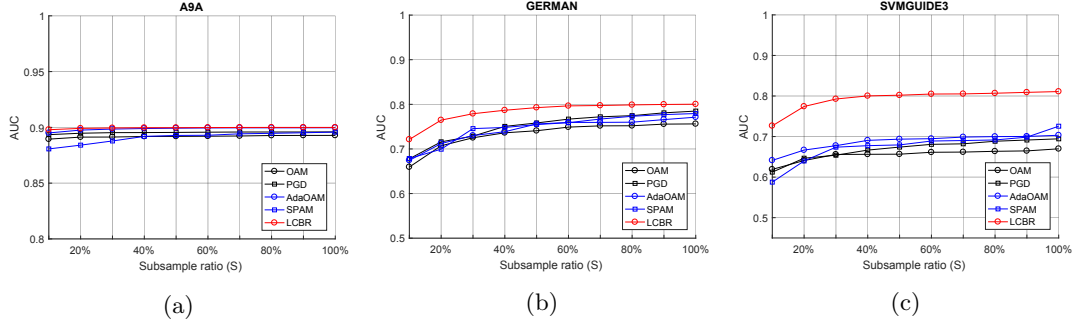


Figure 2: Results of experiments with LIBSVM datasets.

We show the results in Figure 2. The x-axis correspond to the number of subsamples used as a percentage of total number of samples available. For example, **A9A** dataset contains approximately $33K$ data points. A subsample ratio of 50% implies we use approximately $16.5K$ random samples from this dataset. For the **A9A** dataset all algorithms have good generalization, with an AUC of approximately 90%. For the **GERMAN** and **SVMGUIDE3** datasets, on the other hand, there is a significant gap between LCBR and the SGD-based competitors. Here the SGD and learning rate free approach of LCBR proves useful and it achieves better generalization for the given number of samples. Therefore the sample complexity of LCBR is favorable.

5.3 Algorithmic Complexity

For the algorithms presented in this paper, it can be seen that the computational bottleneck is at the accumulation phase (cf. Algorithms 1 and 2). For BBR the cost of this step is $O(N_1 N_0 D^2)$ and for LCBR this is $O(S D^2)$. Typically $S \ll N_1 N_0$ which can save significant computation. On the other hand, the storage for both BBR and LCBR is $O(D^2)$ as it requires storing the covariance matrix. This quadratic dependence on dimension can be an issue when D is too large. In this case, a sparse approximation of Σ_S can be necessary. However, note that this storage bottleneck is also present in the algorithms we compared [4, 7, 8]. In contrast, SPAM [19] can be implemented in $O(D)$ space; however that algorithm requires knowledge of expectations with respect to the true conditionals, which is unknown in practice. Our experiments show that LCBR can still achieve better performance.

6 Conclusion

We have considered the problem of bipartite ranking, where the empirical AUC loss is replaced with the pairwise squared loss. Different from the previous work—which was based on SGD—we proposed a low sample cost bipartite ranking algorithm (LCBR), and showed that the number

of samples required for good performance obeys $S \ll N_1 N_0$. Experiments show that LCBR quickly achieves similar performance with BBR where the number of samples are several order of magnitudes lower. Experiments against state-of-the-art bipartite ranking algorithms also show that LCBR can achieve better generalization with a smaller subsample set. In a longer version of the paper we will also consider extending these results to random feature spaces, which include random kernel features and random neural networks.

7 Appendix: Proof of Lemma 2

For the proofs we will use the shorthand $\mathbf{x}_s = \mathbf{x}_{i_s}^1 - \mathbf{x}_{j_s}^0$.

(i) We recall the Matrix Bernstein Inequality for a $D \times D$ symmetric, random matrix \mathbf{Z} and threshold γ :

$$P(\|\mathbf{Z}\|_2 > \gamma) \leq 2D \exp\left(-\frac{\gamma^2/2}{\mathbb{V}(\mathbf{Z}) + L\gamma/3}\right) \quad (25)$$

where L is a norm bound on summands.

First apply a spectral norm bound to the quadratic expression: $\sup_{\mathbf{w}_1 \in \mathcal{W}} |\Delta_{\Sigma}| \leq W_*^2 \|\Sigma_S - \Sigma_N\|_2$. We now take $\mathbf{Z} = \Sigma_S - \Sigma_N$. The spectral norm can be bounded based on the argument in [30]. We can decompose \mathbf{Z} into a sum: $\mathbf{Z} = \sum_{s=1}^S (1/S)[\mathbf{x}_s \mathbf{x}_s^\top - \Sigma_N]$. We denote each summand by $\mathbf{E}_s = (1/S)[\mathbf{x}_s \mathbf{x}_s^\top - \Sigma_N]$. It then follows from triangle inequality that

$$\|\mathbf{E}_s\|_2 \leq \frac{1}{S} [\|\mathbf{x}_s \mathbf{x}_s^\top\|_2 + \|\Sigma_N\|_2] \leq \frac{2}{S} \|\mathbf{x}_s\|_2^2 \leq \frac{8X_*^2}{S} \quad (26)$$

where the second inequality follows from Jensen's inequality. Since the subsampling is conditional on \mathcal{N} , Σ_N is constant with respect to the subsampled pairs, and each summand is centered and i.i.d. The variance of the sum decomposes as $\mathbb{V}(\mathbf{Z}) = \|\sum_{s \in S} \mathbb{E}[\mathbf{E}_s^2]\|_2$. For a single summand the second moment can be bounded as

$$\mathbb{E}[\mathbf{E}_s^2] = \frac{1}{S^2} \mathbb{E}[(\mathbf{x}_s \mathbf{x}_s^\top - \Sigma_N)^2] = \frac{1}{S^2} \left[\mathbb{E}[\|\mathbf{x}_s\|_2^2 \mathbf{x}_s \mathbf{x}_s^\top] - \Sigma_N^2 \right] \preceq \frac{4X_*^2}{S^2} \Sigma_N \quad (27)$$

from which the variance inequality $\mathbb{V}(\mathbf{Z}) \leq (4X_*^2/S) \|\Sigma_N\|_2$ follows. Substituting these to Eq. (25) with $\gamma = \epsilon/W_*^2$ yields the result.

(ii) We recall the following concentration inequality for i.i.d. and bounded random vectors with mean $\bar{\mathbf{x}}$ [12]:

$$P\left(\left\|\frac{1}{S} \sum_{s=1}^S \mathbf{x}_s - \bar{\mathbf{x}}\right\|_2 \geq \gamma\right) \leq \exp\left\{-\frac{1}{2} \left(\frac{\gamma\sqrt{S}}{L} - 1\right)^2\right\}. \quad (28)$$

From the following inequalities

$$\sup_{\mathbf{w}_1, \mathbf{w}_2 \in \mathcal{W}} |\Delta_\sigma(\mathbf{w}_1, \mathbf{w}_2)| \leq \sup_{\mathbf{w}_1, \mathbf{w}_2} \|\mathbf{w}_1 - \mathbf{w}_2\|_S \|\boldsymbol{\mu}_N - \boldsymbol{\mu}_S\|_2 \leq 2W_* \left\| \frac{1}{S} \sum_{s=1}^S \mathbf{x}_s - \boldsymbol{\mu}_N \right\|_2 \quad (29)$$

the desired result is obtained by setting $\gamma = \epsilon/2W_*$ and $L = X_*$ in Eq. (28). \square

References

- [1] H. B. McMahan, G. Holt, D. Sculley, M. Young, D. Ebner, J. Grady, L. Nie, T. Phillips, E. Davydov, D. Golovin, S. Chikkerur, D. Liu, M. Wattenberg, A. M. Hrafnkelsson, T. Boulos, and J. Kubica, “Ad click prediction: a view from the trenches,” in *Conference on Knowledge Discovery and Data Mining*, 2013.
- [2] J. Hu, H. Yang, M. Lyu, I. King, and A. So, “Online Nonlinear AUC Maximization for Imbalanced Data Sets,” *IEEE Transactions on Neural Networks and Learning Systems*, 2016.
- [3] X. Hong, S. Chen, and C. J. Harris, “A Kernel-Based Two-Class Classifier for Imbalanced Data Sets,” *IEEE Transactions on Neural Networks*, 2007.
- [4] P. Zhao, S. C. H. Hoi, R. Jin, and T. Yang, “Online AUC Maximization,” in *International Conference on Machine Learning*, 2011.
- [5] W. Gao, R. Jin, S. Zhu, and Z.-H. Zhou, “One-Pass AUC Optimization,” in *International Conference on Machine Learning*, 2013.
- [6] W. Gao and Z.-H. Zhou, “On the Consistency of AUC Pairwise Optimization,” in *International Joint Conference on Artificial Intelligence*, 2015.
- [7] M. Boissier, S. Lyu, Y. Ying, and D.-X. Zhou, “Fast Convergence of Online Pairwise Learning Algorithms,” in *International Conference on Artificial Intelligence and Statistics*, 2016.
- [8] Y. Ding, P. Zhao, S. C. H. Hoi, and Y. S. Ong, “An Adaptive Gradient Method for Online AUC Maximization,” in *Association for Advancement of Artificial Intelligence*, 2015.
- [9] P. Kar, B. Sriperumbudur, P. Jain, and H. Karnick, “On the Generalization Ability of Online Learning Algorithms for Pairwise Loss Functions,” in *International Conference on Machine Learning*, 2013.
- [10] M. Mohri, A. Rostamizadeh, and A. Talwalkar, *Foundations of Machine Learning*. MIT Press, 2012.
- [11] S. Agarwal, T. Graepel, R. Herbrich, S. Har-Peled, and D. Roth, “Generalization Bounds for the Area Under the ROC Curve,” *Journal of Machine Learning Research*, 2005.

- [12] A. Rahimi and B. Recht, “Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning,” in *Advances in Neural Information Processing Systems*, 2008.
- [13] Y. Wang, R. Khairon, D. Pechyony, and R. Jones, “Generalization Bounds for Online Learning Algorithms with Pairwise Loss Functions,” in *Conference on Learning Theory*, 2012.
- [14] S. Clemencon, G. Lugosi, and N. Vayatis, “Ranking and empirical minimization of U-Statistics,” *The Annals of Statistics*, 2008.
- [15] T. Peel, A. Sandrine, and L. Ralaivola, “Empirical Bernstein Inequalities for U-Statistics,” in *Advances in Neural Information Processing Systems*, 2010.
- [16] Y. Lei, S.-B. Lin, and K. Tang, “Generalization Bounds for Regularized Pairwise Learning,” in *International Joint Conference on Artificial Intelligence*, 2018.
- [17] S. Agarwal and P. Niyogi, “Generalization Bounds for Ranking Algorithms via Algorithmic Stability,” *Journal of Machine Learning Research*, 2009.
- [18] C. Cortes, M. Mohri, and A. Rastogi, “Magnitude-Preserving Ranking Algorithms,” in *International Conference on Machine Learning*, 2007.
- [19] M. Natole Jr., Y. Ying, and S. Lyu, “Stochastic Proximal Algorithms for AUC Maximization,” in *International Conference on Machine Learning*, 2018.
- [20] Y. Ying, L. Wen, and S. Lyu, “Stochastic Online AUC Maximization,” in *Advances in Neural Information Processing Systems*, 2016.
- [21] P. L. Bartlett, M. I. Jordan, and J. D. McAuliffe, “Convexity, classification, and risk bounds,” *Journal of the American Statistical Association*, 2006.
- [22] S. Agarwal, “Surrogate Regret Bounds for Bipartite Ranking via Strongly Proper Losses,” *Journal of Machine Learning Research*, 2014.
- [23] W. Kotlowski, K. J. Dembczynski, and E. Huellermeier, “Bipartite ranking through minimization of univariate loss,” in *International Conference on Machine Learning*, 2011.
- [24] Y. Ding, C. Liu, P. Zhao, and S. C. Hoi, “Large Scale Kernel Methods for Online AUC Maximization,” in *International Conference on Data Mining (ICDM)*, 2017.
- [25] J. Hu, H. Yang, M. Lyu, I. King, and A. So, “Kernelized Online Imbalanced Learning with Fixed Budgets,” in *Association for Advancement of Artificial Intelligence (AAAI)*, 2015.
- [26] H. Narasimhan and S. Agarwal, “On the Relationship Between Binary Classification, Bipartite Ranking, and Binary Class Probability Estimation,” in *Advances in Neural Information Processing Systems*, 2013.
- [27] S. Agarwal, “Generalization bounds for some ordinal regression algorithms,” in *Proceedings of the 19th International Conference on Algorithmic Learning Theory*, 2008.

- [28] A. Rajkumar and S. Agarwal, “When Can We Rank Well from Comparisons of $O(n \log n)$ Non-Actively Chosen Pairs?” in *Conference on Learning Theory*, 2016.
- [29] H. V. Poor, *An Introduction to Signal Detection and Estimation*. Springer, 1998.
- [30] J. Tropp, “An introduction to matrix concentration inequalities,” *Foundations and Trends in Machine Learning*, 2015.