

Stochastic Hard Thresholding Algorithms for AUC Maximization

Zhenhuan Yang^{*}, Baojian Zhou[†], Yunwen Lei[‡], Yiming Ying^{*}

^{*}University at Albany, Albany, USA

[†]Stony Brook University, Stony Brook, USA

[‡]School of Computer Science, University of Birmingham, Birmingham, UK

Email: zyang6@albany.edu, baojian.zhou@cs.stonybrook.edu, yunwen.lei@htomail.com, yying@albany.edu

Abstract—In this paper, we aim to develop stochastic hard thresholding algorithms for the important problem of AUC maximization in imbalanced classification. The main challenge is the pairwise loss involved in AUC maximization. We overcome this obstacle by reformulating the U-statistics objective function as an empirical risk minimization (ERM), from which a stochastic hard thresholding algorithm (SHT-AUC) is developed. To our best knowledge, this is the first attempt to provide stochastic hard thresholding algorithms for AUC maximization with a per-iteration cost $\mathcal{O}(bd)$ where d and b are the dimension of the data and the minibatch size, respectively. We show that the proposed algorithm enjoys the linear convergence rate up to a tolerance error. In particular, we show, if the data is generated from the Gaussian distribution, then its convergence becomes slower as the data gets more imbalanced. We conduct extensive experiments to show the efficiency and effectiveness of the proposed algorithms.

Index Terms—Area Under the ROC Curve (AUC), sparse learning, stochastic hard thresholding, imbalanced classification

I. INTRODUCTION

Recently, there are a considerable amount of work on developing efficient algorithms for optimizing the Area under the ROC curve (AUC) score. It is a widely used performance measure for imbalanced data classification [1]–[4] which arises from applications including anomaly detection, information retrieval to cancer diagnosis. In particular, the work [4] showed that the AUC score is, in general, a better measure than accuracy for evaluating the predictive performance of many data mining algorithms.

In particular, the work by [5], [6] employed the cutting plane method and gradient descent algorithm, respectively. [7] developed the Nesterov’s accelerated gradient algorithms [8] for optimizing the multivariate performance measures [5]. The work of [9] used ideas of active learning to design heuristic algorithms for AUC maximization. Such optimization algorithms train the model on the whole training data which are not scale well to the high-dimensional data. Stochastic gradient descent (SGD) algorithms are widely used for high dimensional and large-scale data analysis due to its cheap per-iteration cost. In this aspect, variants of stochastic (online) gradient descent algorithms have been developed for AUC maximization. Specifically, [10]–[12] proposed to a variant of online (projected) gradient descent method. At time t , these methods need to compare the current example with previous ones, which have high per-iteration $\mathcal{O}(td)$ for d -dimensional

data. There are various techniques such as using the buffering set to alleviate the bottleneck but the size of the buffer set needs to be large in order to guarantee a good generalization. The appealing work by [13] observed in the case of the least square loss that the updates of these algorithms only rely on the covariance matrix where the per-iteration cost is of $\mathcal{O}(d^2)$. For high dimensional data, it used an appealing low-rank matrix to approximate the covariance matrix in order to reduce the per-iteration costs which may not be an ideal solution. Hence, such algorithms have an expensive per-iteration cost, making them not amenable for high dimensional data analysis. There are some recent work on nonlinear AUC maximization methods such as [14]–[16]. Recently, [17] reformulated the AUC maximization problem as a stochastic saddle point (min-max) problem (e.g. [18]), from which a stochastic primal-dual gradient algorithm was proposed. This algorithm successfully reduced the per-iteration cost to $\mathcal{O}(d)$. [19] followed this saddle point formulation for AUC maximization with ℓ_1 constraints and proposed a fast multi-stage SGD algorithm. In [20], fast SGD-type algorithms were developed for more general strongly convex regularization.

For high-dimensional data analysis, an underlying hypothesis is the sparsity of the data representation [21]–[25]. To obtain a sparse solution, many algorithms have been developed among which the prominent one is based on variants of ℓ_1 -norm constraints (regularization) which includes group lasso [26], [27], tree structured group lasso [28], [29] etc. Such approaches are convex and can be solved efficiently by convex optimization. Concurrently, sparse learning for AUC maximization has been developed in [19], [20], [30] using ℓ_1 regularization, where stochastic primal-dual gradient-type algorithms (SGD) have been developed. However, as many researchers observed [31]–[33], ℓ_1 -based stochastic algorithms are appealing convex approach which may be hard to preserve a truly sparse solution. In contrast, the greedy pursuit based on the sparse ℓ_0 constraints can recover the sparse structure well, among which the most prominent one is gradient hard thresholding [34], [35], [35]–[38]. In addition, compared to the convex ℓ_1 -norm based methods, hard thresholding algorithms are always orders of magnitude computationally more efficient for large-scale problems [39].

However, the existing hard thresholding algorithms are developed for the classical regression and classification where

the loss is pointwise, i.e. it depends on one data point. These algorithms can not directly apply to the setting of AUC maximization as its objective function is in the form of U-statistics [40] where the pairwise loss function depends on a pair of data points.

In this paper, we aim to develop stochastic hard thresholding algorithm for the problem of AUC maximization in imbalanced classification. The main challenge is the pairwise loss involved in AUC maximization. We overcome this obstacle by leveraging the ideas from [17], [20] by reformulating the U-statistics objective function as a standard empirical risk minimization (ERM). In particular, the reformulated AUC objective does not necessarily possess the strong convexity property as a whole. Instead, it is assumed that the objective function obeys the restricted strong convexity and restricted smoothness (RCS/RSS) [41], [42]. The main contribution of the paper is summarized as follows

- We reformulate the empirical AUC objective in the form of U-statistics as an ERM objective, from which a stochastic hard thresholding algorithm (referred to as SHT-AUC) is developed. To our best knowledge, this is the first attempt to provide stochastic hard thresholding algorithms for AUC maximization with a per-iteration cost $\mathcal{O}(bd)$ where d and b are the dimension of the data and the size of minibatch, respectively.

- We show that the proposed algorithm enjoys the linear convergence up to a tolerance error under RCS/RSS properties. We then characterize the RCS/RSS properties in AUC context. In particular, we show, if the data is generated from the Gaussian distribution, that its convergence becomes slower as the data gets more imbalanced, i.e. the imbalance ratio is getting smaller.

- We conduct extensive experiments to validate the proposed algorithm (SHT-AUC) on both simulated and real-world datasets. Our experiments show that the proposed algorithms outperform the existing algorithms in terms of AUC score and the ability of selecting meaningful features.

Outline of the paper. The rest of this paper are organized as follows. In Section 2, we reformulate the objective function of AUC maximization, and present the Stochastic Hard Thresholding Algorithm for AUC maximization (i.e. SHT-AUC). In Section 3, we present its convergence rate and discuss the implication of the theoretical results. In Section 4, we perform experiments on both simulation and real-world datasets to validate the proposed algorithm.

The detailed proofs for the theoretical results, and the source code of all methods and datasets are available at <https://github.com/baojianzhou/sparse-auc>.

II. PROBLEM FORMULATION AND PROPOSED ALGORITHM

In this section, we introduce necessary notations, formulate the problems of AUC maximization, and present the stochastic hard thresholding algorithm for AUC maximization (SHT-AUC).

A. Preliminaries

Given an integer $n \geq 1$, we define $[n] = \{1, \dots, n\}$. The standard Euclidean norm of vector $\mathbf{v} = (v_1, \dots, v_d)^\top \in \mathbb{R}^d$ is denoted by $\|\mathbf{v}\|_2 = \sqrt{\sum_{i=1}^d v_i^2}$. For any $\mathbf{v}, \mathbf{w} \in \mathbb{R}^d$, the inner product is given by $\langle \mathbf{v}, \mathbf{w} \rangle = \sum_{i=1}^d v_i w_i$. The support set of \mathbf{v} , i.e. indices of non-zeros, is denoted by $\text{supp}(\mathbf{v})$ whose cardinality is written as $\|\mathbf{v}\|_0$. For any integer $d > 0$, suppose that Ω is a subset of $[d]$. Then for any vector $\mathbf{v} \in \mathbb{R}^d$, we define $\mathcal{P}_\Omega(\cdot)$ as the orthogonal projection to the support set Ω which is defined by $(\mathcal{P}_\Omega(\mathbf{v}))_i = v_i$ if $i \in \Omega$ and 0 otherwise. In particular, let Γ be the support set indexing the k largest absolute components of \mathbf{v} . In this way, the hard thresholding operator is given by

$$\mathcal{H}_k(\mathbf{v}) = \mathcal{P}_\Gamma(\mathbf{v}). \quad (1)$$

Let \mathcal{X} be a domain in \mathbb{R}^d and $\mathcal{Y} = \{\pm 1\}$. Assume that the training data $\mathcal{S} = \{\mathbf{z}_i = (\mathbf{x}_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i \in [n]\}$ is drawn i.i.d from an unknown distribution on $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$. For each $1 \leq i \leq n$, if $y_i = 1$ we say \mathbf{z}_i is a positive example otherwise it is a negative example. Let n_+ denote the number of positive examples and n_- denote the number of negative examples, and define $r = \frac{n_+}{n_-}$ as the *imbalanced ratio*. Without loss of generality, we assume $n_- \geq n_+$, i.e. $r \leq 1/2$.

Definition of AUC. AUC score [3], [43] measures the probability for a randomly drawn positive instance to have a higher decision value than a randomly sampled negative instance. Specifically, for any \mathbf{w} , the AUC score on the data \mathcal{S} is defined by

$$\text{AUC}(\mathbf{w}) = \frac{1}{n_+ n_-} \sum_{i,j=1}^n \mathbb{I}_{[\mathbf{w}^\top(x_i - x_j) > 0]} \mathbb{I}_{[y_i=1]} \mathbb{I}_{[y_j=-1]}, \quad (2)$$

where $\mathbb{I}_{[\cdot]}$ is the indicator function which is 1 for the true event and 0 otherwise. The higher the AUC score is, the better performance of the linear function parametrized by \mathbf{w} will be. Maximizing the AUC score is equivalent to minimizing $1 - \text{AUC}(\mathbf{w}) = \frac{1}{n_+ n_-} \sum_{i,j=1}^n \mathbb{I}_{[\mathbf{w}^\top(x_i - x_j) \leq 0]} \mathbb{I}_{[y_i=1]} \mathbb{I}_{[y_j=-1]}$. In practice, the discontinuous indicator function $\mathbb{I}_{[\mathbf{w}^\top(x_i - x_j) \leq 0]}$ is replaced by a relaxed convex function. As done in [13], [17], [19], [20], in this paper we restrict our attention to the least square loss, i.e. replacing $\mathbb{I}_{[\mathbf{w}^\top(x_i - x_j) \leq 0]}$ by $(1 - \mathbf{w}^\top(x_i - x_j))^2$.

Now sparse AUC maximization with ℓ_0 constraints is given by

$$\min_{\|\mathbf{w}\|_0 \leq k} \frac{1}{n_+ n_-} \sum_{i,j=1}^n (1 - \mathbf{w}^\top(x_i - x_j))^2 \mathbb{I}_{[y_i=1]} \mathbb{I}_{[y_j=-1]}. \quad (3)$$

The objective function $F(\mathbf{w}) = \frac{1}{n_+ n_-} \sum_{i,j=1}^n (1 - \mathbf{w}^\top(x_i - x_j))^2 \mathbb{I}_{[y_i=1]} \mathbb{I}_{[y_j=-1]}$ is the average of pairwise losses and has the form of U-statistics [40].

Objective and Challenges. Our main objective in this paper is to develop efficient stochastic optimization algorithms for the sparse AUC maximization formulation (3) which is scalable to large scale and high-dimensional imbalanced data.

One possible approach is to directly apply stochastic hard thresholding algorithms [34], [35], [44] to the setting of AUC maximization by regarding pairs of examples as individual ones, i.e. at each time we randomly samples a pair of examples or a minibatch of pairs to update model parameters. However, this means that one pass of the data, i.e. passing all pairs, will require n passes of the original dataset \mathcal{S} which makes it not suitable for large-scale and high-dimensional data analysis. To address this challenge, we show in the following subsections that this multiple passes can be avoided by reformulating the minimization problem of pairwise U-statistics objective function $F(\mathbf{w})$ as a novel ERM formulation. From this new reformulation, we can develop efficient stochastic hard thresholding algorithms for AUC maximization.

B. Equivalent Reformulation

Inspired by the work [20], [30], [45], we will formulate the U-statistics objective function in (3) as an ERM objective function, i.e. singled-summed objective function. For this purpose, let the positive and negative sample mean be respectively denoted by $\bar{\mathbf{x}}_+ = \frac{\sum_{i=1}^n \mathbf{x}_i \mathbb{I}_{[y_i=1]}}{n_+}$, $\bar{\mathbf{x}}_- = \frac{\sum_{i=1}^n \mathbf{x}_i \mathbb{I}_{[y_i=-1]}}{n_-}$. Then, we have the following proposition.

Proposition 1. *The empirical AUC objective function $F(\mathbf{w})$ given by (3) can be reformulated as*

$$F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \tilde{f}(\mathbf{w}; \mathbf{z}_i) \quad (4)$$

where $\tilde{f}(\mathbf{w}; \mathbf{z}_i) = \frac{1}{r} (\mathbf{w}^\top (\mathbf{x}_i - \bar{\mathbf{x}}_+))^2 \mathbb{I}_{[y_i=1]} + \frac{1}{1-r} (\mathbf{w}^\top (\mathbf{x}_i - \bar{\mathbf{x}}_-))^2 \mathbb{I}_{[y_i=-1]} + 1 + 2\mathbf{w}^\top (\bar{\mathbf{x}}_- - \bar{\mathbf{x}}_+) + (\mathbf{w}^\top (\bar{\mathbf{x}}_- - \bar{\mathbf{x}}_+))^2$.

The proof of Proposition 1 is inspired by [17], [20]. However, the original proofs there need to introduce three auxiliary variables. Our proof is much simpler and straightforward without introducing auxiliary variables.

Assume that n can be divided by m and let the block size $b = n/m$. Then, in order to apply minibatch updates, let $\{B_i : i \in [m]\}$ denote non-overlapping subsets of \mathcal{S} , each of which is of size b . Therefore, the U-statistics form in the AUC maximization problem (3), with $m = n/b$, can be formulated as the following ERM with sparse constraints:

$$\mathbf{w}_* = \arg \min_{\|\mathbf{w}\|_0 \leq k_*} F(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m f_{B_i}(\mathbf{w}). \quad (5)$$

where $f_{B_i}(\mathbf{w}) = \frac{1}{b} \sum_{j \in B_i} \tilde{f}(\mathbf{w}; \mathbf{z}_j)$ is a block objective.

C. The SHT-AUC Igorithm

From the formulation (5), we are ready to present the Stochastic Hard Thresholding AUC Optimization algorithm which is referred to as SHT-AUC.

The pseudo-code is given by Algorithm 1 which is taken from [37], [46], [47]. It can be regarded as "expansive" projected SGD with projections to the ℓ_0 constraints. Specifically, at each iteration, it randomly selects i_t from $[m]$ with probability $\frac{1}{m}$, and hence the minibatch B_{i_t} . Then, the current model parameter \mathbf{w}_t is updated using projected gradient descent

Algorithm 1 SHT-AUC Algorithm

Input: Relaxed sparsity level k , step size γ , initial classifier \mathbf{w}_0 such that $\|\mathbf{w}_0\|_0 \leq k$
Compute: $\bar{\mathbf{x}}_+$ and $\bar{\mathbf{x}}_-$
for $t = 0$ to $T - 1$ **do**
 Randomly selected $i_t \in [m]$
 $\mathbf{w}_{t+1} = \mathcal{H}_k(\mathbf{w}_t - \gamma \nabla f_{B_{i_t}}(\mathbf{w}_t))$
end for
Output: \mathbf{w}_T

based on the gradient $\nabla f_{B_{i_t}}$, which is the hard thresholding operator given by (1). The main computation is $\mathcal{O}(bd)$ with the gradient and $\mathcal{O}(d)$ with the hard thresholding. Hence the per-iteration cost is $\mathcal{O}(bd)$. To further explain why the Hard Thresholding Operator $\mathcal{H}_k(\mathbf{w}) : \mathbb{R}^d \mapsto \mathbb{R}^d$ in Algorithm 1 only needs time complexity $\mathcal{O}(d)$, firstly we consider a common choice to fulfill \mathcal{H}_k – use sorting algorithm of \mathbf{w} with the time complexity $\mathcal{O}(d \log(d))$ in expectation, such as quick sort, then pick the largest k entries. In fact, given $\mathbf{w} \in \mathbb{R}^d$ and the sparsity $k \ll d$, the Hard Thresholding Operator can be computed as following

$$\mathcal{H}_k(\mathbf{w}) = \mathbf{w} \cdot \mathbb{I}_{|w| \geq |w_{\tau_k}|} \mathbf{1}_d, \quad (6)$$

where τ_k is the index of k -th largest magnitude among $|w_1|, |w_2|, \dots, |w_d|$ and all operations in (6) are element-wise. Clearly, if we know w_{τ_k} in advance, we can get the solution of the operator in $\Theta(d)$. To find w_{τ_k} , we choose the Floyd-Rivest algorithm [48] with time complexity $\mathcal{O}(d)$ in expectation. Algorithm 2 is the pseudo-code of the Floyd-Rivest method. One can immediately find $|w_{\tau_k}| = |w_k|$ after call **Select**($\mathbf{w}, 0, d - 1, k$).

Another note is here the sparsity level k in SHT-AUC is not necessary to be k_* . In particular, the flexible choice of $k \geq k_*$ follows the appealing work [37], [49], which will allow a relaxed projection to the ℓ_0 constraints, and therefore lead to tighter bounds as we can see below.

III. CONVERGENCE ANALYSIS

In this section, we turn to the convergence analysis of SHT-AUC algorithm. The convergence typically need the following standard assumptions.

Assumption 1. *The function $F(\mathbf{w})$ satisfies the ρ_k^- -restricted strong convexity (RSC) condition if there exists a positive constant ρ_k^- such that*

$$F(\mathbf{w}') - F(\mathbf{w}) - \langle \nabla F(\mathbf{w}), \mathbf{w}' - \mathbf{w} \rangle \geq \frac{\rho_k^-}{2} \|\mathbf{w}' - \mathbf{w}\|_2^2 \quad (7)$$

for any \mathbf{w} and \mathbf{w}' such that $|\text{supp}(\mathbf{w}) \cup \text{supp}(\mathbf{w}')| \leq k$.

Assumption 2. *For all $1 \leq i \leq m$, the function $f_{B_i}(\mathbf{w})$ satisfies the ρ_k^+ -restricted strong smoothness (RSS) condition if there exists a positive constant ρ_k^+ such that*

$$\|\nabla f_{B_i}(\mathbf{w}) - \nabla f_{B_i}(\mathbf{w}')\|_2 \leq \rho_k^+ \|\mathbf{w} - \mathbf{w}'\|_2 \quad (8)$$

for all vectors \mathbf{w} and \mathbf{w}' such that $|\text{supp}(\mathbf{w}) \cup \text{supp}(\mathbf{w}')| \leq k$.

Algorithm 2 $\text{Select}(w, l, r, k)$: Floyd-Rivest Algorithm ([48])

```

1: while  $r > l$  do
2:   if  $r - l > 600$  then
3:      $n = r - l + 1$ ;  $i = k - l + 1$ ;
4:      $z = \ln(N)$ ;  $s = 0.5 * \exp(2 * z/3)$ ;
5:      $sd = 0.5 * \sqrt{z * s * (n - s)/n * \text{sign}(i - n/2)}$ ;
6:      $ll = \max(l, k - i * s/n + sd)$ ;
7:      $rr = \min(r, k + (n - 1) * s/n + sd)$ ;
8:     Select( $w, ll, rr, k$ );
9:   end if
10:   $t = w_k$ ;  $i = l$ ;  $j = r$ ;  $\text{swap}(w_l, w_k)$ 
11:  while  $i < j$  do
12:     $\text{swap}(w_i, w_j)$ ;  $i = i + 1$ ;  $j = j - 1$ ;
13:    while  $w_i < t$  do
14:       $i = i + 1$ ;
15:    end while
16:    while  $w_j < t$  do
17:       $j = j - 1$ ;
18:    end while
19:  end while
20:  if  $w_l = t$  then
21:     $\text{swap}(w_l, w_j)$ ;
22:  else
23:     $j = j + 1$ ;  $\text{swap}(w_j, w_r)$ ;
24:  end if
25:  if  $j \leq k$  then
26:     $l = j + 1$ ;
27:  end if
28:  if  $k \leq j$  then
29:     $r = j - 1$ ;
30:  end if
31: end while

```

RSC/RSS properties are firstly introduced in [41]. Since it captures sparsity of many functions, it has been widely used for designing sparsity constrained algorithms [34], [35], [37], [49], [50]. Here, we define the k -restricted condition number to be $\rho_k = \rho_k^+ / \rho_k^-$.

In the sequel, we will first state the convergence results related to the RSC and RSS conditions. Then we will prove that the objective function of AUC maximization defined by $F(w)$ satisfies the RSC and RSS conditions and discuss their implications on the convergence of SHT-AUC .

A. General Convergence Results

To state the convergence of SHT-AUC recall that k_* is the desired sparsity level and k is the relaxed sparsity level. We are now ready to state the general convergence result of SHT-AUC .

Theorem 1. *Let w_* be a k_* -sparse vector of interest, and w_0 be the initial solution. Consider the problem (5) with sparsity level k such that $d \gg k > (\rho_{2k+k_*}^2 - \rho_{2k+k_*})k_*$. Select $\gamma = \frac{1}{\rho_{2k+k_*}^+}$ and let $\nu = 1 + k_*/k + \sqrt{k_*/k}$, we have,*

$$\mathbb{E} \|w_{t+1} - w_*\|_2 \leq \kappa^{t+1} \|w_0 - w_*\|_2 + \frac{\sigma_{w_*}}{1 - \kappa} \quad (9)$$

where the expectation is taken over all choices of random variables i_0, \dots, i_t . Here

$$\kappa = \sqrt{\nu(1 - 1/\rho_{2k+k_*})} < 1 \quad (10)$$

is the convergence parameter and

$$\sigma_{w_*} = \frac{\gamma}{m} \sqrt{\nu} \sum_{i=1}^m \max_{|\Omega| \leq 2k+k_*} \|\mathcal{P}_\Omega(\nabla f_{B_i}(w_*))\|_2 \quad (11)$$

is the tolerance error parameter.

Theorem 1 shows that Algorithm 1 still possibly enjoys linear convergence up to the tolerance error $\sigma_{w_*}/(1 - \kappa)$.

As indicated in Theorem 1, both the convergence rate and the error are depending on the restricted condition number ρ_{2k+k_*} . In the next subsection, we will characterize ρ_{2k+k_*} in terms of imbalance ratio r .

B. Estimation of RCS and RSS Conditions

In this subsection, we will estimate the RSC and RSS conditions, and the condition number for AUC maximization. Combining this with Theorem 1, we will discuss the implications of these estimations, particularly on the effect of imbalance ratio $r = \frac{n_+}{n_-}$ on the convergence of SHT-AUC .

For our analysis, we assume each $x_i \in \mathbb{R}^d$ are i.i.d Gaussian random vectors from $\mathcal{N}(0, \Sigma)$ with covariance matrix Σ and its diagonal elements satisfying $\Sigma_{jj} \leq 1$. We also define a shorthand notation $\lambda = \lambda_{\min}(\Sigma^{1/2})$. Now we have the following theorem.

Theorem 2. *Consider objective function of AUC maximization given by (5) and Algorithm 1 with sparsity level $k < d$. The RCS/RSS condition is satisfied and we have following results: With probability at least $1 - \exp(-n^+/72) - 2(2k + k_*)/d$, there holds*

$$\rho_{2k+k_*}^- = \left(\frac{1}{2} \lambda - 6\sqrt{2} \sqrt{\frac{(2k + k_*) \log d}{rn}} \right)^2 - \frac{32(2k + k_*) \log(d)}{3rn}. \quad (12)$$

With probability at least $1 - (2k + k_*)/2d$, we have

$$\rho_{2k+k_*}^+ = \frac{16(2k + k_*) \log(d) (\frac{1}{2} \log(b) + \log(d))}{r}. \quad (13)$$

Imbalance ratio on the contraction coefficient. Since we focus on the case of large scale problem, we can assume the number of total examples n is large enough (mainly negative examples n_-) such that ρ_k^- is positive. We can write the

restricted condition number ρ_k as a function of the imbalance ratio r ,

$$\rho_k(r) = \frac{16}{ar + b\sqrt{r} + c}, \quad (14)$$

where the coefficients are $a = \frac{\lambda^2}{4k \log(d) (\frac{1}{2} \log(b) + \log(d))}$, $b = \frac{6\sqrt{2}\lambda}{\sqrt{nk \log(d) (\frac{1}{2} \log(b) + \log(d))}}$ and $c = \frac{184}{3n (\frac{1}{2} \log(b) + \log(d))}$. The bottom of (14) is a concave quadratic function of \sqrt{r} with its minimum attaining at the axis of symmetry: $\sqrt{r_*} = \frac{12\sqrt{2}\sqrt{k \log(d) (\frac{1}{2} \log(b) + \log(d))}}{\lambda\sqrt{n}}$. Since we consider the regime when n (or n_-) sufficiently large, i.e. the axis of symmetry is close to 0. Therefore $\rho_k(r)$ can be regarded as a monotonically decreasing function of \sqrt{r} . Recall in Theorem 1 equation (10), $\kappa = \Omega(\sqrt{1 - 1/\rho_{2k+k_*}})$ is monotonically increasing with respect to ρ_{2k+k_*} . Therefore κ is also a monotonically decreasing function of r .

Imbalance ratio on the tolerance error. Recall that the step-size in Theorem 1 is chosen as $\gamma = 1/\rho_{2k+k_*}^+$. Hence, the tolerance parameter in Theorem 1 equation (11) is of the form $\sigma_{\mathbf{w}_*} = \Omega(1/\rho_{2k+k_*}^+)$. Now, combine the discussion on the contraction parameter κ , the total tolerance error, after simplification, is of the form

$$\frac{\sigma_{\mathbf{w}_*}}{1 - \kappa} \geq \frac{c_1 r}{1 - c_2 \sqrt{1 - r}} \quad (15)$$

for some constant c_4 and c_5 . See the exact values in the detailed proofs. Combining this with the above discussion on the relation between r and κ , we can conclude that *the more imbalance the data is, the slower the convergence is, and the larger tolerance error is*, which matches the empirical experience in the subsequent section.

IV. EXPERIMENTS

To validate the effectiveness of our proposed SHT-AUC and test our theory, we apply it to both synthetic and real-world datasets.

Baseline methods.¹ We consider six baseline methods which can be divided into two kinds. The first kind is methods that directly optimize the AUC objective. It includes: SOLAM, a Stochastic OnLine algorithm for AUC Maximization proposed in [17]; SPAM-based, a stochastic proximal algorithm for AUC maximization designed in [20]. Based on different regularizations, we refer SPAM using ℓ_1 and ℓ^2 as SPAM- ℓ_1 , SPAM- ℓ^2 respectively; FSAUC, a Fast Stochastic algorithm for true AUC maximization as proposed in [19]. The second kind is algorithms that optimize the logistic loss with ℓ_0 -norm constraint. We consider two popular methods of this type including STOIHT, a Stochastic Iterative Hard Thresholding method defined [46] and HSG-HT, a Hybrid Stochastic Gradient Hard Thresholding [35] algorithm.

Evaluation Metrics. One of the main goals is to testify the effectiveness of optimizing AUC score and the feature

¹We did not consider methods such as OAM [12] and OPAUC [13] due to their inferior performance on both run time and AUC score reported in [17], [19].

selection ability. We use AUC score [3] for the classification performance and use F1 score for the feature selection. The F1 score with respect to \mathbf{w}_t and \mathbf{w}_* is defined as

$$F1(\mathbf{w}_t, \mathbf{w}_*) = \frac{2 \text{Pre}(\mathbf{w}_t, \mathbf{w}_*) \cdot \text{Rec}(\mathbf{w}_t, \mathbf{w}_*)}{\text{Pre}(\mathbf{w}_t, \mathbf{w}_*) + \text{Rec}(\mathbf{w}_t, \mathbf{w}_*)},$$

where $\text{Pre}(\mathbf{w}_t, \mathbf{w}_*) = \frac{|\text{supp}(\mathbf{w}_*) \cap \text{supp}(\mathbf{w}_t)|}{\|\mathbf{w}_t\|_0}$ and $\text{Rec}(\mathbf{w}_t, \mathbf{w}_*) = \frac{|\text{supp}(\mathbf{w}_*) \cap \text{supp}(\mathbf{w}_t)|}{\|\mathbf{w}_*\|_0}$. We also use the Jaccard Index as an alternative metric to evaluate the feature selection ability. Jaccard Index (JI) with respect to \mathbf{w}_t and \mathbf{w}_* is defined as

$$JI(\mathbf{w}_t, \mathbf{w}_*) = \frac{|\text{supp}(\mathbf{w}_*) \cap \text{supp}(\mathbf{w}_t)|}{|\text{supp}(\mathbf{w}_*) \cup \text{supp}(\mathbf{w}_t)|}.$$

A. Synthetic Datasets

Data generation. We first generate simulation datasets with the data size $n = 1000$ and dimension $d = 1000$. This simulation is motivated from the task of disease outbreak detection [51]. More specifically, for each of the datasets, each training sample $\mathbf{x} \in \mathbb{R}^d$ and $y \in \{\pm 1\}$. All entries of each training sample are from $\mathcal{N}(0, 1)$ while each positive training sample is generated according to: $\mathbf{x}_i \sim \mathcal{N}(\mu, 1)$ if $i \in S$ and $\mathbf{x}_i \sim \mathcal{N}(0, 1)$ if $i \notin S$. Here we fix $\mu = 0.3$ and S is a subset of ‘‘important features’’ as the ground truth features. They are randomly selected from $\{0, 1, 2, \dots, 999\}$ and the size of S is treated as the true sparsity k_* . We generate datasets that have different true sparsity $k_* \in \{20 : 20 : 80\}$ and different imbalance ratios $r \in \{0.05 : 0.05 : 0.5\}$.

Parameter Tuning. For SOLAM, it has two parameters, the bound $R \in 10^{[-1:1:5]}$ on \mathbf{w} and the initial learning rate $\xi \in [1 : 9 : 100]$ as suggested; For the SPAM- ℓ_1 method, it has ℓ_1 -regularization parameter $\beta_1 \in 10^{[-5:1:2]}$; Similarly, SPAM- ℓ^2 has the ℓ^2 -regularization parameter $\beta_2 \in 10^{[-5:1:2]}$ or both. The initial learning rate ξ of SPAM- ℓ_1 and SPAM- ℓ_2 is the same as ξ in SOLAM; For FSAUC, the initial step size η_1 is tuned from $2^{[-10:1:10]}$ and the bound parameter R of \mathbf{w} is the same as SOLAM. For three non-convex methods, the sparsity parameter k of SHT-AUC, STOIHT, and HSG-HT is tuned from $k \in \{10 : 10 : 100\}$. The number of blocks of SHT-AUC and STOIHT is from $\{1, 2, 4, 8, 10\}$.

	$k_* = 20$	$k_* = 40$	$k_* = 60$	$k_* = 80$
SHT-AUC	.551±.107	.675±.068	.766±.074	.820±.061
SPAM- ℓ^1	.560±.087	.621±.094	.697±.118	.763±.128
SPAM- ℓ^2	.537±.095	.597±.110	.653±.141	.752±.135
FSAUC	.571±.107	.654±.083	.754±.079	.820±.071
SOLAM	.523±.102	.628±.077	.732±.092	.740±.139
StoiHT	.538±.096	.604±.087	.659±.091	.719±.092
HSG-HT	.484±.094	.593±.089	.661±.116	.759±.083

TABLE I
AVERAGED AUC ON FOUR SYNTHETIC DATASETS.

Generalization Performance and Feature Selection. Table I reports the averaged AUC of four datasets with imbal-

	F1 score				Jaccard Index			
	$k_* = 20$	$k_* = 40$	$k_* = 60$	$k_* = 80$	$k_* = 20$	$k_* = 40$	$k_* = 60$	$k_* = 80$
SHT-AUC	.209±.046	.365±.078	.382±.053	.450±.072	.126±.047	.200±.043	.275±.051	.311±.032
SPAM- ℓ^1	.058±.053	.182±.137	.147±.102	.177±.126	.028±.040	.087±.076	.078±.060	.101±.059
SPAM- ℓ^2	.037±.019	.060±.053	.100±.085	.159±.065	.017±.021	.029±.020	.040±.034	.065±.033
FSAUC	.100±.075	.202±.114	.222±.096	.320±.102	.037±.033	.117±.071	.146±.060	.210±.069
SOLAM	.044±.021	.088±.039	.125±.064	.171±.064	.024±.029	.049±.032	.073±.027	.100±.072
StoIHT	.089±.037	.163±.069	.231±.054	.237±.067	.051±.033	.093±.028	.122±.036	.146±.045
HSG-HT	.089±.042	.157±.076	.228±.061	.249±.066	.046±.037	.096±.031	.127±.049	.171±.042

TABLE II
AVERAGED F1 SCORES AND JACCARD INDEX ON FOUR SYNTHETIC DATASETS.

anced ratio $r = 0.05^2$. First of all, SHT-AUC gives the best AUC score for $k_* = 40, 60, 80$ and gives competitive AUC score when $k_* = 20$. In fact the AUC scores of SHT-AUC and FSAUC are competitive with each other. One of the reasons could be that both of them have a sparse projection at each iteration. Secondly, the AUC scores of SPAM- ℓ^2 and SOLAM are inferior to SHT-AUC, SPAM- ℓ_1 , and FSAUC. This is because these two are not sparse-inducing methods hence not suitable for sparse learning problem. Last but not least, ℓ_0 -based methods including STOIHT and HSG-HT have lower AUC scores since these two are not for AUC optimization. It shows that our algorithm generalizes well by solving (3) when the ground truth w_* is sparse.

Table II reports the average F1 score and Jaccard Index of four datasets with imbalanced ratio $r = 0.05$. In both metric, our method SHT-AUC are significantly better than any other algorithms. This impact is two-fold. Firstly, SHT-AUC is better than other ℓ_1 based stochastic AUC maximization algorithms. This is consistent with the fact ℓ_1 based stochastic algorithms may be hard to preserve a truly sparse solution [31]–[33]. And it shows the advantage of using ℓ_0 based stochastic algorithm as SHT-AUC. Secondly, SHT-AUC is better than STOIHT and HSG-HT. The advantage of directly optimizing AUC compared with using Empirical Risk Minimization, i.e. logistic loss, when the dataset is imbalanced can also be found in [43]. Our findings prove this well.

In summary, the simulation results indicate that SHT-AUC has better tradeoff between the AUC optimization and feature selection among all methods when the data is imbalanced and the ground truth is sparse.

Effect of Imbalanced Ratio on Convergence and Performance. To demonstrate the impact of imbalance ratio r on the convergence of SHT-AUC, we apply SHT-AUC on datasets with different imbalance ratios $r = 0.05, 0.25, 0.50$. Figure 1 reports the number of epochs against the AUC score, with sparsity level $k_* = 20, 40, 60, 80$ and fix $k = k_*$ and batch size $b = 50$. Note the AUC scores are scaled in order to get better visualization. We can observe that when $r = 0.5$ the SHT-AUC converges after 150 epochs, but when $r = 0.05$, SHT-AUC does not converge even after 300 epochs. This

²AUC scores are calculated on testing dataset. We found that SPAM-based, SOLAM and FSAUC do not produce sparse solutions. Instead, we truncate all entries in w_t to 0 if the magnitude of these entries are not larger than 0.001.

results proves our theoretical analysis in Section III-B, i.e. when data is more imbalanced, the convergence is slower. It also matches ones empirical expectation.

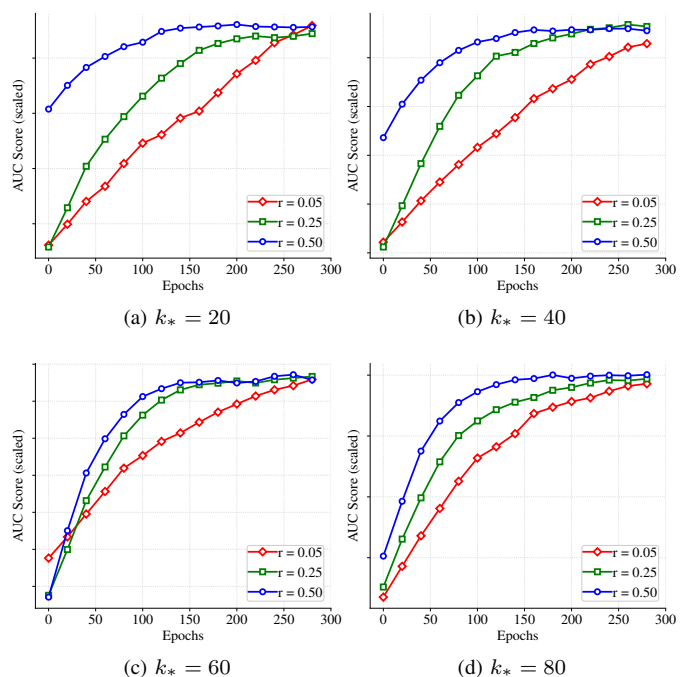


Fig. 1. Convergence plot with different imbalance ratio r .

To further investigate how the imbalance ratio r affects the performance, we apply all methods with different imbalance ratios on $k_* = 20$ dataset. As shown in Figure 2, for SHT-AUC, the more data is imbalanced, the worse the AUC and F1 scores are. This again proves our theoretical analysis in Section III-B. Other AUC maximization algorithms also show the same phenomenon, but they are lack of similar analysis on imbalance ratio. Moreover, the performance of SHT-AUC, SPAM- ℓ_1 , SPAM- ℓ_1/ℓ^2 , and FSAUC are at the same tier. The results of SPAM- ℓ^2 and SOLAM are inferior to the hard thresholding-based and ℓ_1 based methods. The reason is these two methods do not explore sparsity. More interestingly, compared with the methods (i.e. STOIHT and HSG-HT) for Empirical Risk Minimization, the AUC optimization-based methods outperform these two by a large margin when the dataset is more imbalanced. This testifies that minimizing

the empirical risk loss may not lead to the best possible AUC values as stated in [43].

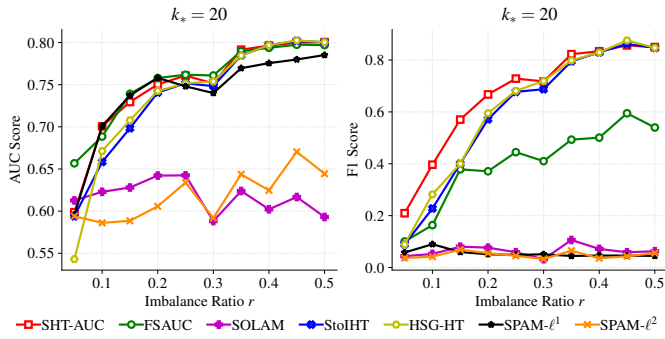


Fig. 2. The left: AUC score as a function of the imbalance ratio r . The right: F1 score as a function of the imbalance ratio r .

B. Gene identification on two real-world datasets

We test our method on two real-world high-dimensional datasets, the leukemia dataset [22] and the colon cancer dataset [52]. The leukemia dataset consists of 72 samples where each positive sample (47 in total) is a patient has acute lymphoblastic leukemia and each negative sample (25 in total) is a patient has acute myeloid leukemia. Each training example has 7,129 genes. The colon cancer dataset has 62 training samples with 40 positive samples (patients who have tumor tissues) and 22 negative (patients who are normal). Each training sample consists 2,000 gene markers. Our goal is to classify these patients at the same time to select genes related with these two disease. As shown in Table IV and V, we choose a subset of ground truth of cancer related genes from [53] and compare the gene selection ability of different methods.

Parameter Tuning. For SHT-AUC, STOIHT, and HSG-HT, the sparsity parameter k is tuned from $\{1, 5, 10, \dots, 50, 60, \dots, 100, 200, \dots, 500\}$. For SPAM- ℓ^1 and SPAM- ℓ^1/ℓ^2 , we choose the ℓ_1 -regularization parameter $\lambda_{\ell_1} \in [0.07, 0.00001]$ such that models are from sparsest models to dense models. FSAUC has a parameter R to control ℓ_1 ball, we choose $R \in [0.00001, 10000]$ such that models are from sparsest models to dense models too. SPAM- ℓ_2 and SOLAM are two non-sparse methods. We randomly shuffle the dataset 20 times which form 20 trials. For each trial, we use 5-fold cross-validation to train all methods. The block size b of SHT-AUC and STOIHT is tuned from 1 to 40 and sparsity k is from 5 to 500.

	Colon Cancer	Leukemia
SHT-AUC	.8777 ± .1114	.9963 ± .0098
SPAM- ℓ^1	.8409 ± .1646	.9812 ± .0602
SPAM- ℓ^2	.8304 ± .1478	.9812 ± .0604
FSAUC	.7907 ± .2143	.9708 ± .0730
SOLAM	.8089 ± .1752	.9751 ± .0773
StoIHT	.8647 ± .1339	.9947 ± .0138
HSG-HT	.8759 ± .1246	.9898 ± .0218

TABLE III
AVERAGE AUC SCORE ON REAL DATASETS



Fig. 3. AUC score as a function of sparsity k .

Generalization Performance. We compare our method SHT-AUC on AUC score with seven baseline methods on colon cancer dataset. As shown in Table III, our algorithm SHT-AUC has highest AUC score among all methods. It again shows the advantage of our algorithm in maximizing AUC with sparsity constraint. Interestingly, Empirical Risk Minimization algorithms STOIHT and HSG-HT also give comparable AUC scores. The explanation is that both datasets are not severely imbalanced. The performance of SPAM, FSAUC and SOLAM are second tier. The results show the advantage of ℓ_0 -based over ℓ_1 -based stochastic algorithms in AUC maximization.

Figure 3 shows the AUC score against the relaxed sparsity level k^3 . In Figure 3a, SHT-AUC, STOIHT and HSG-HT reach their highest AUC scores when k is moderately larger than k_* . This matches the condition in Theorem 1. In Figure 3b, the AUC scores are saturated after certain k with respect to different algorithms. The reason is Leukemia dataset has much higher dimension than Colon Cancer dataset. Hence k is still in a reasonable range, and it still matches the condition in Theorem 1.

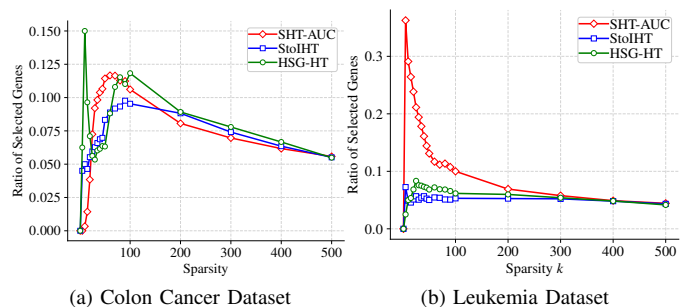


Fig. 4. Percentage of related genes as a function as sparsity k on the colon cancer and leukemia dataset. The ratio of selected genes is defined by the number of selected cancer-related genes divided by total number of genes found (corresponding to total number of non-zeros).

Feature Selection Ability. To further investigate how is the gene selection ability for different methods on these datasets, we measure the gene selection ability as the ratio of related genes selected: the number of genes overlapped with the genes

³SOLAM, FSAUC and SPAM are drawn as constant lines with best performance

ID	Marker	ID	Marker
1	Myeloperoxidase	17	Probable protein disulfide isomerase ER-60 precursor
2	CD13	18	CD34
3	CD33	19	CD24
4	HOXA9 Homeo box A9	20	60S ribosomal protein L23
5	MYBL2	21	5-aminolevulinic acid synthase
6	CD19	22	HLA class II histocompatibility antigen
7	CD10 (CALLA)	23	Epstein-Barr virus small RNA-associated protein
8	TCL1 (T cell leukemia)	24	HNRPA1 Heterogeneous nuclear ribonucleoprotein A1
9	C-myb	25	Azurocidin
10	Deoxyhypusine synthase	26	Red cell anion exchanger (EPB3, AE1, Band 3)
11	KIAA0220	27	Topoisomerase II beta
12	G-gamma globin	28	Probable G protein-coupled receptor LCR1 homolog
13	Delta-globin	29	Int-6
14	Brain-expressed HHCPA78 homolog	30	Alpha-tubulin
15	Myeloperoxidase	31	Terminal transferase
16	NPM1 Nucleophosmin	32	Glycophorin B precursor

TABLE IV
MARKERS RELATED WITH ACUTE MYELOID LEUKEMIA AND ACUTE LYMPHOBLASTIC LEUKEMIA.

ID	Marker	ID	Marker
1	Phospholipase A2	16	Splicing factor (CC1.4)
2	Keratin 6 isoform	17	Nucleolar protein (B23)
3	Protein-tyrosine phosphatase PTP-H1	18	Lactate dehydrogenase-A (LDH-A)
4	Transcription factor IIIA	19	Guanine nucleotide-binding protein G(OLF)
5	Viral (v-raf) oncogene homolog 1	20	LI-cadherin
6	Dual specificity mitogen-activated protein kinase kinase 1	21	Lysozyme
7	Transmembrane carcinoembryonic antigen	22	Prolyl 4-hydroxylase (P4HB)
8	Oncoprotein 18	23	Eukaryotic initiation factor 4AII
9	Phosphoenolpyruvate carboxykinase	24	Interferon-inducible protein 1-8D
10	Extracellular signal-regulated kinase 1	25	Dipeptidase
11	26 kDa cell surface protein TAPA-1	26	Heat shock 27 kDa protein
12	Id1	27	Tyrosine-protein kinase receptor TIE-1 precursor
13	Interferon-inducible protein 9-27	28	Mitochondrial matrix protein P1 precursor
14	Nonspecific crossreacting antigen	29	Eukaryotic initiation factor EIF-4A homolog
15	cAMP response element regulatory protein (CREB2)		

TABLE V
MARKERS RELATED WITH COLON CANCER AS SHOWN IN [53]

defined in Table IV and V divided by total number of genes found. i.e. let w_t be the algorithm output and w_* be the groundtruth, the ratio is defined as

$$\text{Ratio}(w_t, w_*) = \frac{|\text{supp}(w_*) \cap \text{supp}(w_t)|}{|\text{supp}(w_*)|}.$$

We report our results in Figure 4. In Figure 4a, HSG-HT returns the best percentage when the sparsity is $k = 5$, but when $k = 29$, which is the number of related genes, SHT-AUC can achieve the highest percentage of related genes. Hence SHT-AUC and HSG-HT are comparable in general. This might due to the fact the colon cancer dataset is a relatively balanced dataset, with $r = 0.355$. In Figure 4b, SHT-AUC recovers the best percentage also when sparsity is the number of related genes, i.e. $k = 32$, and significantly outperforms HSG-HT and STOIHT. This phenomenon highlights the advantage of maximizing AUC rather than accuracy under imbalanced classification setting.

In summary, SHT-AUC enjoys better generalization performance on real-world high-dimensional datasets, while maintain a more robust feature selection ability against state-of-art algorithms.

V. CONCLUSION

In this paper, we proposed stochastic hard thresholding algorithm for AUC maximization with sparse ℓ_0 constraints in imbalanced classification. In particular, we formulated the U-statistic objective function of AUC maximization as an ERM objective function. This new reformulation facilitated the design of stochastic hard thresholding algorithm for AUC maximization. The proposed algorithm, SHT-AUC, enjoys a cheap $\mathcal{O}(bd)$ per-iteration cost, making it amenable for high-dimensional data analysis. We proved that under RCS/RSS conditions, SHT-AUC enjoys a linear convergence rate up to a tolerance error. We also showed, under Gaussian assumptions on the data, the RCS/RSS conditions can be satisfied and how the convergence rate and tolerance error are affected by the imbalance ratio. Our experiments validated our theoretical findings while SHT-AUC is shown to have a very good property in feature selection against state-of-the-art algorithms.

VI. ACKNOWLEDGEMENT

This work is supported by NSF IIS-1816227 and IIS-2008532. The work of Yunwen Lei is supported by the National Natural Science Foundation of China (Grant Nos. 61806091).

REFERENCES

- [1] A. Bradley, "The use of the area under the ROC curve in the evaluation of machine learning algorithms," *Pattern recognition*, vol. 30, no. 7, pp. 1145–1159, 1997.
- [2] T. Fawcett, "An introduction to ROC analysis," *Pattern recognition letters*, vol. 27, no. 8, pp. 861–874, 2006.
- [3] J. A. Hanley and B. J. McNeil, "The meaning and use of the area under a receiver operating characteristic (roc) curve," *Radiology*, vol. 143, no. 1, pp. 29–36, 1982.
- [4] J. Huang, J. Lu, and C. X. Ling, "Comparing naive bayes, decision trees, and svm with auc and accuracy," in *Third IEEE International Conference on Data Mining*, 2003, pp. 553–556.
- [5] T. Joachims, "A support vector method for multivariate performance measures," in *International Conference on Machine Learning*. ACM, 2005, pp. 377–384.
- [6] A. Herschtal and B. Raskutti, "Optimising area under the ROC curve using gradient descent," in *International Conference on Machine Learning*. ACM, 2004, p. 49.
- [7] X. Zhang, A. Saha, and S. V. N. Vishwanathan, "Smoothing multivariate performance measures," *Journal of Machine Learning Research*, vol. 13, pp. 3623–3680, 2012.
- [8] Y. Nesterov, "A method for solving the convex programming problem with convergence rate $o(1/k^2)$," in *Dokl. Akad. Nauk SSSR*, 1983, pp. 543–547.
- [9] M. Culver, D. Kun, and S. Scott, "Active learning to maximize area under the roc curve," in *Sixth International Conference on Data Mining (ICDM'06)*, 2006, pp. 149–158.
- [10] Y. Wang, R. Kharon, D. Pechyony, and R. Jones, "Generalization bounds for online learning algorithms with pairwise loss functions," in *Proceedings of the 25th Annual Conference on Learning Theory*, ser. Proceedings of Machine Learning Research, S. Mannor, N. Srebro, and R. C. Williamson, Eds., vol. 23. Edinburgh, Scotland: PMLR, 25–27 Jun 2012, pp. 13.1–13.22. [Online]. Available: <http://proceedings.mlr.press/v23/wang12.html>
- [11] P. Kar, B. K. Sriperumbudur, P. Jain, and H. C. Karnick, "On the generalization ability of online learning algorithms for pairwise loss functions," in *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28*, ser. ICML '13. JMLR.org, 2013, pp. III–441–III–449.
- [12] P. Zhao, S. C. H. Hoi, R. Jin, and T. Yang, "Online auc maximization," in *Proceedings of the 28th International Conference on International Conference on Machine Learning*, ser. ICML '11. Madison, WI, USA: Omnipress, 2011, pp. 233–240.
- [13] W. Gao, R. Jin, S. Zhu, and Z.-H. Zhou, "One-pass auc optimization," in *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28*, ser. ICML'13. JMLR.org, 2013, p. III–906–III–914.
- [14] S. Gultekin, A. Saha, A. Ratnaparkhi, and J. Paisley, "Mba: mini-batch auc optimization," *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [15] M. Khalid, I. Ray, and H. Chitsaz, "Scalable nonlinear auc maximization methods," in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 2018, pp. 292–307.
- [16] M. Liu, Z. Yuan, Y. Ying, and T. Yang, "Stochastic auc maximization with deep neural networks," *International Conference on Learning Representations (ICLR)*, 2020.
- [17] Y. Ying, L. Wen, and S. Lyu, "Stochastic online auc maximization," in *Advances in neural information processing systems*, 2016, pp. 451–459.
- [18] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro, "Robust stochastic approximation approach to stochastic programming," *SIAM Journal on optimization*, vol. 19, no. 4, pp. 1574–1609, 2009.
- [19] M. Liu, X. Zhang, Z. Chen, X. Wang, and T. Yang, "Fast stochastic auc maximization with $o(1/n)$ -convergence rate," in *International Conference on Machine Learning*, 2018, pp. 3195–3203.
- [20] M. Natole, Jr., Y. Ying, and S. Lyu, "Stochastic proximal algorithms for AUC maximization," in *Proceedings of the 35th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, J. Dy and A. Krause, Eds., vol. 80. Stockholm Sweden: PMLR, 10–15 Jul 2018, pp. 3710–3719. [Online]. Available: <http://proceedings.mlr.press/v80/natole18a.html>
- [21] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik, "Gene selection for cancer classification using support vector machines," *Machine learning*, vol. 46, no. 1-3, pp. 389–422, 2002.
- [22] T. R. Golub, D. K. Slonim, P. Tamayo, C. Huard, M. Gaasenbeek, J. P. Mesirov, H. Coller, M. L. Loh, J. R. Downing, M. A. Caligiuri, C. D. Bloomfield, and E. S. Lander, "Molecular classification of cancer: class discovery and class prediction by gene expression monitoring," *Science*, vol. 286, no. 5439, pp. 531–537, Oct. 1999.
- [23] T. Hromádka, M. R. DeWeese, and A. M. Zador, "Sparse representation of sounds in the unanesthetized auditory cortex," *PLoS biology*, vol. 6, no. 1, 2008.
- [24] E. J. Candès and M. B. Wakin, "An introduction to compressive sampling," *IEEE signal processing magazine*, vol. 25, no. 2, pp. 21–30, 2008.
- [25] J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, and Y. Ma, "Robust face recognition via sparse representation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 31, no. 2, pp. 210–227, 2008.
- [26] M. Yuan and Y. Lin, "Model selection and estimation in regression with grouped variables," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 68, no. 1, pp. 49–67, 2006.
- [27] B. A. Turlach, W. N. Venables, and S. J. Wright, "Simultaneous variable selection," *Technometrics*, vol. 47, no. 3, pp. 349–363, 2005.
- [28] L. Jacob, G. Obozinski, and J.-P. Vert, "Group lasso with overlap and graph lasso," in *Proceedings of the 26th annual international conference on machine learning*, 2009, pp. 433–440.
- [29] R. Jenatton, J.-Y. Audibert, and F. Bach, "Structured variable selection with sparsity-inducing norms," *Journal of Machine Learning Research*, vol. 12, no. Oct, pp. 2777–2824, 2011.
- [30] Y. Lei and Y. Ying, "Stochastic proximal auc maximization," *arXiv preprint arXiv:1906.06053*, 2019.
- [31] J. Duchi and Y. Singer, "Efficient online and batch learning using forward backward splitting," *Journal of Machine Learning Research*, vol. 10, no. Dec, pp. 2899–2934, 2009.
- [32] J. Langford, L. Li, and T. Zhang, "Sparse online learning via truncated gradient," *Journal of Machine Learning Research*, vol. 10, no. Mar, pp. 777–801, 2009.
- [33] L. Xiao, "Dual averaging methods for regularized stochastic learning and online optimization," *Journal of Machine Learning Research*, vol. 11, no. Oct, pp. 2543–2596, 2010.
- [34] N. Nguyen, D. Needell, and T. Woolf, "Linear convergence of stochastic iterative greedy algorithms with sparse constraints," *IEEE Transactions on Information Theory*, vol. 63, no. 11, pp. 6869–6895, 2017.
- [35] P. Zhou, X. Yuan, and J. Feng, "Efficient stochastic gradient hard thresholding," in *Advances in Neural Information Processing Systems*, 2018, pp. 1985–1994.
- [36] T. Murata and T. Suzuki, "Sample efficient stochastic gradient iterative hard thresholding method for stochastic sparse linear regression with limited attribute observation," in *Advances in Neural Information Processing Systems*, 2018, pp. 5313–5322.
- [37] J. Shen and P. Li, "A tight bound of hard thresholding," *Journal of Machine Learning Research*, vol. 18, no. 208, pp. 1–42, 2018. [Online]. Available: <http://jmlr.org/papers/v18/shen16-299.html>
- [38] B. Liu, X.-T. Yuan, L. Wang, Q. Liu, and D. N. Metaxas, "Dual iterative hard thresholding: From non-convex sparse minimization to non-smooth concave maximization," in *International Conference on Machine Learning*, 2017, pp. 2179–2187.
- [39] J. A. Tropp and S. J. Wright, "Computational methods for sparse solution of linear inverse problems," *Proceedings of the IEEE*, vol. 98, no. 6, pp. 948–958, 2010.
- [40] S. Cléménçon, G. Lugosi, N. Vayatis *et al.*, "Ranking and empirical minimization of u-statistics," *The Annals of Statistics*, vol. 36, no. 2, pp. 844–874, 2008.
- [41] S. Negahban, P. Ravikumar, M. J. Wainwright, and B. Yu, "A unified framework for high-dimensional analysis of m-estimators with decomposable regularizers," in *Proceedings of the 22nd International Conference on Neural Information Processing Systems*, ser. NIPS'09. Red Hook, NY, USA: Curran Associates Inc., 2009, p. 1348–1356.
- [42] A. Agarwal, S. Negahban, M. J. Wainwright *et al.*, "Fast global convergence of gradient methods for high-dimensional statistical recovery," *The Annals of Statistics*, vol. 40, no. 5, pp. 2452–2482, 2012.
- [43] C. Cortes and M. Mohri, "Auc optimization vs. error rate minimization," in *Advances in neural information processing systems*, 2004, pp. 313–320.
- [44] J. Shen, P. Li, and H. Xu, "Online low-rank subspace clustering by basis dictionary pursuit," in *International Conference on Machine Learning*, 2016, pp. 622–631.

- [45] Y. Ying and D.-X. Zhou, "Online regularized classification algorithms," *IEEE Transactions on Information Theory*, vol. 52, no. 11, pp. 4775–4788, 2006.
- [46] N. Nguyen, D. Needell, and T. Woolf, "Linear convergence of stochastic iterative greedy algorithms with sparse constraints," *IEEE Transactions on Information Theory*, vol. 63, no. 11, pp. 6869–6895, Nov 2017.
- [47] D. Needell and J. A. Tropp, "Cosamp: Iterative signal recovery from incomplete and inaccurate samples," *Applied and computational harmonic analysis*, vol. 26, no. 3, pp. 301–321, 2009.
- [48] M. Blum, R. W. Floyd, V. Pratt, R. L. Rivest, and R. E. Tarjan, "Time bounds for selection," *J. Comput. Syst. Sci.*, vol. 7, no. 4, pp. 448–461, Aug. 1973. [Online]. Available: [https://doi.org/10.1016/S0022-0000\(73\)80033-9](https://doi.org/10.1016/S0022-0000(73)80033-9)
- [49] P. Jain, A. Tewari, and P. Kar, "On iterative hard thresholding methods for high-dimensional m-estimation," in *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 1*, ser. NIPS '14. Cambridge, MA, USA: MIT Press, 2014, pp. 685–693.
- [50] E. R. Elenberg, R. Khanna, A. G. Dimakis, S. Negahban *et al.*, "Restricted strong convexity implies weak submodularity," *The Annals of Statistics*, vol. 46, no. 6B, pp. 3539–3568, 2018.
- [51] E. Arias-Castro, E. J. Candès, and A. Durand, "Detection of an anomalous cluster in a network," *The Annals of Statistics*, pp. 278–304, 2011.
- [52] U. Alon, N. Barkai, D. Notterman, K. Gish, S. Ybarra, D. Mack, and A. Levine, "Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays," *Proceedings of the National Academy of Sciences*, vol. 96, no. 12, pp. 6745–6750, Jun. 1999.
- [53] S. Agarwal and S. Sengupta, "Ranking genes by relevance to a disease," in *Proceedings of the 8th annual international conference on computational systems bioinformatics*, 2009.
- [54] G. Raskutti, M. J. Wainwright, and B. Yu, "Minimax rates of estimation for high-dimensional linear regression over ℓ_q -balls," *arXiv e-prints*, p. arXiv:0910.2042, Oct 2009.
- [55] P. C. Bellec, G. Lecué, and A. B. Tsybakov, "Slope meets Lasso: improved oracle bounds and optimality," *arXiv e-prints*, p. arXiv:1605.0865, May 2016.

VII. SUPPLEMENTARY MATERIAL

In this Supplementary Material, we provide the detailed proofs for Proposition 1, Theorems 1 and 2.

A. Proof of Proposition 1

Proof. The objective function of AUC maximization given by (4) can be write in three terms,

$$\begin{aligned}
F(\mathbf{w}) &= \frac{1}{n_+} \frac{1}{n_-} \sum_{i=1}^n \sum_{j=1}^n (1 - \mathbf{w}^\top (\mathbf{x}_i - \mathbf{x}_j))^2 \mathbb{I}_{[y_i=1]} \mathbb{I}_{[y_j=-1]} \\
&= \frac{1}{n_+} \frac{1}{n_-} \sum_{i=1}^n \sum_{j=1}^n (1 + \mathbf{w}^\top (\bar{\mathbf{x}}_- - \bar{\mathbf{x}}_+) - \mathbf{w}^\top (\mathbf{x}_i - \bar{\mathbf{x}}_+) + \mathbf{w}^\top (\mathbf{x}_j - \bar{\mathbf{x}}_-))^2 \mathbb{I}_{[y_i=1]} \mathbb{I}_{[y_j=-1]} \\
&= \frac{1}{n_+} \frac{1}{n_-} \underbrace{\sum_{i=1}^n \sum_{j=1}^n (1 + \mathbf{w}^\top (\bar{\mathbf{x}}_- - \bar{\mathbf{x}}_+))^2 \mathbb{I}_{[y_i=1]} \mathbb{I}_{[y_j=-1]}}_I \\
&\quad + \frac{1}{n_+} \frac{1}{n_-} \underbrace{\sum_{i=1}^n \sum_{j=1}^n (\mathbf{w}^\top (\mathbf{x}_i - \bar{\mathbf{x}}_+) - \mathbf{w}^\top (\mathbf{x}_j - \bar{\mathbf{x}}_-))^2 \mathbb{I}_{[y_i=1]} \mathbb{I}_{[y_j=-1]}}_{II} \\
&\quad + \frac{1}{n_+} \frac{1}{n_-} \underbrace{\sum_{i=1}^n \sum_{j=1}^n 2 (1 + \mathbf{w}^\top (\bar{\mathbf{x}}_- - \bar{\mathbf{x}}_+)) (\mathbf{w}^\top (\mathbf{x}_i - \bar{\mathbf{x}}_+) - \mathbf{w}^\top (\mathbf{x}_j - \bar{\mathbf{x}}_-)) \mathbb{I}_{[y_i=1]} \mathbb{I}_{[y_j=-1]}}_{III}.
\end{aligned}$$

It suffices to estimate the above terms one by one. To this end, the first term has $n_+ n_-$ same terms, so

$$I = (1 + \mathbf{w}^\top (\bar{\mathbf{x}}_- - \bar{\mathbf{x}}_+))^2 = 1 + 2\mathbf{w}^\top (\bar{\mathbf{x}}_- - \bar{\mathbf{x}}_+) + (\mathbf{w}^\top (\bar{\mathbf{x}}_- - \bar{\mathbf{x}}_+))^2.$$

For the second term, notice that the cross term

$$\begin{aligned}
&\frac{1}{n_+} \frac{1}{n_-} \sum_{i=1}^n \sum_{j=1}^n 2\mathbf{w}^\top (\mathbf{x}_i - \bar{\mathbf{x}}_+) \mathbf{w}^\top (\mathbf{x}_j - \bar{\mathbf{x}}_-) \mathbb{I}_{[y_i=1]} \mathbb{I}_{[y_j=-1]} \\
&= 2 \left(\frac{1}{n_+} \sum_{i=1}^n \mathbf{w}^\top (\mathbf{x}_i - \bar{\mathbf{x}}_+) \mathbb{I}_{[y_i=1]} \right) \left(\frac{1}{n_-} \sum_{j=1}^n \mathbf{w}^\top (\mathbf{x}_j - \bar{\mathbf{x}}_-) \mathbb{I}_{[y_j=-1]} \right) \\
&= 2 (\mathbf{w}^\top (\bar{\mathbf{x}}_+ - \bar{\mathbf{x}}_+)) (\mathbf{w}^\top (\bar{\mathbf{x}}_- - \bar{\mathbf{x}}_-)) \\
&= 0,
\end{aligned}$$

we have

$$\begin{aligned}
II &= \frac{1}{n_+} \frac{1}{n_-} \sum_{i=1}^n \sum_{j=1}^n (\mathbf{w}^\top (\mathbf{x}_i - \bar{\mathbf{x}}_+))^2 \mathbb{I}_{[y_i=1]} \mathbb{I}_{[y_j=-1]} + \frac{1}{n_+} \frac{1}{n_-} \sum_{i=1}^n \sum_{j=1}^n (\mathbf{w}^\top (\mathbf{x}_j - \bar{\mathbf{x}}_-))^2 \mathbb{I}_{[y_i=1]} \mathbb{I}_{[y_j=-1]} \\
&= \frac{1}{n_+} \sum_{i=1}^n (\mathbf{w}^\top (\mathbf{x}_i - \bar{\mathbf{x}}_+))^2 \mathbb{I}_{[y_i=1]} + \frac{1}{n_-} \sum_{j=1}^n (\mathbf{w}^\top (\mathbf{x}_j - \bar{\mathbf{x}}_-))^2 \mathbb{I}_{[y_j=-1]} \\
&= \frac{1}{n_+} \sum_{i=1}^n (\mathbf{w}^\top (\mathbf{x}_i - \bar{\mathbf{x}}_+))^2 \mathbb{I}_{[y_i=1]} + \frac{1}{n_-} \sum_{i=1}^n (\mathbf{w}^\top (\mathbf{x}_i - \bar{\mathbf{x}}_-))^2 \mathbb{I}_{[y_i=-1]}.
\end{aligned}$$

For the third term, by a similar argument of the cross term in the second term, we have $III = 0$. Now the equation (4) holds by proper scaling. \square

B. Proof of Theorem 1

Before we introduce the proof of the theorem we need several lemmas.

The first lemma was originally proved in [34]. It provides an estimate for our convergence analysis. Recall that we assume $\{f_{B_i}(\mathbf{w})\}_{i=1}^m$ satisfies the RSS and $F(\mathbf{w}) = \sum_{i=1}^m f_{B_i}(\mathbf{w})$ satisfies the RSC.

Lemma 1. Let i be an index selected with probability $1/n$ from the set $[n]$. For any fixed sparse vectors \mathbf{w} and \mathbf{w}' , let Ω be a set such that $\text{supp}(\mathbf{w}) \cup \text{supp}(\mathbf{w}') \in \Omega$ and denote $s = |\Omega|$. We have

$$\mathbb{E}_i \|\mathbf{w}' - \mathbf{w} - \gamma \mathcal{P}_\Omega (\nabla f_{B_i}(\mathbf{w}') - \nabla f_{B_i}(\mathbf{w}))\|_2 \leq \sqrt{1 - (2 - \gamma \rho_s^+) \gamma \rho_s^-} \|\mathbf{w}' - \mathbf{w}\|_2 \quad (16)$$

The second lemma provides a refined bound on the deviation of the thresholded variable, which is originally proved in [37].

Lemma 2. Let $\mathbf{w} \in \mathbb{R}^d$ be an arbitrary vector and $\mathbf{w}^* \in \mathbb{R}^d$ be any k^* -sparse vector. For any $k \geq k^*$, we have the following bound:

$$\|\mathcal{H}_k(\mathbf{w}) - \mathbf{w}^*\|_2 \leq \sqrt{1 + \nu} \|\mathbf{w} - \mathbf{w}^*\|_2, \quad \nu = \frac{\mu + \sqrt{(4 + \mu)\mu}}{2}, \quad \mu = \frac{\min\{k^*, d - k\}}{k - k^* + \min\{k^*, d - k\}}.$$

Proof of Theorem 1. By specifying $\Omega = \text{supp}(\mathbf{w}_{t+1}) \cup \text{supp}(\mathbf{w}_t) \cup \text{supp}(\mathbf{w}_*)$ and notice $|\Omega| \leq 2k + k_*$, it follows that

$$\mathcal{H}_{2k+k_*}(\widehat{\mathbf{w}}_t) = \mathcal{H}_{2k+k_*}(\mathcal{P}_\Omega(\widehat{\mathbf{w}}_t)).$$

Thus, by the updating rule in Algorithm 1 and Lemma 2 we have,

$$\begin{aligned} \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2 &= \|\mathcal{H}_{2k+k_*}(\mathcal{P}_\Omega(\widehat{\mathbf{w}}_t)) - \mathbf{w}_*\|_2 \\ &\leq \sqrt{1 + \nu} \|\mathcal{P}_\Omega(\widehat{\mathbf{w}}_t) - \mathbf{w}_*\|_2 \\ &= \sqrt{1 + \nu} \|\mathbf{w}_t - \mathbf{w}_* - \gamma \mathcal{P}_\Omega(\nabla f_{B_{i_t}}(\mathbf{w}_t))\|_2 \\ &\leq \sqrt{1 + \nu} (\|\gamma \mathcal{P}_\Omega(\nabla f_{B_{i_t}}(\mathbf{w}_*))\|_2 + \|\mathbf{w}_t - \mathbf{w}_* - \gamma \mathcal{P}_\Omega(\nabla f_{B_{i_t}}(\mathbf{w}_t) - \nabla f_{B_{i_t}}(\mathbf{w}_*))\|_2) \end{aligned}$$

where the second inequality holds because $\text{supp}(\mathbf{w}_t - \mathbf{w}_*) \subseteq \Omega$, the second inequality holds because of triangle inequality.

Denote I_t as the set containing all indices i_1, i_2, \dots, i_t randomly selected at or before step t of the algorithm: $I_t = \{i_1, \dots, i_t\}$. It is clear that I_t determines the solutions $\mathbf{w}_1, \dots, \mathbf{w}_{t+1}$. We also denote the conditional expectation $\mathbb{E}_{i_t|I_{t-1}} \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2 \triangleq \mathbb{E}_{i_t}(\|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2 | I_{t-1})$. Now taking the conditional expectation on both sides of the above inequality we obtain

$$\begin{aligned} \mathbb{E}_{i_t|I_{t-1}} \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2 &\leq \sqrt{1 + \nu} (\mathbb{E}_{i_t|I_{t-1}} \|\mathbf{w}_t - \mathbf{w}_* - \gamma \mathcal{P}_\Omega(\nabla f_{B_{i_t}}(\mathbf{w}_t) - \nabla f_{B_{i_t}}(\mathbf{w}_*))\|_2 \\ &\quad + \mathbb{E}_{i_t|I_{t-1}} \|\gamma \mathcal{P}_\Omega(\nabla f_{B_{i_t}}(\mathbf{w}_*))\|_2). \end{aligned}$$

Conditioning on I_{t-1} , \mathbf{w}_t can be seen as a fixed vector. We apply the inequality (16) of Lemma 1, we get

$$\begin{aligned} \mathbb{E}_{i_t|I_{t-1}} \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2 &\leq \sqrt{(1 + \nu) (1 - (2\gamma - \gamma^2 \rho_{2k+k_*}^+) \rho_{2k+k_*}^-)} \|\mathbf{w}_t - \mathbf{w}_*\|_2 \\ &\quad + \sqrt{1 + \nu} \gamma \mathbb{E}_{i_t} \|\mathcal{P}_\Omega(\nabla f_{B_{i_t}}(\mathbf{w}_*))\|_2 \\ &\leq \kappa \|\mathbf{w}_t - \mathbf{w}_*\|_2 + \sigma_{\mathbf{w}_*}, \end{aligned}$$

where κ and $\sigma_{\mathbf{w}_*}$ are defined in Theorem 1. Taking the expectation on both sides with respect to I_{t-1} yields

$$\mathbb{E}_{I_t} \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2 \leq \kappa \mathbb{E}_{I_{t-1}} \|\mathbf{w}_t - \mathbf{w}_*\|_2 + \sigma_{\mathbf{w}_*}.$$

Applying this result recursively over t iterations yields the desired result:

$$\begin{aligned} \mathbb{E}_{I_t} \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2 &\leq \kappa^{t+1} \|\mathbf{w}_0 - \mathbf{w}_*\|_2 + \sum_{j=0}^t \kappa^j \sigma_{\mathbf{w}_*} \\ &\leq \kappa^{t+1} \|\mathbf{w}_0 - \mathbf{w}_*\|_2 + \frac{1}{1 - \kappa} \sigma_{\mathbf{w}_*}. \end{aligned}$$

□

C. Proof of Theorems 2

In order to prove Theorems 2, we need to following lemmas. Firstly, we introduce a lemma which is originally proved in [54]. This lemma captures the lower bound and upper bound of Gaussian random design matrix.

Lemma 3. Consider a random design matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ formed by drawing each row $\mathbf{x}_i \in \mathbb{R}^d$ i.i.d. from an $N(0, \Sigma)$ distribution. Then for some positive constants c_1, c_2, c_3 and c_4 , we have for all $\mathbf{v} \in \mathbb{R}^d$ and $\mathbf{v} \in B_0(2k)$,

$$\frac{\|\mathbf{X}\mathbf{v}\|_2}{\sqrt{n}} \geq \left(\frac{\|\Sigma^{1/2}\mathbf{v}\|_2}{2\|\mathbf{v}\|_2} - 6\sqrt{2} \sqrt{\frac{\rho(\Sigma)k \log d}{n}} \right) \|\mathbf{v}\|_2 \quad (17)$$

with probability $1 - \exp(-n/72)$.

Secondly we include a elementary bound on the sum of ordered Gaussian variables, which is originally proved in [55].

Lemma 4. *Let g_1, \dots, g_d be zero-mean Gaussian random variables with variance at most σ^2 . Denote by $(g_{(1)}, \dots, g_{(d)})$ be a non-increasing rearrangement of $(|g_1|, \dots, |g_d|)$. Then*

$$\mathbb{P} \left(\frac{1}{k\sigma^2} \sum_{j=1}^k g_{(j)}^2 > t \log \left(\frac{2d}{k} \right) \right) \leq \left(\frac{2d}{k} \right)^{1-\frac{3t}{8}} \quad (18)$$

for all $t > 0$ and $k \in \{1, \dots, d\}$.

Proof of Theorem 2. Since F it is a quadratic function of \mathbf{w} , we have

$$F(\mathbf{w}') - F(\mathbf{w}) - \langle \nabla F(\mathbf{w}), \mathbf{w}' - \mathbf{w} \rangle = \frac{1}{2} (\mathbf{w}' - \mathbf{w})^\top \nabla^2 F(\mathbf{w}) (\mathbf{w}' - \mathbf{w}).$$

By the definition of F in equation (4), we have

$$\begin{aligned} \nabla^2 F &= \frac{1}{n_+} \frac{1}{n_-} \sum_{i=1}^n \sum_{j=1}^n (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top \mathbb{I}_{[y_i=1]} \mathbb{I}_{[y_j=-1]} \\ &= \frac{1}{n^+} \frac{1}{n^-} \sum_{i=1}^{n^+} \sum_{j=1}^{n^-} (\mathbf{x}_i^+ - \mathbf{x}_j^-)(\mathbf{x}_i^+ - \mathbf{x}_j^-)^\top \\ &= \frac{1}{n^+} \sum_{i=1}^{n^+} \mathbf{x}_i^+ (\mathbf{x}_i^+)^\top + \frac{1}{n^-} \sum_{j=1}^{n^-} \mathbf{x}_j^- (\mathbf{x}_j^-)^\top - \frac{1}{n^+} \frac{1}{n^-} \sum_{i=1}^{n^+} \sum_{j=1}^{n^-} \mathbf{x}_i^+ (\mathbf{x}_j^-)^\top \\ &\quad - \frac{1}{n^+} \frac{1}{n^-} \sum_{i=1}^{n^+} \sum_{j=1}^{n^-} \mathbf{x}_j^- (\mathbf{x}_i^+)^\top \\ &= \frac{1}{n^+} (\mathbf{X}^+)^\top \mathbf{X}^+ + \frac{1}{n^-} (\mathbf{X}^-)^\top \mathbf{X}^- - \bar{\mathbf{x}}^+ (\bar{\mathbf{x}}^-)^\top - \bar{\mathbf{x}}^- (\bar{\mathbf{x}}^+)^\top \end{aligned}$$

By Lemma 3 we have for all $\mathbf{v} \in \mathbb{R}^d$ and $\mathbf{v} \in B_0(2k)$, \mathbf{X}^+ satisfies

$$\frac{\mathbf{v}^\top (\mathbf{X}^+)^\top \mathbf{X}^+ \mathbf{v}}{n^+} \geq \left(\frac{\|\Sigma^{1/2} \mathbf{v}\|_2}{2\|\mathbf{v}\|_2} - 6\sqrt{2} \sqrt{\frac{\rho(\Sigma)k \log d}{n^+}} \right)^2 \|\mathbf{v}\|_2^2, \quad (19)$$

with probability $1 - \exp(-n^+/72)$. And \mathbf{X}^- satisfies

$$\frac{\mathbf{v}^\top (\mathbf{X}^-)^\top \mathbf{X}^- \mathbf{v}}{n^-} \geq \left(\frac{\|\Sigma^{1/2} \mathbf{v}\|_2}{2\|\mathbf{v}\|_2} - 6\sqrt{2} \sqrt{\frac{\rho(\Sigma)k \log d}{n^-}} \right)^2 \|\mathbf{v}\|_2^2, \quad (20)$$

with probability $1 - \exp(-n^-/72)$.

Notice that

$$\begin{aligned} \mathbf{v}^\top \left(\bar{\mathbf{x}}^+ (\bar{\mathbf{x}}^-)^\top + \bar{\mathbf{x}}^- (\bar{\mathbf{x}}^+)^\top \right) \mathbf{v} &= 2 (\mathbf{v}^\top \bar{\mathbf{x}}^+) (\mathbf{v}^\top \bar{\mathbf{x}}^-) \\ &\leq 2 \|\mathcal{H}_k(\bar{\mathbf{x}}^+)\|_2 \|\mathcal{H}_k(\bar{\mathbf{x}}^-)\|_2 \|\mathbf{v}\|_2^2 \end{aligned}$$

for any $\|\mathbf{v}\|_0 \leq k$. Now by Lemma 4 with $t = 16/3$ we have with probability $1 - k/2d$,

$$\begin{aligned} \|\mathcal{H}_k(\bar{\mathbf{x}}^+)\|_2^2 &\leq \frac{16\rho(\Sigma)k \log(d)}{3n_+} \\ \|\mathcal{H}_k(\bar{\mathbf{x}}^-)\|_2^2 &\leq \frac{16\rho(\Sigma)k \log(d)}{3n_-} \end{aligned}$$

Therefore

$$\mathbf{v}^\top \left(\bar{\mathbf{x}}^+ (\bar{\mathbf{x}}^-)^\top + \bar{\mathbf{x}}^- (\bar{\mathbf{x}}^+)^\top \right) \mathbf{v} \leq \frac{32}{3} \frac{\rho(\Sigma)k \log(d)}{\sqrt{n_+ n_-}} \|\mathbf{v}\|_2^2 \quad (21)$$

Now combine Equation (19), (20) and (21) and rearrange, we have for any $w, w' \in B_0(k)$,

$$\begin{aligned}
& \frac{1}{2}(\mathbf{w}' - \mathbf{w})^\top \nabla^2 F(\mathbf{w}' - \mathbf{w}) \\
& \geq \left(\left(\frac{\|\Sigma^{1/2}(\mathbf{w}' - \mathbf{w})\|_2}{2\|\mathbf{w}' - \mathbf{w}\|_2} - 6\sqrt{2}\sqrt{\frac{\rho(\Sigma)k \log d}{n^+}} \right)^2 - \frac{32}{3} \frac{\rho(\Sigma)k \log(d)}{\sqrt{n_+n_-}} \right) \|\mathbf{w}' - \mathbf{w}\|_2^2 \\
& \geq \left(\left(\frac{1}{2}\lambda_{\min}(\Sigma^{1/2}) - 6\sqrt{2}\sqrt{\frac{\rho(\Sigma)k \log d}{rn}} \right)^2 - \frac{32}{3} \frac{\rho(\Sigma)k \log(d)}{\sqrt{r(1-r)n}} \right) \|\mathbf{w}' - \mathbf{w}\|_2^2 \\
& \geq \left(\left(\frac{1}{2}\lambda_{\min}(\Sigma^{1/2}) - 6\sqrt{2}\sqrt{\frac{\rho(\Sigma)k \log d}{rn}} \right)^2 - \frac{32}{3} \frac{\rho(\Sigma)k \log(d)}{rn} \right) \|\mathbf{w}' - \mathbf{w}\|_2^2
\end{aligned}$$

with probability $(1 - \exp(-n^+/72))(1 - k/2d)$.

Now we turn to the estimate of ρ_s^\dagger . Note that the restricted smoothness condition in equation (8) is equivalent to

$$f_{B_i}(\mathbf{w}') - f_{B_i}(\mathbf{w}) - \langle \nabla f_{B_i}(\mathbf{w}), \mathbf{w}' - \mathbf{w} \rangle \leq \frac{\rho_k^\dagger}{2} \|\mathbf{w} - \mathbf{w}'\|^2$$

for all vectors \mathbf{w} and \mathbf{w}' such that $|\text{supp}(\mathbf{w}) \cup \text{supp}(\mathbf{w}')| \leq k$. Also

$$f_{B_i}(\mathbf{w}') - f_{B_i}(\mathbf{w}) - \langle \nabla f_{B_i}(\mathbf{w}), \mathbf{w}' - \mathbf{w} \rangle = \frac{1}{2}(\mathbf{w}' - \mathbf{w})^\top \nabla^2 f_{B_i}(\mathbf{w}' - \mathbf{w}),$$

since f_{B_i} is a quadratic function of \mathbf{w} . And

$$\begin{aligned}
\nabla^2 f_{B_i} &= \frac{1}{b} \sum_{j \in B_i} \left(\frac{1}{r} (\mathbf{x}_j - \bar{\mathbf{x}}_+) (\mathbf{x}_j - \bar{\mathbf{x}}_+)^\top \mathbb{I}_{[y_j=1]} + \frac{1}{1-r} (\mathbf{x}_j - \bar{\mathbf{x}}_+) (\mathbf{x}_j - \bar{\mathbf{x}}_-)^\top \mathbb{I}_{[y_j=-1]} \right. \\
&\quad \left. + (\bar{\mathbf{x}}_+ - \bar{\mathbf{x}}_+) (\bar{\mathbf{x}}_- - \bar{\mathbf{x}}_+)^\top \right) \\
&= \frac{1}{b} \sum_{j \in B_i^+} \frac{1}{r} (\mathbf{x}_j^+ - \bar{\mathbf{x}}_+) (\mathbf{x}_j^+ - \bar{\mathbf{x}}_+)^\top + \frac{1}{b} \sum_{j \in B_i^-} \frac{1}{1-r} (\mathbf{x}_j^- - \bar{\mathbf{x}}_+) (\mathbf{x}_j^- - \bar{\mathbf{x}}_-)^\top \\
&\quad + (\bar{\mathbf{x}}_- - \bar{\mathbf{x}}_+) (\bar{\mathbf{x}}_- - \bar{\mathbf{x}}_+)^\top
\end{aligned}$$

Let $\tilde{\mathbf{x}}_j \in \{\mathbf{x}_j^+ - \bar{\mathbf{x}}_+, \mathbf{x}_j^- - \bar{\mathbf{x}}_-, \bar{\mathbf{x}}_- - \bar{\mathbf{x}}_+\}$, since $\|\mathbf{v}\|_0 \leq k$, we have

$$\mathbf{v} \tilde{\mathbf{x}}_j \tilde{\mathbf{x}}_j^\top \mathbf{v} = (\mathbf{v}^\top \tilde{\mathbf{x}}_j)^2 \leq \|\mathcal{H}_k(\tilde{\mathbf{x}}_j)\|_2^2 \|\mathbf{v}\|_2^2.$$

Also notice that by the properties of mean and variance, we have $\mathbf{x}_j^+ - \bar{\mathbf{x}}_+ \sim N\left(0, \frac{n_+-1}{n_+} \Sigma\right)$, $\mathbf{x}_j^- - \bar{\mathbf{x}}_- \sim N\left(0, \frac{n_- - 1}{n_-} \Sigma\right)$ and $\bar{\mathbf{x}}_- - \bar{\mathbf{x}}_+ \sim N\left(0, \left(\frac{1}{n_+} + \frac{1}{n_-}\right) \Sigma\right)$. Hence, by Lemma 4, pick $t = 16 \log(2\sqrt{bd})/3 \log(2d)$ we have with probability $1 - k/2bd$,

$$\|\mathcal{H}_k(\mathbf{x}_j^+ - \bar{\mathbf{x}}_+)\|_2^2 \leq \frac{16}{3} \rho(\Sigma) \left(\frac{n_+ - 1}{n_+}\right) k \log(d) \left(\frac{1}{2} \log(b) + \log(d)\right),$$

$$\|\mathcal{H}_k(\mathbf{x}_j^- - \bar{\mathbf{x}}_-)\|_2^2 \leq \frac{16}{3} \rho(\Sigma) \left(\frac{n_- - 1}{n_-}\right) k \log(d) \left(\frac{1}{2} \log(b) + \log(d)\right),$$

$$\|\mathcal{H}_k(\bar{\mathbf{x}}_+ - \bar{\mathbf{x}}_-)\|_2^2 \leq \frac{16}{3} v \left(\frac{1}{n_+} + \frac{1}{n_-}\right) k \log(d) \left(\frac{1}{2} \log(b) + \log(d)\right).$$

Define the imbalance ratio in each batch $r_i = \frac{b_i^+}{b}$, therefore we have for any $\mathbf{w}, \mathbf{w}' \in B_0(k)$ with probability $1 - k/2d$,

$$\begin{aligned}
& \max_i \left\{ \frac{1}{2} (\mathbf{w} - \mathbf{w}')^\top \nabla^2 f_{B_i} (\mathbf{w} - \mathbf{w}') \right\} \\
& \leq \max_i \left\{ \frac{16}{3} \left(\rho(\Sigma) \frac{r_i(n_+ - 1)}{rn_+} + \rho(\Sigma) \frac{(1 - r_i)(n_- - 1)}{(1 - r)n_-} + \rho(\Sigma) \frac{1}{n_+} + \rho(\Sigma) \frac{1}{n_-} \right) k \log(d) \right. \\
& \quad \left. \times \left(\frac{1}{2} \log(b) + \log(d) \right) \|\mathbf{w} - \mathbf{w}'\|_2^2 \right\} \\
& \leq \max_i \left\{ \frac{16}{3} \rho(\Sigma) \left(\frac{r_i}{r} + \frac{1 - r_i}{1 - r} + \frac{1}{rn} + \frac{1}{(1 - r)n} \right) k \log(d) \left(\frac{1}{2} \log(b) + \log(d) \right) \|\mathbf{w} - \mathbf{w}'\|_2^2 \right\} \\
& \leq \frac{16}{3} \rho(\Sigma) \left(\frac{1}{r} + \frac{1}{rn} + \frac{1}{(1 - r)n} \right) k \log(d) \left(\frac{1}{2} \log(b) + \log(d) \right) \|\mathbf{w} - \mathbf{w}'\|_2^2 \\
& \leq \frac{16\rho(\Sigma)k \log(d) \left(\frac{1}{2} \log(b) + \log(d) \right)}{r} \|\mathbf{w} - \mathbf{w}'\|_2^2
\end{aligned}$$

where in the third and last inequality we use the fact that $1 - r \geq r$. □

D. Full formula of $\frac{\sigma_{\mathbf{w}_*}}{1 - \kappa}$ in Section III-B

The derivation of $\sigma_{\mathbf{w}_*}$ is similar to the derivation of ρ_s^+ . Hence it is omitted here. The full formula is given as follow.

$$\frac{\sigma_{\mathbf{w}_*}}{1 - \kappa} = \frac{4r \|\mathbf{w}_*\|_2 + \sqrt{\frac{r}{2n\rho(\Sigma)(2k+k_*) \log(d)}}}{1 - \sqrt{(1 + \nu) \left(1 - \frac{3\lambda^2}{128k \log(d)} r + \left(\frac{9\sqrt{2}\lambda}{16\sqrt{k \log(d)n_-}} + \frac{1}{n} \right) \sqrt{r} - \frac{27}{4n} \right)}}$$