

Iterative Thresholding for Non-Linear Learning in the Strong ϵ -Contamination Model

Arvind Rathnashyam
RPI Math and CS, rathna@rpi.edu

Alex Gittens
RPI CS, gittea@rpi.edu

Abstract

We study the problem of learning single neurons when both the labels and the covariates are possibly corrupted adversarially with a gradient-descent iterative thresholding algorithm. We assume the data is sampled from a ground truth distribution,

$$y = \psi(\mathbf{x}^\top \mathbf{w}^*) + \xi$$

where ξ is Gaussian noise and σ is an activation-function. We consider activations with three different properties, (i) Lipschitz, (ii) lower bounded first derivative, and (iii) second derivative has finitely many zeros. Functions such as tanh, sigmoid, leaky-ReLU, and ReLU are functions with a subset of the properties.

We also study the linear regression problem when $\psi(x) = x$. We improve upon previous approximation bounds for gradient based iterative thresholding algorithms [BJK15, SS19] and show with high probability a $O(\sigma\epsilon \log \epsilon^{-1})$ approximation upper bound, matching the best known approximation bound for iterative thresholding algorithms [ADKS22] while improving run-time. We extend previous works by studying the case when \mathbf{x} is sampled from a sub-Gaussian distribution with a general covariance matrix Σ .

1 Introduction

There has been extensive study of algorithms to learn the target distribution from a Huber ϵ -Contaminated Model for a Generalized Linear Model (GLM), [DKK⁺19, ADKS22, LBSS21, OZS20, FB81] as well as for linear regression [BJKK17, MGJK19]. Robust Statistics has been studied extensively [DK23] for problems such as high-dimensional mean estimation [PBR19, CDGS20] and Robust Covariance Estimation [CDGW19, FWZ18]. Recently, there has been an interest in solving robust machine learning problems by gradient descent [PSBR18, DKK⁺19]. Subquantile minimization aims to address the shortcomings of standard ERM in applications of noisy/corrupted data [KLA18, JZL⁺18]. In many real-world applications, the covariates have a non-linear dependence on labels [AMMIL12, Section 3.4]. In which case it is suitable to transform the covariates to a different space utilizing kernels [HSS08]. Therefore, in this paper we consider the problem of Robust Learning for Kernel Learning.

Definition 1 (Strong ϵ -Contamination Model [HR09]). *Given a corruption parameter $0 \leq \epsilon < 0.5$, a data matrix, X and labels \mathbf{y} . An adversary is allowed to inspect all samples and modify ϵn samples arbitrarily. The algorithm is then given the ϵ -corrupted data matrix X and ϵ -corrupted labels vector \mathbf{y} as training data.*

Current approaches for robust learning across various machine learning tasks often use gradient descent over a robust objective, [LBSS21]. These robust objectives tend to not be convex and therefore do not have a strong analysis on the error bounds for general classes of models.

We similarly propose a robust objective which has a nonconvex-concave objective. This objective function has also been proposed recently in [HYwL20] where there has been an analysis in the Binary Classification Task. We show Subquantile Minimization reduces to the same objective function given in [HYwL20].

The study of Kernel Learning in the Gaussian Design is quite popular, [CLKZ21, Dic16]. In [CLKZ21], the feature space, $\phi(\mathbf{x}_i) \sim \mathcal{N}(0, \Sigma)$ where Σ is a diagonal matrix of dimension p , where p can be infinite. We will now give our formal definition of the dataset.

Definition 2 (Corruption Model). *Let \mathcal{P} be a distribution over \mathbb{R}^d such that $\mathcal{P}_\# \phi$ is a centered distribution in the Hilbert Space \mathcal{H} with trace-class covariance operator Σ and trace-class sub-Gaussian proxy Γ such that $\Sigma \preceq c\Gamma$. The original dataset is denoted as \hat{P} , the adversary is able to observe \hat{P} and arbitrarily corrupts ϵn samples denoted as Q such that $|Q| = \epsilon n$. The remaining uncorrupted samples are denoted as P such that $|P| = n(1 - \epsilon)$. Together $X \triangleq P \cup Q$ represents the given dataset.*

We will now give one of the first results proving the effectiveness of Iterative Thresholding in Learning Problems.

Theorem 3 (Theorem 5 in [BJK15]). *Let X be a sub-Gaussian data matrix, and $\mathbf{y} = X^\top \mathbf{w}^* + \mathbf{b}$ where \mathbf{b} represents the corruption. Then there exists an algorithm such that $\|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2 \leq \epsilon$ after $t = O\left(\left(\log\left(\frac{\|\mathbf{b}\|_2}{\sqrt{n}}\right)\right) \frac{1}{\epsilon}\right)$ iterations.*

More recently, Awasthi et al. [ADKS22] studied the iterative trimmed maximum likelihood estimator. In their algorithm, at each step they find \mathbf{w}^* which minimizes the elements in the trimmed set. We will give their formal theorem result.

Theorem 4 (Theorem 4.2 in [ADKS22]). *Let $P = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ be the data generated by a Gaussian regression model defined as $y_i = \mathbf{x}_i^\top \mathbf{w}^* + \eta_i$ where $\eta_i \sim \mathcal{N}(0, \sigma^2)$ and \mathbf{x}_i are sampled from a sub-Gaussian Distribution with second-moment matrix I . Suppose the dataset has ϵ -fraction of label corruption and $n = \Omega\left(\frac{d + \log(1/\delta)}{\epsilon^2}\right)$. Then there exists an algorithm that returns $\hat{\mathbf{w}}$ such that with probability $1 - \delta$,*

$$\|\hat{\mathbf{w}} - \mathbf{w}^*\|_2 = O(\sigma \epsilon \log(1/\epsilon))$$

Our first result recovers this result for vectorized-regression. We will now give our results for the Kernelized GLM problem.

1.1 Contributions

Our main contribution is the approximation bounds for Subquantile Minimization for various non-linear learning problems from the iterative thresholding algorithm given in Algorithm 1. Our proof techniques extend [BJK15, SS19, ADKS22] as we suppose the adversary also corrupts the covariates. To our knowledge, we are also the first to theoretically study iterative thresholding for non-linear learning algorithms beyond the generalized linear model.

Reference	Approximation	Runtime	Model
[BJK15]	$O(\sigma)$	$O\left(N^2 d \log\left(\frac{1}{\sqrt{n}} \frac{\ \mathbf{b}\ _2}{\epsilon}\right)\right)$	Regression; Label
[SS19]	$O(\sigma)$	$O\left(N d \log\left(\frac{\ \mathbf{w}^*\ _2 + \sigma^2}{\sigma}\right)\right)$	Regression; Label
[ADKS22]	$O(\sigma \epsilon \log \epsilon^{-1})$	$O\left((N d^2 + d^3) \frac{1}{\sigma \epsilon^2}\right)$	Regression; Label
Corollary 18	$O(\sigma \epsilon \log \epsilon^{-1})$	$O\left(N d^2 \log\left(\frac{\ \mathbf{w}^*\ }{\sigma \epsilon \log \epsilon^{-1}}\right)\right)$	Regression; Label
Theorem 17	$O(\sigma \epsilon \sqrt{C_Q K} \log \epsilon^{-1})$	$O\left(N K d \log\left(\frac{\ W^*\ _F}{\sigma \epsilon \log \epsilon^{-1}}\right)\right)$	K -variate Regression
Theorem 19	$O(L^2(\sigma + \epsilon \sqrt{C_Q} \log \epsilon^{-1}))$	$O\left(N d \log\left(\frac{\ f^*\ _{\mathcal{H}}}{\sigma \epsilon \log \epsilon^{-1}}\right)\right)$	GLM (<i>In progress</i>)
Theorem 21	$O\left(\sigma \kappa^2(\Sigma) C_\psi^{-1} C_K \sqrt{\epsilon \log \epsilon^{-1}}\right)$	$O\left(N d \log\left(\frac{\ \mathbf{w}^*\ _2 + (d/\delta)}{\epsilon}\right)\right)$	Leaky ReLU Neuron
Theorem 23	$O\left(\sigma \kappa^2(\Sigma) C_\psi^{-1} C_K d \epsilon \log \epsilon^{-1}\right)$	$O\left(N d \log\left(\frac{\ \mathbf{w}^*\ _2 + (d/\delta)}{\epsilon}\right)\right)$	ReLU Neuron

Table 1: Summary of related work on Iterative Thresholding Algorithms for Learning in the Huber- ϵ Contamination Model and our contributions. We assume the good data is sampled from a sub-Gaussian distribution with second-moment matrix, Σ , and sub-Gaussian norm C_K and dimension d . We assume the variance of the optimal estimator is σ . The Leaky-ReLU function is given as $\max\{C_\psi x, x\}$.

We are able to remove the super-linear dependence on d by leveraging Gradient Descent. Comparing to [BJK15], our bound does not depend on the norm of the noise.

2 Preliminaries

Notation. We denote $[T]$ as the set $\{1, 2, \dots, T\}$. We define $(x)^+ \triangleq \max(0, x)$ as the Rectified Linear Unit (ReLU) function. We say $y = O(x)$ if there exists x_0 s.t. for all $x \geq x_0$ there exists C s.t. $y \leq Cx$. We say $y = \Omega(x)$ if there exists x_0 s.t. for all $x \geq x_0$ there exists C s.t. $y \geq Cx$. We denote $a \vee b \triangleq \max(a, b)$ and $a \wedge b \triangleq \min(a, b)$. We define \mathbb{S}^{d-1} as the sphere $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| = 1\}$. We will typically denote capital roman letters A, B, C as fixed constants, lower-case roman letters f, g, h as functions, bold-face lower-case letters $\mathbf{x}, \mathbf{y}, \mathbf{z}$ as vectors, and bold-face, upper-case letters P, Q, R as matrices. We denote the Hadamard product between two vectors of the same size as $\mathbf{x} \circ \mathbf{y}$ such that for any vectors $(\mathbf{x} \circ \mathbf{y})_i = x_i y_i$.

Matrices. For a matrix A , let $\lambda_{\max}(A)$ and $\lambda_{\min}(A)$ represent the maximum and minimum eigenvalues of A , respectively. We use the following matrix norms for a matrix $A \in \mathbb{R}^{m \times n}$,

$$\text{Spectral Norm: } \|A\| = \max_{\mathbf{x} \in \mathbb{S}^{m-1}} \|A\mathbf{x}\| = \sigma_1(A)$$

$$\text{Trace Norm: } \text{Tr}(A) = \sum_{i \in [m \wedge n]} \sigma_i(A)$$

$$\text{Frobenius Norm: } \|A\|_F^2 = \text{Tr}(A^\top A) = \sum_{i \in [m \wedge n]} \sigma_i^2(A)$$

Let $\text{vec} : \mathbb{R}^{n \times k} \rightarrow \mathbb{R}^{nk}$ represent the vectorization of a matrix to a vector placing its columns one by one into a vector. We then have the useful facts,

Lemma 5. Suppose $A, B \in \mathbb{R}^{m \times n}$, then

$$\langle A, B \rangle_{\text{Tr}} = \langle \text{vec}(A), \text{vec}(B) \rangle$$

Let $\otimes : \mathbb{R}^{N \times K} \times \mathbb{R}^{L \times M} \rightarrow \mathbb{R}^{NL \times KM}$ represent the Kronecker delta product between two matrices, this gives us the following relation,

Lemma 6. *Suppose A, B, C are size compatible matrices, then*

$$\text{vec}(ABC) = (C^\top \otimes A)\text{vec}(B)$$

Probability. We know discuss the probability theory concepts used throughout the paper. We consider the general sub-Gaussian design. The sub-Gaussian design is highly prevalent in the study of robust statistics [JLT20].

Definition 7 (Sub-Gaussian Distribution). *We say a vector \mathbf{x} is sampled from a sub-Gaussian distribution with second-moment matrix Σ and sub-Gaussian norm K , if $\mathbf{E}[\mathbf{x}\mathbf{x}^\top] = \Sigma$ and*

$$\mathbf{E}[\exp(t\mathbf{x}^\top \mathbf{v})] \leq \exp\left(\frac{t^2 K^2 \mathbf{v}^\top \Sigma \mathbf{v}}{2}\right) \text{ for } \mathbf{v} \in \mathbb{S}^{d-1}, t \in \mathbb{R}$$

A scalar random variable X is sub-Gaussian with sub-Gaussian norm ν if for all $p \in \mathbb{N}$,

$$\|X\|_{L_p} = (\mathbf{E}|X|^p)^{1/p} \leq \nu \sqrt{p}$$

We often work with the products of sub-Gaussian random variables, which by the following indicate they are sub-Exponential.

Lemma 8 (Lemma 2.7.7 in [Ver20]). *Let X, Y be sub-Gaussian random variables, then XY is sub-Exponential, furthermore,*

$$\|XY\|_{\psi_1} \leq \|X\|_{\psi_2} \|Y\|_{\psi_2}$$

To give probabilistic bounds on the concentration of sub-Exponential random variables, we often utilize Bernstein's Theorem.

Lemma 9 (Theorem 1.13 in [RH23]). *Let X_1, \dots, X_N be independent sub-Exponential random variables such that $\|X_i\|_{\psi_1} = \nu$ for $i \in [N]$, then*

$$\Pr\left\{\frac{1}{N} \sum_{i \in [N]} X_i \geq t\right\} \leq \exp\left(-\frac{N}{2} \left(\frac{t^2}{\nu^2} \wedge \frac{t}{\nu}\right)\right)$$

2.1 Related Work

The idea of iterative thresholding algorithms for robust learning tasks dates back to 1806 by Legendre [Leg06]. Iterative thresholding have been studied theoretically and tested empirically in various machine learning domains [HYW⁺23, MGJK19].

[BJK15] study iterative thresholding for least squares regression / sparse recovery. In particular, one of their contributions is a gradient descent algorithm, TORRENT, when the covariates are sampled from a sub-Gaussian distribution. Their approximation bound in Theorem 5 relies on the fact that $\lambda_{\min}(\Sigma) = \lambda_{\max}(\Sigma)$ and with sufficiently large data and sufficiently small ϵ , $\kappa(X) \searrow 1$. Bhatia et al. also study a full solve algorithm, where after each thresholding step and obtaining $(1 - \epsilon)N$ samples, they set $\mathbf{w}^{(t)}$ to the minimizer of the squared loss over the $(1 - \epsilon)N$ points and refer to this algorithm as TORRENT-FC. They study this algorithm in the presence of both adversarial and intrinsic noise in the optimal estimator. Their analysis in Corollary 11 gives $O(\sigma)$ error when the intrinsic noise is sub-Gaussian with norm σ^2 .

[SS19] study iterative thresholding for learning Generalized Linear Models (GLMs). In both the linear and non-linear case, they present results on linear convergence. Their results imply a bound of $O(\sigma)$ in the linear case. They furthermore provide experimental evidence of the success of iterative thresholding when applied to neural networks.

More recently, [ADKS22] studied iterative trimmed maximum likelihood estimator for General Linear Models. They prove the best known bounds for iterative thresholding algorithms in the linear regression case. The algorithm studied by [ADKS22] utilizes a filtering algorithm with the sketching matrix $\Sigma^{-1/2}$ so the columns of X are sampled from a sub-Gaussian Distribution with multivariate proxy I before running the iterative thresholding procedure.

3 Subquantile Minimization

We propose to optimize over the subquantile of the risk. The p -quantile of a random variable, U , is given as $\mathcal{Q}_p(U)$, this is the largest number, t , such that the probability of $U \leq t$ is at least p .

$$\mathcal{Q}_p(U) \leq t \iff \mathbf{Pr}\{U \leq t\} \geq p$$

The p -subquantile of the risk is then given by

$$L_p(U) = \frac{1}{p} \int_0^p \mathcal{Q}_p(U) dq = \mathbf{E}[U|U \leq \mathcal{Q}_p(U)] = \max_{t \in \mathbb{R}} \left\{ t - \frac{1}{p} \mathbf{E}[(t - U)^+] \right\}$$

Given a function to minimize, \mathcal{L} , the variational problem is given as:

$$\min_{f \in \mathcal{K}} \max_{t \in \mathbb{R}} \left\{ g(t, \mathbf{w}) \triangleq t - \frac{1}{(1-\epsilon)N} \cdot \sum_{i=1}^n (t - \mathcal{L}(\mathbf{w}; \mathbf{x}_i, y_i))^+ \right\}$$

where t is the p -quantile of the empirical risk. Note that for a fixed t therefore the objective is not concave with respect to \mathbf{w} . Thus, to solve this problem we use the iterations from Equation 11 in [RHL⁺20]. Our update steps are,

$$t^{(k+1)} = \arg \max_{t \in \mathbb{R}} g(\mathbf{w}^{(k)}, t) \quad (3.1)$$

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \eta \nabla_f g(\mathbf{w}^{(k)}, t^{(k+1)}) \quad (3.2)$$

3.1 Reduction to Iterative Thresholding

To consider theoretical guarantees of Subquantile Minimization, we first analyze the inner and outer optimization problems. We first analyze kernel learning in the presence of corrupted data. Next, we provide error bounds for the two most important kernel learning problems, kernel ridge regression, and kernel classification. Now we will give our first result regarding kernel learning in the Huber ϵ -contamination model. Now we will analyze the two-step minimax optimization steps described in Equations (3.1) and (3.2).

Lemma 10. *Let $\mathcal{R} : \mathcal{H} \times \mathbb{R} \rightarrow \mathbb{R}$ be a loss function (not necessarily convex). Let $\mathbf{x}_{[i]}$ represent the point with the i -th smallest loss w.r.t \mathcal{R} . If we denote $\hat{v}_i \triangleq \mathcal{R}(\mathbf{w}; \mathbf{x}_{[i]}, y_{[i]})$, it then follows $\hat{v}_{(1-\epsilon)N} \in \arg \max_{t \in \mathbb{R}} g(\mathbf{w}, t)$.*

Proof. First we can note, the max value of t for g is equivalent to the min value of t for the convex w.r.t t function $-g$. We can now find the Fermat Optimality Conditions for g .

$$\partial(-g(t, f)) = \partial \left(-t + \frac{1}{(1-\epsilon)N} \sum_{i=1}^N (t - \hat{v}_i)^+ \right) = -1 + \frac{1}{(1-\epsilon)N} \sum_{i=1}^{(1-\epsilon)N} \begin{cases} 1 & \text{if } t > \hat{v}_i \\ 0 & \text{if } t < \hat{v}_i \\ [0, 1] & \text{if } t = \hat{v}_i \end{cases}$$

We observe when setting $t = \hat{v}_{(1-\epsilon)N}$, it follows that $0 \in \partial(-g(t, f))$. This is equivalent to the $(1-\epsilon)$ -quantile of the Empirical Risk. \blacksquare

From Lemma 10, we see that t will be greater than or equal to the errors of exactly $(1-\epsilon)N$ points. Thus, we are continuously updating over the $(1-\epsilon)N$ minimum errors.

Lemma 11. *Let $\hat{v}_i \triangleq \mathcal{R}(\mathbf{w}; \mathbf{x}_{[i]}, y_{[i]})$, if we choose $t^{(k+1)} = \hat{v}_{(1-\epsilon)N}$ as by Lemma 10, it then follows $\nabla_f g(t^{(k)}, \mathbf{w}^{(k)}) = \frac{1}{(1-\epsilon)N} \sum_{i=1}^{(1-\epsilon)N} \nabla_{\mathbf{w}} \mathcal{R}(\mathbf{w}^{(k)}; \mathbf{x}_{[i]}, y_{[i]})$.*

Proof. By our choice of $t^{(k+1)}$, it follows,

$$\begin{aligned} \partial_f g(t^{(k+1)}, \mathbf{w}^{(k)}) &= \partial_f \left(t^{(k+1)} - \frac{1}{(1-\epsilon)N} \sum_{i=1}^N (t^{(k+1)} - \mathcal{R}(\mathbf{w}^{(k)}; \mathbf{x}_{[i]}, y_{[i]}))^+ \right) \\ &= -\frac{1}{(1-\epsilon)N} \sum_{i=1}^{(1-\epsilon)N} \partial_f (t^{(k+1)} - \mathcal{R}(\mathbf{w}^{(k)}; \mathbf{x}_{[i]}, y_{[i]}))^+ \end{aligned}$$

$$= \frac{1}{(1-\epsilon)N} \sum_{i=1}^n \nabla_f \mathcal{R}(\mathbf{w}^{(k)}; \mathbf{x}_{[i]}, y_{[i]}) \begin{cases} 1 & \text{if } t > \hat{\nu}_i \\ 0 & \text{if } t < \hat{\nu}_i \\ [0, 1] & \text{if } t = \hat{\nu}_i \end{cases}$$

Now we note $\hat{\nu}_{(1-\epsilon)N} \leq t^{(k+1)} \leq \hat{\nu}_{(1-\epsilon)N+1}$. Then, we have

$$\partial_f g(t^{(k+1)}, \mathbf{w}^{(k)}) \ni \frac{1}{(1-\epsilon)N} \sum_{i=1}^{(1-\epsilon)N} \nabla_f \mathcal{R}(\mathbf{w}^{(k)}; \mathbf{x}_{[i]}, y_{[i]})$$

This concludes the proof. ■

Theorem 12. *Suppose there exists constants such that the Hessian of any $(1-\epsilon)$ -subset of the data w.r.t to \mathcal{L} is μ -strongly convex and L -smooth. Then there exists an algorithm that returns a stationary point of $g(t, \mathbf{w})$ in $O(\text{poly}(\mu, L))$.*

Proof. The proof is by construction, consider the algorithm described in Equations (3.1) and (3.2). From the proof of Lemma 10, $t = \hat{\nu}_{(1-\epsilon)N}$ is a maximizer of Equation (3.2). We can then see that $g(t^{(k+1)}, \mathbf{w}^{(k+1)}) \leq g(t^{(k)}, \mathbf{w}^{(k)})$. ■

4 Convergence

In this section we give the algorithm for subquantile minimization. We will start with the simple case of vectorized regression as a warm-up to our general proof technique. We then move to the GLM with kernel learning. Finally, we give our results for one-hidden layer neural networks. We will now give the algorithm for Subquantile Minimization with Gradient Descent. We first give the non-linear functions we will be learning.

Property 13. *The non-linear function, $\psi : \mathbb{R} \rightarrow \mathbb{R}$ has lower bounded derivative, $\psi'(x) \geq C_{\lfloor \psi \rfloor}$ for all $x \in \mathbb{R}$.*

Property 14. *The non-linear function, $\psi : \mathbb{R} \rightarrow \mathbb{R}$ is $C_{\lceil \psi \rceil}$ -Lipschitz.*

Property 15. *The second-derivative, $\psi'' : \mathbb{R} \rightarrow \mathbb{R}$ is non-zero at finitely many points.*

4.1 Algorithm

We will first define the thresholding operator to simplify the notation in our formal algorithm.

Definition 16 (Hard Thresholding Operator in [BJK15]). *For any vector $\mathbf{v} \in \mathbb{R}^n$, let $\sigma_{\mathbf{v}}$ be the permutation that orders elements in ascending order, i.e. $\mathbf{v}_{\sigma_{\mathbf{v}}(1)} \leq \mathbf{v}_{\sigma_{\mathbf{v}}(2)} \leq \dots \leq \mathbf{v}_{\sigma_{\mathbf{v}}(n)}$. Then for any $k \leq n$, the hard thresholding operator is defined as,*

$$\text{HT}(\mathbf{v}; k) = \{i \in [n] : \sigma_{\mathbf{v}}^{-1}(i) \leq k\}$$

We now give the gradient descent algorithm which we will study for the remainder of this paper.

Runtime. In each iteration we calculate the ℓ_2 error for N points, in total $O(Nd)$. For the Hard Thresholding step, it suffices to find the $n(1-\epsilon)$ -th largest element, we can run a selection algorithm in worst-case time $O(N \log N)$, then partition the data in $O(N)$. The run-time for calculating the gradient and updating $\mathbf{w}^{(t)}$ is dominated by the matrix multiplication in $X_{S^{(t)}} X_{S^{(t)}}^T$ which can be done in $O(Nd^2)$. Then considering the choice of T , we have the algorithm runs in time $O\left(N^2 d \log\left(\frac{\|\mathbf{w}^*\|_2}{\sigma \epsilon \log \epsilon^{-1}}\right)\right)$ to obtain $O(\sigma \epsilon \log \epsilon^{-1})$ ℓ_2 -approximation error.

Algorithm 1 Subquantile Minimization for Learning a Neuron

input: Possibly corrupted $X \in \mathbb{R}^{d \times N}$ with outputs $\mathbf{y} \in \mathbb{R}^N$, activation function σ , and corruption parameter $\epsilon = O(1)$.

output: ϵ -Approximate solution $\mathbf{w} \in \mathbb{R}^d$ to minimize $\|\mathbf{w} - \mathbf{w}^*\|_2$.

- 1: $\mathbf{w}^{(0)} \leftarrow \mathbf{0}$
- 2: $\eta \leftarrow 0.1\kappa^{-2}(\Sigma)$
- 3: $T \leftarrow O\left(\kappa^2(\Sigma) \log\left(\frac{\|\mathbf{w}^*\|_2}{\epsilon}\right)\right)$
- 4: **for** $t \in [T]$ **do**
- 5: $\nu_i^{(t)} = (\sigma(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - y_i)^2$
- 6: $S^{(t)} \leftarrow \text{HT}(\nu^{(t)}, (1 - \epsilon)N)$
- 7: $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; S^{(t)})$

return: $\mathbf{w}^{(T)}$

4.2 Proof Sketch

In this section we will give a general sketch of our proofs, from which all the individual theorems will be based upon. Let $t \in [T]$, we then have,

$$\begin{aligned} \|\mathbf{w}^{(t+1)} - \mathbf{w}^*\| &= \|\text{Proj}_\Theta[\mathbf{w}^{(t)} - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; S^{(t)})] - \mathbf{w}^*\|_2 \\ &\stackrel{(i)}{\leq} \|\mathbf{w}^{(t)} - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; S^{(t)}) - \mathbf{w}^*\|_2 \\ &\leq \|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; S^{(t)} \cap P)\| + \|\nabla \mathcal{R}(\mathbf{w}^{(t)}; S^{(t)} \cap Q)\|_2 \end{aligned}$$

In the above, (i) follows from the assumption that $\mathbf{w}^* \in \Theta$, therefore the projection operator is on to a convex set and is thus must be closer or at the same ℓ_2 norm distance to \mathbf{w}^* . We analyze the first term through its square,

$$\|\mathbf{w}^{(t)} - \mathbf{w}^* - \nabla \mathcal{R}(\mathbf{w}^{(t)}; S^{(t)} \cap P)\|^2 = \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2 - 2\eta \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \nabla \mathcal{R}(\mathbf{w}^{(t)}; S^{(t)} \cap P) \rangle + \eta^2 \|\nabla \mathcal{R}(\mathbf{w}^{(t)}; S^{(t)} \cap P)\|^2$$

In this step the finer details of the particular proof will differ, however the structure remains the same. We will prove there exists some $c_1 > 0$ and $c_2 \geq c_1$ such that,

$$\langle \mathbf{w}^{(t)} - \mathbf{w}^*, \nabla \mathcal{R}(\mathbf{w}^{(t)}; S^{(t)} \cap P) \rangle \geq (1 - 2\epsilon)c_1 \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2$$

and

$$\|\nabla \mathcal{R}(\mathbf{w}^{(t)}; S^{(t)} \cap P)\| \leq (1 - \epsilon)c_2 \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2$$

Then, solving a simple quadratic equation, we have that $\eta^2 \beta \leq \eta c_1 c_3$ for a $c_3 \in (0, 1)$ when we choose $\eta \leq \frac{c_1 c_3}{c_2}$ and we are able to eliminate the norm of the gradient squared term. We must now control the corrupted gradient term. The key idea is to note that from the optimality of the sub-quantile set,

$$\sum_{i \in S^{(t)} \cap Q} \mathcal{L}(\mathbf{w}^{(t)}; \mathbf{x}_i, y_i) \leq \sum_{i \in P \setminus S^{(t)}} \mathcal{L}(\mathbf{w}^{(t)}; \mathbf{x}_i, y_i)$$

We then prove the existence of a constant c_4 such that,

$$\|\nabla \mathcal{R}(\mathbf{w}^{(t)}; S^{(t)} \cap Q)\|_2 \leq \epsilon c_4 \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2$$

Then, combining our results, we end up with a linear convergence of the form,

$$\|\mathbf{w}^{(t+1)} - \mathbf{w}^*\|_2 \leq \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2 (1 - \eta(1 - 2\epsilon)c_1 + \eta \epsilon c_4)$$

4.3 Warm-up: Multivariate Linear Regression

We will first present our results for the well-studied problem of linear regression in the Huber- ϵ contamination model. Our results will extend the results in [BJKK17] Theorem 5 and [ADKS22] Lemma A.1 by including covariate corruption without requiring a filtering algorithm, variance in the optimal estimator, and non-identity second-moment matrix of the uncorrupted data. The loss function for the multivariate linear regression problem for $W \in \mathbb{R}^{k \times d}$, $X \in \mathbb{R}^{d \times n}$, and $Y \in \mathbb{R}^{k \times n}$.

$$\mathcal{L}(W; X, Y) = \|WX - Y\|_F^2$$

Theorem 17. *Let $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{d \times N}$ be the data matrix and $Y = [\mathbf{y}_1, \dots, \mathbf{y}_n] \in \mathbb{R}^{K \times N}$ be the output, such that for $i \in P$, \mathbf{x}_i are sampled from a sub-Gaussian distribution with sub-Gaussian norm L and second-moment matrix Σ . Suppose for $i \in P$ the output is given as $\mathbf{y}_i = W\mathbf{x}_i + \mathbf{e}_i$ where for $j \in [K]$, $[\mathbf{e}_i]_j \sim \mathcal{N}(0, \sigma^2)$. Then after $O\left(\kappa(\Sigma) \log\left(\frac{\|W^*\|_F}{\epsilon}\right)\right)$ gradient descent iterations, $N \geq \frac{(1/\delta)K \text{Tr}(\Sigma)}{1600\epsilon^2 \log \epsilon^{-1} \lambda_{\max}(\Sigma)}$, and learning rate $\eta = 0.1\lambda_{\max}^{-2}(\Sigma)$, with probability exceeding $1 - \delta$, Algorithm 1 returns $W^{(T)}$ such that*

$$\|W^{(T)} - W^*\|_F \leq \epsilon + O\left(\sigma\epsilon\sqrt{KC_3 \log \epsilon^{-1}}\right)$$

Proof. The proof is deferred to § A.1. ■

We are able to recover the result of Lemma 4.2 in [ADKS22] when $K = 1$ and the covariates (corrupted and un-corrupted) are sampled from a sub-Gaussian distribution with second-moment matrix I . The full solve algorithm studied in [ADKS22] returns a $O(\sigma\epsilon \log \epsilon^{-1})$ in time $O(\frac{1}{\epsilon^2}(Nd^2 + d^3))$ and the same algorithm studied in [BJK15], TORRENT-FC obtains $O(\sigma)$ approximation error in run-time $O\left(\log\left(\frac{1}{\sqrt{n}} \frac{\|\mathbf{b}\|}{\sigma\epsilon \log \epsilon^{-1}}\right)(Nd^2 + d^3)\right)$, with the gradient descent based approach, we are able to improve the runtime to $O\left(\log\left(\frac{\|\mathbf{w}^*\|}{\sigma\epsilon \log \epsilon^{-1}}\right)Nd^2\right)$ for the same approximation bound. By no longer requiring the full-solve, we are able to remove super-linear relation to d . In comparison to [BJK15], we do not have dependence on the noise vector \mathbf{b} , which can have very large norm in relation to the norm of \mathbf{w}^* . Our proof is also a significant improvement over the presentation given in Lemma 5 of [SS19] as under the same conditions, we give more than the linear convergence, but we show linear convergence is possible on any second-moment matrix of the good covariates and covariate corruption, and then develop concentration inequality bounds to match the best known result for iteratively trimmed estimators. We will formalize our results into a corollary to give a more representative comparison in the literature.

Corollary 18. *Let $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{d \times N}$ be the data matrix and $\mathbf{y} = [y_1, \dots, y_n] \in \mathbb{R}^N$ be the output, such that for $i \in P$, \mathbf{x}_i are sampled from a sub-Gaussian distribution with sub-Gaussian norm L and second-moment matrix I . Suppose for $i \in P$ the output is given as $\mathbf{y}_i = \mathbf{x}_i^\top \mathbf{w}^* + \xi_i$ where $\xi_i \sim \mathcal{N}(0, \sigma^2)$. Then after $O\left(\log\left(\frac{\|\mathbf{w}^*\|_2}{\epsilon}\right)\right)$ gradient descent iterations, $N \geq \frac{(d/\delta)}{800\epsilon^2}$, and learning rate $\eta = 0.1$, with probability exceeding $1 - \delta$, Algorithm 1 returns $\mathbf{w}^{(T)}$ such that*

$$\|\mathbf{w}^{(T)} - \mathbf{w}^*\|_2 \leq \epsilon + O\left(\sigma\epsilon\sqrt{KC_3 \log \epsilon^{-1}}\right)$$

Suppose for $i \in Q$, \mathbf{x}_i are sampled from a sub-Gaussian distribution with sub-Gaussian Norm K and second-moment matrix I . Then,

$$\|\mathbf{w}^{(T)} - \mathbf{w}^*\|_2 \leq \epsilon + O\left(\sigma\epsilon \log \epsilon^{-1}\right)$$

The second relation given in Corollary 18 matches the best known bound for robust linear regression with iterative thresholding. The first relation given in Corollary 18 is the extension to handle second-moment matrices which do not have unitary condition number as well as corrupted covariates.

4.4 Learning Sigmoidal Neurons

We next study the problem of learning GLMs following the model given in § 5 of [SS19]. The error function for the the Kernelized GLM problem is given by the following equation for a single training pair $(\mathbf{x}_i, y_i) \sim \mathcal{D}$ in the kernelized case.

$$\mathcal{L}(f; \mathbf{x}_i, y_i) = (\sigma(f(\mathbf{x}_i)) - y_i)^2$$

Theorem 19. Let $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{d \times N}$ be the data matrix and $\mathbf{y} = [y_1, \dots, y_n]$ be the output, such that for $i \in \mathcal{P}$, \mathbf{x}_i are sampled from a sub-Gaussian distribution s.t. $\phi(\mathbf{x}_i) = X_i$ has sub-Gaussian norm L and second-moment matrix Σ , and $y_i = \omega(\mathbf{x}_i^\top \mathbf{w}^* + \xi_i)$ for ξ_i sampled from a centered sub-Gaussian distribution with sub-Gaussian norm ν . Suppose the link function $\omega : \mathbb{R} \rightarrow \mathbb{R}$, has bounded gradient s.t. $C_1 \leq \omega'(x) \leq C_2$ for absolute constants $C_1, C_2 > 0$ for all $x \in \mathbb{R}$. Then after $O\left(\kappa(\Sigma) \log\left(\frac{\|\mathbf{f}^*\|_{\mathcal{H}}}{\epsilon}\right)\right)$ gradient descent iterations, then with probability exceeding $1 - \delta$,

$$\|\mathbf{f}^{(T)} - \mathbf{f}^*\|_{\mathcal{H}} \leq \epsilon + O\left(C_1^{-1} C_2^2 \cdot \epsilon \sqrt{\lambda_{\max}(\Sigma) C_3 \log \epsilon^{-1}}\right) + O(\lambda_{\max}(\Sigma) \lambda_{\min}^{-2}(\Sigma) \cdot \sigma \cdot (C_1^{-2} C_2^2))$$

when $N \geq \frac{1}{\lambda_{\min}^2(\Sigma)} \cdot \left(8C_K \cdot d + \frac{2}{c_K} \cdot \log(2/\delta)\right)$.

Proof. The proof is deferred to § B.1.1. ■

Theorem 20. Let $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{d \times N}$ be the data matrix and $\mathbf{y} = [y_1, \dots, y_n]$ be the output, such that for $i \in \mathcal{P}$, \mathbf{x}_i are sampled from a sub-Gaussian distribution s.t. $\phi(\mathbf{x}_i) = X_i$ has sub-Gaussian norm L and second-moment matrix Σ , and $y_i = \omega(f^*(\mathbf{x}_i)) + \xi_i$ for ξ_i sampled from a centered sub-Gaussian distribution with sub-Gaussian norm ν . Suppose the link function $\omega : \mathbb{R} \rightarrow \mathbb{R}$, has bounded gradient s.t. $C_1 \leq \omega'(x) \leq C_2$ for absolute constants $C_1, C_2 > 0$ for all $x \in \mathbb{R}$. Then after $O\left(\kappa(\Sigma) \log\left(\frac{\|\mathbf{f}^*\|_{\mathcal{H}}}{\epsilon}\right)\right)$ gradient descent iterations, then with probability exceeding $1 - \delta$,

$$\|\mathbf{f}^{(T)} - \mathbf{f}^*\|_{\mathcal{H}} \leq \epsilon + O\left(C_1^{-1} C_2^2 \cdot \epsilon \sqrt{\lambda_{\max}(\Sigma) C_3 \log \epsilon^{-1}}\right) + O(\lambda_{\max}(\Sigma) \lambda_{\min}^{-2}(\Sigma) \cdot \sigma \cdot (C_1^{-2} C_2^2))$$

when $N \geq \frac{1}{\lambda_{\min}^2(\Sigma)} \cdot \left(8C_K \cdot d + \frac{2}{c_K} \cdot \log(2/\delta)\right)$.

Proof. The proof is deferred to § B.1.1. ■

When considerin the linear regression case

4.5 Learning Leaky-ReLU Neural Networks

We will now consider functions with Properties 13, 14, and 15. Functions that follow these properties include the Leaky-ReLU. In our proof we are able to leverage the fact that the second derivative is zero almost surely.

Theorem 21. Let $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{N \times d}$ be the data matrix and $\mathbf{y} = [y_1, \dots, y_n]^\top$ be the output, such that for $i \in \mathcal{P}$, $\mathbf{x}_i \sim \mathcal{P}$ are sampled from a sub-Gaussian distribution with sub-Gaussian norm K and second-moment matrix Σ , and $y_i = \psi(\mathbf{x}_i^\top \mathbf{w}^*) + \xi_i$ for $\xi_i \sim \mathcal{N}(0, \sigma^2)$ where $\psi(x) = \max\{C_\psi x, x\}$. Then after $O\left(C_\psi^{-2} \kappa(\Sigma) \log\left(\frac{\|\mathbf{w}^*\|_{\mathcal{F}}}{\epsilon}\right)\right)$ gradient descent iterations and $\epsilon \leq \frac{C_\psi^2 \lambda_{\min}(\Sigma)}{\sqrt{32C_Q \lambda_{\max}(\Sigma)}}$, with probability exceeding $1 - \delta$, Algorithm 1 with learning rate $\eta = O(\kappa^{-2}(\Sigma))$ returns $\mathbf{w}^{(T)}$ such that

$$\|\mathbf{w}^{(T)} - \mathbf{w}^*\|_2 \leq \epsilon + O\left(\kappa^2(\Sigma) C_\psi^{-1} C_K C_\sigma d \epsilon \log \epsilon^{-1}\right)$$

Proof. The proof is deferred to § B.2.1. ■

4.6 Learning ReLU Neural Networks

We will now consider the problem of learning ReLU neural networks. We first give a preliminary result for randomized initialization.

Lemma 22 (Theorem 3.4 in [DLT⁺18]). Suppose $\mathbf{w}^{(0)}$ is sampled uniformly from a p -dimensional ball with radius $\alpha \|\mathbf{w}^*\|$ such that $\alpha \leq \sqrt{\frac{1}{2\pi p}}$, then with probability at least $\frac{1}{2} - \alpha \sqrt{\frac{\pi p}{2}}$

$$\|\mathbf{w}^{(0)} - \mathbf{w}^*\|_2 \leq \sqrt{1 - \alpha^2} \|\mathbf{w}^*\|_2$$

From this result we are able to derive probabilistic guarantees on the convergence of learning ReLU neuron.

Theorem 23. Let $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$ be the data matrix and $\mathbf{y} = [y_1, \dots, y_n]^\top$ be the output, such that for $i \in P$, \mathbf{x}_i are sampled from a sub-Gaussian distribution with sub-Gaussian norm K and second-moment matrix Σ . Suppose for $i \in P$, the output is given as $y_i = \psi(\mathbf{x}_i^\top \mathbf{w}^*) + \xi_i$ for $\xi_i \sim \mathcal{N}(0, \sigma^2)$ and $\psi = \max\{0, x\}$ is the ReLU function. Then after $O(\Xi)$ gradient descent iterations and $n = \Omega(\Xi)$, then with probability exceeding $1 - \delta$, Algorithm 1 with learning rate $\eta = O(\Xi)$ returns $\mathbf{w}^{(T)}$ such that $\|\mathbf{w}^{(T)} - \mathbf{w}^*\|_2 \leq \varepsilon + O(\Xi)$.

Proof. The proof is deferred to § B.3.1 ■

5 Discussion

In this paper, we study the theoretical convergence properties of iterative thresholding for non-linear learning problems in the Strong ϵ -contamination model. Our warm-up result for linear regression reduces the runtime while achieving the best known approximation for iterative thresholding algorithms. Many papers have experimentally studied the iterative thresholding estimator in large scale neural networks [SS19, HYW⁺23] and to our knowledge, we are the first paper to make advancements in the theory of iterative thresholding for a general class of activation functions.

References

- [ADKS22] Pranjali Awasthi, Abhimanyu Das, Weihao Kong, and Rajat Sen. Trimmed maximum likelihood estimation for robust generalized linear model. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [AMMIL12] Yaser S Abu-Mostafa, Malik Magdon-Ismael, and Hsuan-Tien Lin. *Learning from data*, volume 4. AMLBook New York, 2012.
- [B⁺15] Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- [BJK15] Kush Bhatia, Prateek Jain, and Purushottam Kar. Robust regression via hard thresholding. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.
- [BJKK17] Kush Bhatia, Prateek Jain, Parameswaran Kamalaruban, and Purushottam Kar. Consistent robust regression. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [CDGS20] Yu Cheng, Ilias Diakonikolas, Rong Ge, and Mahdi Soltanolkotabi. High-dimensional robust mean estimation via gradient descent. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 1768–1778. PMLR, 13–18 Jul 2020.
- [CDGW19] Yu Cheng, Ilias Diakonikolas, Rong Ge, and David P. Woodruff. Faster algorithms for high-dimensional robust covariance estimation. In Alina Beygelzimer and Daniel Hsu, editors, *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*, pages 727–757. PMLR, 25–28 Jun 2019.
- [CLKZ21] Hugo Cui, Bruno Loureiro, Florent Krzakala, and Lenka Zdeborová. Generalization error rates in kernel regression: The crossover from the noiseless to noisy regime. *Advances in Neural Information Processing Systems*, 34:10131–10143, 2021.

- [CLRS22] Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, and Clifford Stein. *Introduction to algorithms*. MIT press, 2022.
- [Dic16] Lee H Dicker. Ridge regression and asymptotic minimax estimation over spheres of growing dimension. 2016.
- [DK23] Ilias Diakonikolas and Daniel M Kane. *Algorithmic high-dimensional robust statistics*. Cambridge University Press, 2023.
- [DKK⁺19] Ilias Diakonikolas, Gautam Kamath, Daniel M. Kane, Jerry Li, Jacob Steinhardt, and Alistair Stewart. Sever: A robust meta-algorithm for stochastic optimization. In *Proceedings of the 36th International Conference on Machine Learning*, ICML ’19, pages 1596–1606. JMLR, Inc., 2019.
- [DLT⁺18] Simon Du, Jason Lee, Yuandong Tian, Aarti Singh, and Barnabas Poczos. Gradient descent learns one-hidden-layer cnn: Don’t be afraid of spurious local minima. In *International Conference on Machine Learning*, pages 1339–1348. PMLR, 2018.
- [FB81] Martin A. Fischler and Robert C. Bolles. Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM*, 24(6):381–395, jun 1981.
- [FWZ18] Jianqing Fan, Weichen Wang, and Yiqiao Zhong. An l_1 eigenvector perturbation bound and its application to robust covariance estimation. *Journal of Machine Learning Research*, 18(207):1–42, 2018.
- [H89] O. Hölder. Ueber einen mittelwerthabsatz. *Nachrichten von der Königl. Gesellschaft der Wissenschaften und der Georg-Augusts-Universität zu Göttingen*, 1889:38–47, 1889.
- [HR09] Peter J. Huber and Elvezio. Ronchetti. *Robust statistics*. Wiley series in probability and statistics. Wiley, Hoboken, N.J., 2nd ed. edition, 2009.
- [HSS08] Thomas Hofmann, Bernhard Schölkopf, and Alexander J. Smola. Kernel methods in machine learning. *The Annals of Statistics*, 36(3):1171 – 1220, 2008.
- [HYW⁺23] Shu Hu, Zhenhuan Yang, Xin Wang, Yiming Ying, and Siwei Lyu. Outlier robust adversarial training. *arXiv preprint arXiv:2309.05145*, 2023.
- [HYwL20] Shu Hu, Yiming Ying, xin wang, and Siwei Lyu. Learning by minimizing the sum of ranked range. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 21013–21023. Curran Associates, Inc., 2020.
- [JLT20] Arun Jambulapati, Jerry Li, and Kevin Tian. Robust sub-gaussian principal component analysis and width-independent Schatten packing. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 15689–15701. Curran Associates, Inc., 2020.
- [JZL⁺18] Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In *ICML*, 2018.
- [KLA18] Ashish Khetan, Zachary C. Lipton, and Anima Anandkumar. Learning from noisy singly-labeled data. In *International Conference on Learning Representations*, 2018.
- [LBSS21] Tian Li, Ahmad Beirami, Maziar Sanjabi, and Virginia Smith. Tilted empirical risk minimization. In *International Conference on Learning Representations*, 2021.
- [Leg06] Adrien M Legendre. *Nouvelles methodes pour la determination des orbites des cometes: avec un supplement contenant divers perfectionnemens de ces methodes et leur application aux deux cometes de 1805*. Courcier, 1806.

- [LM00] Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of statistics*, pages 1302–1338, 2000.
- [MGJK19] Bhaskar Mukhoty, Govind Gopakumar, Prateek Jain, and Purushottam Kar. Globally-convergent iteratively reweighted least squares for robust regression problems. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pages 313–322. PMLR, 16–18 Apr 2019.
- [OZS20] Muhammad Osama, Dave Zachariah, and Petre Stoica. Robust risk minimization for statistical learning from corrupted data. *IEEE Open Journal of Signal Processing*, 1:287–294, 2020.
- [PBR19] Adarsh Prasad, Sivaraman Balakrishnan, and Pradeep Ravikumar. A unified approach to robust mean estimation. *arXiv preprint arXiv:1907.00927*, 2019.
- [PP⁺08] Kaare Brandt Petersen, Michael Syskind Pedersen, et al. The matrix cookbook. *Technical University of Denmark*, 7(15):510, 2008.
- [PSBR18] Adarsh Prasad, Arun Sai Suggala, Sivaraman Balakrishnan, and Pradeep Ravikumar. Robust estimation via robust gradient estimation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 82, 2018.
- [RH23] Philippe Rigollet and Jan-Christian Hütter. High-dimensional statistics. *arXiv preprint arXiv:2310.19244*, 2023.
- [RHL⁺20] Meisam Razaviyayn, Tianjian Huang, Songtao Lu, Maher Nouiehed, Maziar Sanjabi, and Mingyi Hong. Nonconvex min-max optimization: Applications, challenges, and recent theoretical advances. *IEEE Signal Processing Magazine*, 37(5):55–66, 2020.
- [SS19] Yanyao Shen and Sujay Sanghavi. Learning with bad training data via iterative trimmed loss minimization. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5739–5748. PMLR, 09–15 Jun 2019.
- [Ver10] Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- [Ver20] Roman Vershynin. High-dimensional probability. *University of California, Irvine*, 2020.
- [You12] William Henry Young. On classes of summable functions and their fourier series. *Proceedings of the Royal Society of London. Series A, Containing Papers of a Mathematical and Physical Character*, 87(594):225–229, 1912.
- [ZYWG19] Xiao Zhang, Yaodong Yu, Lingxiao Wang, and Quanquan Gu. Learning one-hidden-layer relu networks via gradient descent. In *The 22nd international conference on artificial intelligence and statistics*, pages 1524–1534. PMLR, 2019.

A Proofs for Linear Regression

Notation. We will first give some notational preliminaries. Let $X = P \cup Q$ for $|P| = (1 - \epsilon)N$ and $|Q| = \epsilon N$ represent the sets such that for $i \in P$, (\mathbf{x}_i, y_i) is the good data and for $j \in Q$, (\mathbf{x}_j, y_j) has been given by the adversary. For $t \in [T]$, we denote $S^{(t)}$ as the Subquantile set at iteration t and represents the points. We decompose $S^{(t)} = S^{(t)} \cap P \cup S^{(t)} \cap Q = TP \cup FP$ to represent the *True Positives* and *False Positives*. We also decompose $X \setminus S^{(t)} = (X \setminus S^{(t)}) \cap P \cup (X \setminus S^{(t)}) \cap Q = FN \cup TN$ to represent the *False Negatives* and the *True Negatives*.

A.1 Proof of Theorem 17

Proof. Recall that for any $W \in \mathbb{R}^{k \times d}$, $X \in \mathbb{R}^{d \times n}$, and $Y \in \mathbb{R}^{k \times n}$,

$$\begin{aligned} \mathcal{L}(W; X, Y) &= \|WX - Y\|_F^2 = \text{Tr}(X^\top W^\top WX - X^\top W^\top Y - Y^\top WX + Y^\top Y) \\ &= \text{Tr}(X^\top W^\top WX) + \text{Tr}(Y^\top Y) - 2 \text{Tr}(X^\top W^\top Y) \end{aligned}$$

Then, from [PP⁺08] Equations (102) and (119) (where we set $B = I$ and $C = \mathbf{0}$). We have,

$$\nabla \mathcal{L}(W) = 2(WX - Y)X^\top$$

We first show we can obtain Linear Convergence plus a noise term and a consistency term. We then show the consistency goes to zero with high probability as $N \rightarrow \infty$. Finally, we use the concentration inequalities developed in § C to give a clean error bound.

Step 1: Linear Convergence. We will first show iterative thresholding has linear convergence to optimal. Let $\widehat{W} = \arg \min_{W \in \mathbb{R}^{k \times d}} \mathcal{R}(W; P)$ be the minimizer over the good data. Then, we have

$$\begin{aligned} \|W^{(t+1)} - W^*\|_F &= \|W^{(t)} - W^* - \eta \nabla \mathcal{R}(W^{(t)}; S^{(t)})\|_F \\ &= \|W^{(t)} - W^* - \eta \nabla \mathcal{R}(W^{(t)}; TP) + \eta \nabla \mathcal{R}(W^*; P) - \eta \nabla \mathcal{R}(W^*; P) + \eta \nabla \mathcal{R}(\widehat{W}; P) - \eta \nabla \mathcal{R}(W^{(t)}; FP)\|_F \\ &\leq \|W^{(t)} - W^* + \eta \nabla \mathcal{R}(W^{(t)}; TP) + \eta \nabla \mathcal{R}(W^*; TP)\|_F + \|\eta \nabla \mathcal{R}(W^{(t)}; FP)\|_F + \|\eta \nabla \mathcal{R}(W^*; FN)\|_F \\ &\quad + \|\eta \nabla \mathcal{R}(W^*; P) - \eta \nabla \mathcal{R}(\widehat{W}; P)\|_F \end{aligned} \tag{A.1}$$

We first will upper bound the first term in Equation (A.1) through its square.

$$\begin{aligned} \|W^{(t)} - W^* - \eta \nabla \mathcal{R}(W^{(t)}; TP) + \eta \nabla \mathcal{R}(W^*; TP)\|_F^2 &= \|W^{(t)} - W^*\|_F^2 \\ &\quad - 2\eta \cdot \text{Tr}((W^{(t)} - W^*)^\top (\nabla \mathcal{R}(W^{(t)}; TP) - \nabla \mathcal{R}(W^*; TP))) + \|\eta \nabla \mathcal{R}(W^{(t)}; TP) - \eta \nabla \mathcal{R}(W^*; TP)\|_F^2 \end{aligned} \tag{A.2}$$

We then lower bound the second term in Equation (A.2),

$$\begin{aligned} &2\eta \cdot \text{Tr}((W^{(t)} - W^*)^\top (\nabla \mathcal{R}(W^{(t)}; TP) - \nabla \mathcal{R}(W^*; TP))) \\ &\stackrel{\text{def}}{=} \frac{4\eta}{(1 - \epsilon)N} \cdot \text{Tr}((W^{(t)} - W^*)^\top ((W^{(t)} X_{TP} - Y_{TP}) X_{TP}^\top - E_{TP} X_{TP}^\top)) \\ &\stackrel{(i)}{=} \frac{4\eta}{(1 - \epsilon)N} \cdot \text{Tr}((W^{(t)} - W^*)^\top (W^{(t)} - W^*) X_{TP} X_{TP}^\top) \end{aligned} \tag{A.3}$$

In the above, (i) follows from noting that $Y_{TP} = W^* X_{TP} - E_{TP}$. We can now lower bound the term in Equation (A.3).

$$\begin{aligned} &\frac{4\eta}{(1 - \epsilon)N} \cdot \text{Tr}((W^{(t)} - W^*)^\top (W^{(t)} - W^*) X_{TP} X_{TP}^\top) \\ &\stackrel{(ii)}{=} \frac{4\eta}{(1 - \epsilon)N} \cdot \text{Tr}((W^{(t)} - W^*) X_{TP} X_{TP}^\top (W^{(t)} - W^*)^\top) \\ &\stackrel{(iii)}{=} \frac{4\eta}{(1 - \epsilon)N} \cdot \langle \text{vec}((W^{(t)} - W^*)^\top), \text{vec}(X_{TP} X_{TP}^\top (W^{(t)} - W^*)^\top) \rangle \end{aligned}$$

$$\begin{aligned}
&\stackrel{(iv)}{=} \frac{4\eta}{(1-\epsilon)N} \cdot \langle \text{vec}((W^{(t)} - W^*)^\top), (I \otimes X_{\text{TP}} X_{\text{TP}}^\top) \text{vec}((W^{(t)} - W^*)^\top) \rangle \\
&\stackrel{(v)}{=} \frac{4\eta}{(1-\epsilon)N} \cdot \sum_{k \in [K]} \langle \mathbf{w}_k^{(t)} - \mathbf{w}_k^*, X_{\text{TP}} X_{\text{TP}}^\top (\mathbf{w}_k^{(t)} - \mathbf{w}_k^*) \rangle \\
&\stackrel{(vi)}{\geq} \frac{2\eta}{(1-\epsilon)N} \cdot \sum_{k \in [K]} (\lambda_{\min}(X_{\text{TP}} X_{\text{TP}}^\top) \|\mathbf{w}_k^{(t)} - \mathbf{w}_k^*\|_2^2 + \|X_{\text{TP}} X_{\text{TP}}^\top\|_2^{-1} \|X_{\text{TP}} X_{\text{TP}}^\top (\mathbf{w}_k^{(t)} - \mathbf{w}_k^*)\|_2^2) \\
&= \frac{2\eta}{(1-\epsilon)N} \cdot \lambda_{\min}(X_{\text{TP}} X_{\text{TP}}^\top) \|W^{(t)} - W^*\|_F^2 + \frac{2\eta}{(1-\epsilon)N} \cdot \|X_{\text{TP}} X_{\text{TP}}^\top\|_2^{-1} \|(W^{(t)} - W^*) X_{\text{TP}} X_{\text{TP}}^\top\|_F^2 \quad (\text{A.4})
\end{aligned}$$

In the above, (ii) follows from the cyclic property of the trace, (iii) follows from the relation given in Lemma 5, (iv) holds from the relation given in Lemma 6, the inequality in (vi) follows from Lemma 33, and in (v) we apply Lemma 6, which gives the following equality,

$$(I \otimes X_{\text{TP}} X_{\text{TP}}^\top) \text{vec}((W^{(t)} - W^*)^\top) = \begin{bmatrix} X_{\text{TP}} X_{\text{TP}}^\top & & \\ & \ddots & \\ & & X_{\text{TP}} X_{\text{TP}}^\top \end{bmatrix} \begin{bmatrix} \mathbf{w}_1 - \mathbf{w}_1^* \\ \vdots \\ \mathbf{w}_K - \mathbf{w}_K^* \end{bmatrix} = \begin{bmatrix} X_{\text{TP}} X_{\text{TP}}^\top (\mathbf{w}_1 - \mathbf{w}_1^*) \\ \vdots \\ X_{\text{TP}} X_{\text{TP}}^\top (\mathbf{w}_K - \mathbf{w}_K^*) \end{bmatrix}$$

We now upper bound the second term in Equation (A.1),

$$\begin{aligned}
\|\eta \nabla \mathcal{R}(W^{(t)}; \text{FP})\|_F &\stackrel{\text{def}}{=} \frac{2\eta}{(1-\epsilon)N} \cdot \|(W^{(t)} X_{\text{FP}} - Y_{\text{FP}}) X_{\text{FP}}^\top\|_F \\
&\stackrel{(vii)}{\leq} \frac{2\eta}{(1-\epsilon)N} \cdot \|X_{\text{FP}}\|_2 \|W^{(t)} X_{\text{FP}} - Y_{\text{FP}}\|_F \\
&\stackrel{(viii)}{\leq} \frac{2\eta}{(1-\epsilon)N} \cdot \|X_{\text{FP}}\|_2 \|W^{(t)} X_{\text{FN}} - Y_{\text{FN}}\|_F \\
&\stackrel{(ix)}{\leq} \frac{2\eta}{(1-\epsilon)N} \cdot \|X_{\text{FP}}\|_2 (\|W^{(t)} X_{\text{FN}} - W^* X_{\text{FN}}\|_F + \|E_{\text{FN}}\|_F) \\
&\stackrel{(x)}{\leq} \frac{2\eta}{(1-\epsilon)N} \cdot \|X_{\text{FP}}\|_2 \|X_{\text{FN}}\|_2 \|W^{(t)} - W^*\|_F + \frac{2\eta}{(1-\epsilon)N} \cdot \|X_{\text{FP}}\|_2 \|E_{\text{FN}}\|_F
\end{aligned}$$

In the above, the equalities in (vii) and (x) from the fact that for any two size compatible matrices, A, B , it holds that $\|AB\|_F \leq \|A\|_F \|B\|_2$, (viii) follows from the optimality of the Subquantile set, and (ix) follows from the sub-additivity of the Frobenius norm. We will now upper bound the third term in Equation (A.1),

$$\|\eta \nabla \mathcal{R}(W^*; \text{FN})\|_F \stackrel{\text{def}}{=} \frac{2\eta}{(1-\epsilon)N} \cdot \|(W^* X_{\text{FN}} - Y_{\text{FN}}) X_{\text{FN}}^\top\|_F \leq \frac{2\eta}{(1-\epsilon)N} \cdot \|E_{\text{FN}} X_{\text{FN}}^\top\|_F \quad (\text{A.5})$$

In the above, we use the fact that for any two size compatible matrices, A, B , it holds that $\|AB\|_F \leq \|A\|_F \|B\|_2$.

Step 2: Consistency. We will now upper bound the consistency estimate of the fourth term of Equation (A.1).

$$\|\eta \nabla \mathcal{R}(W^*; \text{P}) - \eta \nabla \mathcal{R}(\widehat{W}; \text{P})\|_F \stackrel{\text{def}}{=} \frac{2\eta}{(1-\epsilon)N} \cdot \|(\widehat{W} - W^*) X_{\text{P}} X_{\text{P}}^\top\|_F \leq \frac{2\eta}{(1-\epsilon)N} \cdot \|E_{\text{P}} X_{\text{P}}^\top\|_F$$

We can then have from Lemma 31,

$$\|E_{\text{P}} X_{\text{P}}^\top\|_F^2 \leq \frac{16}{3} \cdot \sigma(\|X_{\text{P}}\|_F^2 \log((2/\delta)N^2)) \leq \frac{32}{3} \cdot \sigma(N(1-\epsilon)d\lambda_{\max}(\Sigma) \log((2/\delta)N^2))$$

We then have with probability exceeding $1 - \delta$,

$$\frac{2\eta}{(1-\epsilon)N} \cdot \|E_{\text{P}} X_{\text{P}}^\top\|_F \leq \eta \cdot \sqrt{\frac{32\sigma^2 d\lambda_{\max}(\Sigma) \log((2/\delta)N^2)}{3(1-\epsilon)N}}$$

We then have from our choice of $\eta = 0.1\lambda_{\max}^{-2}(\Sigma)$, the third term in Equation (A.1) will be less than the second term in Equation (A.4). We thus obtain from noting that $\sqrt{1-2x} \leq 1-x$ for any $x \leq 1/2$,

$$\begin{aligned} \|W^{(t+1)} - W^*\|_F &\leq \|W^{(t)} - W^*\|_F \left(1 - \frac{2\eta}{(1-\epsilon)N} \cdot \lambda_{\min}(X_{\text{TP}}X_{\text{TP}}^\top) + \frac{2\eta}{(1-\epsilon)N} \cdot \|X_{\text{FP}}\|_2 \|X_{\text{FN}}\|_2\right) \\ &\quad + \frac{2\eta}{(1-\epsilon)N} \cdot \|X_{\text{FP}}\|_2 \|E_{\text{FN}}\|_F + \frac{2\eta}{(1-\epsilon)N} \cdot \|E_{\text{FN}}X_{\text{FN}}^\top\|_F + \eta \cdot \sqrt{\frac{32\sigma^2 d \lambda_{\max}(\Sigma) \log((2/\delta)N^2)}{3(1-\epsilon)N}} \end{aligned}$$

Step 3: Concentration Bounds. From Proposition 25, we obtain with probability exceeding $1 - \delta$ that $\|E_{\text{FN}}\|_F \leq \sigma \sqrt{30NK\epsilon \log \epsilon^{-1}}$. From our assumption on bounded covariate corruptions, we have $\|E_{\text{FN}}\|_F \|X_{\text{FP}}\|_F \leq \sigma \epsilon \sqrt{30NK C_3 \log \epsilon^{-1}}$. From Lemma 31, we have that $\|E_{\text{FN}}X_{\text{FN}}^\top\|_F \leq \sqrt{6K \log N} \|X_{\text{FN}}\|$ when $n \geq (1/\delta)$ with probability exceeding $1 - \delta$. From Lemma 27, we have for $\epsilon \leq \frac{1}{240}\kappa^{-1}(\Sigma)$, the minimum eigenvalue satisfies $\lambda_{\min}(X_{\text{TP}}X_{\text{TP}}^\top) \geq (N/2) \cdot (1-2\epsilon) \cdot \lambda_{\min}(\Sigma) \geq (N/4) \cdot \lambda_{\min}(\Sigma)$. From our assumption of corrupted covariates, we have that $\|X_{\text{FP}}\| \leq \sqrt{\epsilon N C_3}$. We also have from Lemma 27, we have $\|X_{\text{FN}}\| \leq \sqrt{\lambda_{\max}(\Sigma) \cdot (10N\epsilon \log \epsilon^{-1})}$ with high probability. Then when $\epsilon \leq \sqrt{\frac{1}{960C_3} \cdot \kappa^{-1}(\Sigma) \lambda_{\min}(\Sigma)}$, we have that $\|X_{\text{FP}}\|_2 \|X_{\text{FN}}\|_2 \leq (N/8) \lambda_{\min}(\Sigma)$. Then, solving for the induction with an infinite sum, we have after $O\left(\kappa(\Sigma) \cdot \log\left(\frac{\|W^*\|_F}{\epsilon}\right)\right)$ iterations,

$$\begin{aligned} \|W^{(t+1)} - W^*\|_F &\leq \epsilon + \sigma \epsilon \sqrt{4800K C_3 \log \epsilon^{-1}} \cdot \frac{\sqrt{\lambda_{\max}(\Sigma)}}{\lambda_{\min}(\Sigma)} + \frac{\sqrt{60\sigma K \log N \cdot \epsilon \log \epsilon^{-1} \lambda_{\max}(\Sigma)}}{\sqrt{N} \lambda_{\min}(\Sigma)} \\ &\quad + \frac{\sigma}{\lambda_{\min}(\Sigma)} \cdot \sqrt{\frac{(8/\delta)K \text{Tr}(\Sigma)}{N}} \leq \epsilon + \sigma \epsilon \sqrt{43200K C_3 \log \epsilon^{-1}} \cdot \frac{\sqrt{\lambda_{\max}(\Sigma)}}{\lambda_{\min}(\Sigma)} \end{aligned}$$

In the above, the final equality follows when $N \geq \frac{(1/\delta)K \text{Tr}(\Sigma)}{1600\epsilon^2 \log \epsilon^{-1} \lambda_{\max}(\Sigma)} \vee \frac{1}{6400C_3^2 \epsilon^2}$. Our proof is complete. ■

B Proofs for Learning Single Neurons

B.1 Sigmoidal Neurons

B.1.1 Proof of Theorem 19

Proof. From Algorithm 1, we have the gradient update for learning a sigmoidal neuron.

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \frac{2\eta}{(1-\epsilon)N} \cdot \sum_{i \in S^{(t)}} (\sigma(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - y_i) \cdot \sigma'(f^{(t)}(\mathbf{x}_i)) \cdot \mathbf{x}_i$$

Our proof will follow a similar structure to the proof for linear regression. We first show we can obtain linear convergence of $\mathbf{w}^{(t)}$ to \mathbf{w}^* with some error. Then we rigorously analyze the concentration inequalities to give crisp bounds on the upper bound for ϵ and show the noise term is $O(\epsilon \log \epsilon^{-1})$.

Step 1: Linear Convergence.

$$\begin{aligned} \|\mathbf{w}^{(t+1)} - \mathbf{w}^*\|_2 &= \|\mathbf{w}^{(t)} - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; S^{(t)}) - \mathbf{w}^*\|_2 \\ &= \|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{FP})\|_2 \\ &\leq \|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP})\|_2 + \|\eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{FP})\|_2 \end{aligned} \quad (\text{B.1})$$

We upper bound the first term of Equation (B.1) through its square.

$$\|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP})\|_2^2 = \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2 - 2\eta \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) \rangle + \eta^2 \cdot \|\nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP})\|_2^2 \quad (\text{B.2})$$

We now lower bound the second term of Equation (B.2).

$$2\eta \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) \rangle = \frac{4\eta}{(1-\epsilon)N} \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \sum_{i \in \text{TP}} (\sigma(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - y_i) \cdot \sigma'(\mathbf{x}_i^\top \mathbf{w}^{(t)}) \cdot \mathbf{x}_i \rangle$$

$$\begin{aligned}
&= \frac{4\eta}{(1-\epsilon)N} \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \sum_{i \in \text{TP}} (\sigma(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - \sigma(\mathbf{x}_i^\top \mathbf{w}^*) + \xi_i) \cdot \sigma'(\mathbf{x}_i^\top \mathbf{w}^{(t)}) \cdot \mathbf{x}_i \rangle \\
&\geq \frac{4\eta}{(1-\epsilon)N} \cdot C_{[\sigma]}^2 \lambda_{\min}(X_{\text{TP}} X_{\text{TP}}^\top) \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2 - \frac{4\eta}{(1-\epsilon)N} \cdot \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2 \left\| \sum_{i \in \text{TP}} \xi_i \sigma'(\mathbf{x}_i^\top \mathbf{w}^{(t)}) \cdot \mathbf{x}_i \right\|_2 \quad (\text{B.3})
\end{aligned}$$

In the above, in the final relation, we apply Cauchy-Schwarz inequality. Then from an application of Young's Inequality (see Proposition 34), we obtain

$$\begin{aligned}
&\frac{4\eta}{(1-\epsilon)N} \cdot \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2 \left\| \sum_{i \in \text{TP}} \xi_i \sigma'(\mathbf{x}_i^\top \mathbf{w}^{(t)}) \cdot \mathbf{x}_i \right\|_2 \\
&\leq \frac{\eta}{(1-\epsilon)N} \cdot C_{[\sigma]}^2 \lambda_{\min}(X_{\text{TP}} X_{\text{TP}}^\top) \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2 + \frac{4\eta}{(1-\epsilon)N} \cdot \lambda_{\min}^{-1}(X_{\text{TP}} X_{\text{TP}}^\top) \left\| \sum_{i \in \text{TP}} \xi_i \sigma'(\mathbf{x}_i^\top \mathbf{w}^{(t)}) \cdot \mathbf{x}_i \right\|_2^2
\end{aligned}$$

We next upper bound the third term in Equation (B.2).

$$\begin{aligned}
\eta^2 \cdot \|\nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP})\|_2^2 &= \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot \left\| \sum_{i \in \text{TP}} (\sigma(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - \sigma(\mathbf{x}_i^\top \mathbf{w}^*) + \xi_i) \cdot \sigma'(\mathbf{x}_i^\top \mathbf{w}^{(t)}) \cdot \mathbf{x}_i \right\|_2^2 \\
&= \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot \left\| \sum_{i \in \text{TP}} (\mathbf{x}_i^\top \mathbf{w}^{(t)} - \mathbf{x}_i^\top \mathbf{w}^* + \xi_i) \cdot \sigma'(c_i) \cdot \sigma'(\mathbf{x}_i^\top \mathbf{w}^{(t)}) \cdot \mathbf{x}_i \right\|_2^2 \\
&\leq \frac{8\eta^2}{[(1-\epsilon)N]^2} \cdot \left(C_{[\sigma]}^4 \lambda_{\max}^2(X_{\text{TP}} X_{\text{TP}}^\top) \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2 + \left\| \sum_{i \in \text{TP}} \xi_i \cdot \sigma'(\mathbf{x}_i^\top \mathbf{w}^{(t)}) \cdot \mathbf{x}_i \right\|_2^2 \right) \quad (\text{B.4})
\end{aligned}$$

Then, from choosing $\eta \leq \frac{C_{[\sigma]}^2(1-\epsilon)N\lambda_{\min}(X_{\text{TP}} X_{\text{TP}}^\top)}{4C_{[\sigma]}^4\lambda_{\max}^2(X_{\text{TP}} X_{\text{TP}}^\top)}$. We see the first term of Equation (B.4) is less than half the first term in Equation (B.3). We now bound the second term in Equation (B.1).

$$\begin{aligned}
\eta^2 \cdot \|\nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{FP})\|_2^2 &= \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot \left\| \sum_{i \in \text{FP}} (\sigma(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - y_i) \cdot \sigma'(\mathbf{x}_i^\top \mathbf{w}) \cdot \mathbf{x}_i \right\|_2^2 \\
&\leq \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot C_{[\sigma]}^2 \|X_{\text{FP}} X_{\text{FP}}^\top\|_2 \sum_{i \in \text{FP}} (\sigma(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - y_i)^2 \\
&\leq \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot C_{[\sigma]}^2 \|X_{\text{FP}} X_{\text{FP}}^\top\|_2 \sum_{i \in \text{FN}} (\sigma(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - y_i)^2 \\
&= \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot C_{[\sigma]}^2 \|X_{\text{FP}} X_{\text{FP}}^\top\|_2 \sum_{i \in \text{FN}} (\sigma(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - \sigma(\mathbf{x}_i^\top \mathbf{w}^*) + \xi_i)^2 \\
&\leq \frac{8\eta^2}{[(1-\epsilon)N]^2} \cdot C_{[\sigma]}^2 \|X_{\text{FP}} X_{\text{FP}}^\top\|_2 \left(C_{[\sigma]}^2 \cdot \|X_{\text{FN}} X_{\text{FN}}^\top\|_2 \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2 + \|\boldsymbol{\xi}_{\text{FN}}\|_2^2 \right)
\end{aligned}$$

Concluding the step, we have

$$\begin{aligned}
\|\mathbf{w}^{(t+1)} - \mathbf{w}^*\|_2 &\leq \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2 \left(1 - \frac{\eta}{2(1-\epsilon)N} \cdot C_{[\sigma]}^2 \lambda_{\min}(X_{\text{TP}} X_{\text{TP}}^\top) + \frac{\sqrt{8}\eta}{(1-\epsilon)N} \cdot C_{[\sigma]}^2 \|X_{\text{FP}}\|_2 \|X_{\text{FN}}\|_2 \right) \\
&\quad + \left(\frac{\sqrt{8}\eta}{(1-\epsilon)N} + \frac{2\sqrt{\eta}}{\sqrt{(1-\epsilon)N}} \cdot \sigma_{\min}^{-1}(X_{\text{TP}}) \right) \left\| \sum_{i \in \text{TP}} \xi_i \sigma'(c_i) \sigma'(\mathbf{x}_i^\top \mathbf{w}^{(t)}) \mathbf{x}_i \right\|_2 + \frac{\sqrt{8}\eta}{(1-\epsilon)N} \cdot \|X_{\text{FP}}\|_2 \|\boldsymbol{\xi}_{\text{FN}}\|_2
\end{aligned}$$

Step 2: Concentration Bounds. First we will show that $\xi \sigma'(\mathbf{x}^\top \mathbf{w})$ is sub-Gaussian for any $c \in \mathbb{R}$, $\mathbf{w} \in \mathbb{R}^d$ and \mathbf{x} sampled from a sub-Gaussian distribution. Note that ξ is sub-Gaussian s.t. $\|\xi\|_{\psi_2} = C_\sigma$. Then, we have

$$\left(\mathbf{E} |\sigma'(\mathbf{x}^\top \mathbf{w})|^p \right)^{1/p} \leq \left(\sup_{x \in \mathbb{R}} [\sigma'(x)]^{2p} \right)^{1/p} \leq C_{[\psi]}^2$$

We thus have $\|\sigma'(\mathbf{x}^\top \mathbf{w})\|_{\psi_2} \leq C_{[\psi]}^2 C_\sigma$. Then applying Lemma 28, we have with probability at least $1 - \delta$,

$$\left\| \sum_{i \in \text{TP}} \xi_i \cdot \sigma'(c_i) \cdot \sigma'(\mathbf{x}_i^\top \mathbf{w}^{(t)}) \cdot \mathbf{x}_i \right\|_2 \leq NC_\Sigma C_{[\psi]}^2 C_\sigma \left(\frac{2d}{N} \log(5) + \frac{2}{N} \log(1/\delta) + 6\epsilon \log \epsilon^{-1} \right)^{1/2}$$

Recall that $\|X_{\text{FP}}\|_2 \leq \sqrt{N\epsilon C_Q}$. Then, with probability at least $1 - \delta$, we have $\|X_{\text{FN}}\|_2 \|X_{\text{FP}}\|_2 \leq N\epsilon \cdot \sqrt{10C_3 \log \epsilon^{-1}}$. We also have $\|\xi_{\text{FN}}\|_2 \leq \sqrt{30N\epsilon \log \epsilon^{-1}}$. We then find for $\epsilon \leq \frac{C_{[\psi]}^2}{C_{[\psi]}^2} \frac{\lambda_{\min}(\Sigma)}{32\sqrt{80 \log(1/2)}}$ and probability exceeding $1 - \delta$,

$$C_{[\psi]}^2 \|X_{\text{FP}}\|_2 \|X_{\text{FN}}\|_2 \leq \frac{1}{\sqrt{144}} \cdot C_{[\psi]}^2 \lambda_{\min}(\Sigma)$$

Then combining the results with our choice of η , we obtain,

$$\begin{aligned} \|\mathbf{w}^{(t+1)} - \mathbf{w}^*\|_2 &\leq \varepsilon + \frac{1}{\lambda_{\min}(\Sigma)} \cdot C_{[\psi]}^{-2} \left(8 + 16\sqrt{8} \kappa^2(\Sigma) \right) \left(C_\Sigma C_{[\psi]}^2 C_\sigma \left(\frac{2d}{N} \log(5) + \frac{2}{N} \log(1/\delta) + 6\epsilon \log \epsilon^{-1} \right)^{1/2} \right) \\ &\leq \varepsilon + \frac{1}{\lambda_{\min}(\Sigma)} \cdot C_{[\psi]}^{-2} C_{[\psi]}^2 \sigma \left(24\sqrt{96} \kappa^2(\Sigma) \right) \sqrt{\epsilon \log \epsilon^{-1}} \end{aligned}$$

In the above, the final inequality holds when $N \geq \frac{2d \log 5 + \log(1/\delta)}{6\epsilon \log \epsilon^{-1}}$. Our proof is complete. \blacksquare

B.2 Piecewise Linear Unit Neurons

B.2.1 Proof of Theorem 21

Proof. We will decompose the gradient into the good component and corrupted component. The first part of our proof will show that \mathbf{w} moves in the direction of \mathbf{w}^* , then in the second part of the proof we will show the affect of the corrupted gradient. Finally, we combine step 1 and step 2 to show that there exists sufficiently small ϵ such that we can get linear convergence with a small additive error term.

Step 1: Upper bounding the ℓ_2 norm distance between $\mathbf{w}^{(t+1)}$ and \mathbf{w}^* . We have from Algorithm 1,

$$\begin{aligned} \|\mathbf{w}^{(t+1)} - \mathbf{w}^*\|_2 &= \|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \mathbf{S}^{(t)})\|_2 \\ &= \|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) + \eta \nabla \mathcal{R}(\mathbf{w}^*; \text{P}) - \eta \nabla \mathcal{R}(\mathbf{w}^*; \text{P}) - \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{FP})\|_2 \\ &\leq \|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) + \eta \nabla \mathcal{R}(\mathbf{w}^*; \text{TP})\|_2 + \|\eta \nabla \mathcal{R}(\mathbf{w}^*; \text{TP})\|_2 + \|\eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{FP})\|_2 \end{aligned} \quad (\text{B.5})$$

We will first upper bound the first term of Equation (B.5) through its square.

$$\begin{aligned} \|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP})\|_2^2 &= \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2 - 2\eta \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) \rangle + \eta^2 \cdot \|\nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP})\|_2^2 \end{aligned} \quad (\text{B.6})$$

In the above, the a.s. relation follows from Property 15. We will first lower bound the second term of Equation (B.6). We first will quickly give spectral bounds on $\nabla^2 \mathcal{L}(\mathbf{w}; \text{TP})$ where $\mathbf{w} \in \mathbb{R}^d$ is an arbitrary vector.

$$\begin{aligned} \nabla^2 \mathcal{L}(\mathbf{w}; \text{TP}) &= 2 \cdot \sum_{i \in \text{TP}} (\psi(\mathbf{x}_i^\top \mathbf{w}) - y_i) \cdot \psi''(\mathbf{x}_i^\top \mathbf{w}) \cdot \mathbf{x}_i \mathbf{x}_i^\top + 2 \cdot \sum_{i \in \text{TP}} [\psi'(\mathbf{x}_i^\top \mathbf{w})]^2 \cdot \mathbf{x}_i \mathbf{x}_i^\top \\ &\stackrel{\text{a.s.}}{=} 2 \cdot \sum_{i \in \text{TP}} [\psi'(\mathbf{x}_i^\top \mathbf{w})]^2 \cdot \mathbf{x}_i \mathbf{x}_i^\top \end{aligned}$$

We then obtain for any $\mathbf{w} \in \mathbb{R}^d$,

$$2 \cdot C_{[\psi]}^2 \lambda_{\min}(X_{\text{TP}} X_{\text{TP}}^\top) \cdot I \preceq \nabla^2 \mathcal{L}(\mathbf{w}; \text{TP}) \preceq 2 \cdot C_{[\psi]}^2 \lambda_{\max}(X_{\text{TP}} X_{\text{TP}}^\top) \cdot I$$

We can now lower bound the second term of Equation (B.6).

$$\begin{aligned}
& 2\eta \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) - \nabla \mathcal{R}(\mathbf{w}^*; \text{TP}) \rangle \\
&= \frac{4\eta}{(1-\epsilon)N} \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \int_0^1 \nabla^2 \mathcal{R}(\mathbf{w}^* + \theta(\mathbf{w}^* - \mathbf{w}^{(t)}); \text{TP}) d\theta \cdot (\mathbf{w}^{(t)} - \mathbf{w}^*) \rangle \\
&\stackrel{(i)}{\geq} \frac{4\eta}{(1-\epsilon)N} \cdot C_{\lceil \psi \rceil}^2 \lambda_{\min}(X_{\text{TP}} X_{\text{TP}}^\top) \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2
\end{aligned} \tag{B.7}$$

where (i) follows from Property 13. We now will upper bound the second term of Equation (B.6).

$$\begin{aligned}
& \eta^2 \cdot \|\nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) - \nabla \mathcal{R}(\mathbf{w}^*; \text{TP})\|_2^2 \\
&= \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot \left\| \int_0^1 \nabla^2 \mathcal{R}(\mathbf{w}^* + \theta(\mathbf{w}^* - \mathbf{w}^{(t)}); \text{TP}) d\theta \cdot (\mathbf{w}^{(t)} - \mathbf{w}^*) \right\|_2^2 \\
&\stackrel{(ii)}{\leq} \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot C_{\lceil \psi \rceil}^2 \lambda_{\max}^2(X_{\text{TP}} X_{\text{TP}}^\top) \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2
\end{aligned} \tag{B.8}$$

where (ii) follows from the $C_{\lceil \psi \rceil}$ -Lipschitzness of ψ given in Property 14. We then observe that the first term of Equation (B.7) is greater than half the first term in Equation (B.8) when $\eta \leq \frac{C_{\lceil \psi \rceil}^2 (1-\epsilon) N \lambda_{\min}(X_{\text{TP}} X_{\text{TP}}^\top)}{2\lambda_{\max}^2(X_{\text{TP}} X_{\text{TP}}^\top)}$.

Step 2: Upper bounding the corrupted gradient. We now upper bound the third term in Equation (B.6).

$$\begin{aligned}
\eta^2 \cdot \|\nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{FP})\|_{\text{F}}^2 &= \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot \left\| \sum_{i \in \text{FP}} (\psi(\mathbf{x}_i^\top \mathbf{w}) - \psi(\mathbf{x}_i^\top \mathbf{w}^*)) \cdot \psi'(\mathbf{x}_i^\top \mathbf{w}_k^{(t)}) \cdot \mathbf{x}_i \right\|_2^2 \\
&\leq \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot C_{\lceil \psi \rceil}^2 \|X_{\text{FP}} X_{\text{FP}}^\top\|_2 \sum_{i \in \text{FP}} (\psi(\mathbf{x}_i^\top \mathbf{w}) - \psi(\mathbf{x}_i^\top \mathbf{w}^*))^2 \\
&\leq \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot C_{\lceil \psi \rceil}^2 \|X_{\text{FP}} X_{\text{FP}}^\top\|_2 \sum_{i \in \text{FN}} (\psi(\mathbf{x}_i^\top \mathbf{w}) - \psi(\mathbf{x}_i^\top \mathbf{w}^*))^2 \\
&\leq \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot C_{\lceil \psi \rceil}^4 \|X_{\text{FP}} X_{\text{FP}}^\top\|_2 \|X_{\text{FN}} X_{\text{FN}}^\top\|_2 \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2
\end{aligned}$$

In the above, the first inequality follows from Lemma 32, the second inequality follows from the optimality of the Subquantile set, the final inequality follows from the C_2 -Lipschitzness of ψ . We will now show the random variable, $\psi'(\mathbf{x}_i^\top \mathbf{w}^{(t)}) \cdot \mathbf{x}_i$ is sub-Gaussian. Let $\mathbf{v} \in \mathbb{S}^{d-1}$, then since \mathbf{x}_i is sampled from sub-Gaussian distribution with sub-Gaussian norm C_Σ , we then have $\mathbf{v}^\top \mathbf{x}_i$ is sub-Gaussian with norm K . Therefore,

$$(\mathbf{E} |\psi'(\mathbf{x}_i^\top \mathbf{w}) \mathbf{v}^\top \mathbf{x}_i|^p)^{1/p} \leq \sup_{x \in \mathbb{R}} |\psi'(x)| (\mathbf{E} |\mathbf{v}^\top \mathbf{x}_i|^{2p})^{1/2p} \leq C_{\lceil \psi \rceil} C_\Sigma \sqrt{2p}$$

In the above, in the first inequality we used Hölder's Inequality. We therefore have, $\|\psi'(\mathbf{x}_i^\top \mathbf{w}^{(t)}) \cdot \mathbf{x}_i\|_{\psi_2} \leq \sqrt{2} C_{\lceil \psi \rceil} C_\Sigma$. Then from Lemma 28, we have with probability at least $1 - \delta$,

$$\left\| \sum_{i \in \text{TP}} \xi_i \cdot \psi'(\mathbf{x}_i^\top \mathbf{w}^{(t)}) \cdot \mathbf{x}_i \right\|_2 \leq 9NC_\Sigma C_\sigma \left(\frac{2d}{N} \log 12 + \frac{2Rd}{N} \log 12 + \frac{2}{N} \log(1/\delta) + 6\epsilon \log \epsilon^{-1} \right)^{1/2}$$

We now combine Steps 1 and 2 to give the linear convergence result. Noting that $\sqrt{1-2x} \leq 1-x$ when $x \leq 1/2$, we have

$$\begin{aligned}
\|\mathbf{w}^{(t+1)} - \mathbf{w}^*\|_2 &\leq \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2 \left(1 - \frac{2C_\psi^2 \eta}{(1-\epsilon)N} \cdot \lambda_{\min}(X_{\text{TP}} X_{\text{TP}}^\top) + \frac{2\eta}{(1-\epsilon)N} \cdot \|X_{\text{FP}}\|_2 \|X_{\text{FN}}\|_2 \right) \\
&\quad + \frac{15C_\Sigma C_\sigma d \log 5}{\sqrt{(1-\epsilon)N}} + \frac{15C_\Sigma C_\sigma \log(1/\delta)}{\sqrt{(1-\epsilon)N}} + 16C_\Sigma C_\sigma \sqrt{\epsilon \log \epsilon^{-1}}
\end{aligned}$$

Step 3: Concentration Bounds. We will give the relevant probabilistic bounds for the random variables in Steps 1 and 2. From Lemma 27, we have $\|X_{\text{FN}}\|_2 \|X_{\text{FP}}\|_2 \leq \epsilon \sqrt{\lambda_{\max}(\Sigma) \cdot 10C_3N \log \epsilon^{-1}}$ with probability at least $1 - \delta$ when $N \geq \frac{2}{\epsilon} \cdot \left(dC_\Sigma^2 + \frac{\log(2/\delta)}{c_K} \right)$ when $\epsilon \leq \frac{1}{60} \cdot \kappa^{-1}(\Sigma)$. From the same Lemma and under the same data conditions we have $\lambda_{\min}(X_{\text{TP}}X_{\text{TP}}^\top) \geq \frac{1}{4} \cdot \lambda_{\min}(\Sigma)$. Then for $\epsilon \leq \frac{C_\psi^2 \lambda_{\min}(\Sigma)}{\sqrt{32C_3 \lambda_{\max}(\Sigma)}}$, we have $\|X_{\text{FP}}\|_2 \|X_{\text{FN}}\|_2 \leq \frac{1}{2} \cdot \lambda_{\min}(\Sigma)$. We then have, after $O\left(\kappa^2(\Sigma) \log\left(\frac{\|\mathbf{w}^*\|_2 + 10d}{\epsilon}\right)\right)$ iterations with high probability,

$$\begin{aligned} \|\mathbf{w}^{(T)} - \mathbf{w}^*\|_2 &\leq \varepsilon + \frac{1}{N} \cdot 144C_\psi^{-2}\kappa^2(\Sigma)KC_\sigma d \log 5 + 432\kappa^2(\Sigma)C_\psi^{-1}KC_\sigma d \epsilon \log \epsilon^{-1} \\ &= O\left(\kappa^2(\Sigma)C_\psi^{-2}KC_\sigma d \epsilon \log \epsilon^{-1}\right) \end{aligned}$$

In the final inequality above, we set $\varepsilon = O\left(\kappa^2(\Sigma)C_\psi^{-1}KC_\sigma d \epsilon \log \epsilon^{-1}\right)$ for $N \geq \epsilon^{-2}C_\psi^{-2}\kappa^2(\Sigma) \log 5$. Our proof is complete. \blacksquare

B.3 Learning a ReLU Neuron

In this section, we consider functions such as the ReLU. Our high-level analysis will be similar to the previous sub-sections however the details are considerably different and require stronger conditions we can guarantee by randomness.

B.3.1 Proof of Theorem 23

Proof. We will now begin our standard analysis.

$$\begin{aligned} \|\mathbf{w}^{(t+1)} - \mathbf{w}^*\|_2 &= \|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \mathbf{S}^{(t)})\|_2 \\ &= \|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{FP})\|_2 \\ &\leq \|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP})\|_2 + \|\eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{FP})\|_2 \end{aligned} \quad (\text{B.9})$$

We will now upper bound the first term of Equation (B.9) through its square,

$$\begin{aligned} \|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP})\|_2^2 &= \\ &= \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2 - 2\eta \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) \rangle + \eta^2 \cdot \|\nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP})\|_2^2 \end{aligned} \quad (\text{B.10})$$

We will lower bound the second term of Equation (B.10). We will first adopt the notation from [ZYWG19], let $\Sigma_{\text{TP}}(\mathbf{w}, \hat{\mathbf{w}}) = X_{\text{TP}}^\top X_{\text{TP}}^\top \cdot \mathbb{I}\{X_{\text{TP}}^\top \mathbf{w} \geq \mathbf{0}\} \cdot \mathbb{I}\{X_{\text{TP}}^\top \hat{\mathbf{w}} \geq \mathbf{0}\}$, it then follows

$$\begin{aligned} &2\eta \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) \rangle \\ &\stackrel{\text{def}}{=} \frac{4\eta}{(1-\epsilon)N} \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \sum_{i \in \text{TP}} (\psi(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - y_i) \cdot \mathbf{x}_i \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^{(t)} \geq 0\} \rangle \\ &= \frac{4\eta}{(1-\epsilon)N} \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \Sigma_{\text{TP}}(\mathbf{w}^{(t)}, \mathbf{w}^{(t)})\mathbf{w}^{(t)} - \Sigma_{\text{TP}}(\mathbf{w}^{(t)}, \mathbf{w}^*)\mathbf{w}^* \rangle \\ &= \frac{4\eta}{(1-\epsilon)N} \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \Sigma_{\text{TP}}(\mathbf{w}^{(t)}, \mathbf{w}^*)(\mathbf{w}^{(t)} - \mathbf{w}^*) + \Sigma_{\text{TP}}(\mathbf{w}^{(t)}, -\mathbf{w}^*)\mathbf{w}^{(t)} \rangle \\ &\geq \frac{4\eta}{(1-\epsilon)N} \cdot \lambda_{\min}(\Sigma_{\text{TP}}(\mathbf{w}^{(t)}, \mathbf{w}^*)) \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2 \end{aligned}$$

In the above, the final inequality holds from the following relation,

$$\langle \mathbf{w}^{(t)} - \mathbf{w}^*, \Sigma_{\text{TP}}(\mathbf{w}^{(t)}, -\mathbf{w}^*)\mathbf{w}^{(t)} \rangle = \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \sum_{i \in \text{TP}} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{w}^{(t)} \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^{(t)} \geq 0\} \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^* \leq 0\} \rangle$$

$$= \sum_{i \in \text{TP}} (\mathbf{x}_i^\top \mathbf{w}^{(t)} - \mathbf{x}_i^\top \mathbf{w}^*) (\mathbf{x}_i^\top \mathbf{w}^{(t)}) \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^{(t)} \geq 0\} \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^* \leq 0\} \geq 0$$

In the above, in the final relation we can note that when the indicators are positive, it must follow that both $\mathbf{x}_i^\top \mathbf{w}^{(t)}$ is positive and $\mathbf{x}_i^\top \mathbf{w}^{(t)} \geq \mathbf{x}_i^\top \mathbf{w}^*$ as $\mathbf{x}_i^\top \mathbf{w}^* \leq 0$. We have from Weyl's Inequality,

$$\lambda_{\min}(\Sigma_{\text{TP}}(\mathbf{w}^{(t)}, \mathbf{w}^*)) \geq \lambda_{\min}(\mathbf{E}[\Sigma_{\text{TP}}(\mathbf{w}^{(t)}, \mathbf{w}^*)]) - \|\Sigma_{\text{TP}}(\mathbf{w}^{(t)}, \mathbf{w}^*) - \mathbf{E}[\Sigma_{\text{TP}}(\mathbf{w}^{(t)}, \mathbf{w}^*)]\|_2$$

Let $\Omega = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{x}^\top \mathbf{w}^{(t)} \geq 0, \mathbf{x}^\top \mathbf{w}^* \geq 0\}$, then

$$\begin{aligned} \mathbf{E}[\Sigma_{\text{TP}}(\mathbf{w}^{(t)}, \mathbf{w}^*)] &= \sum_{i \in \text{TP}} \mathbf{E}[\mathbf{x}_i \mathbf{x}_i^\top \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^{(t)} \geq 0\} \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^* \geq 0\}] \\ &\stackrel{(i)}{\succeq} N(1 - 2\epsilon) \cdot (\pi - \Theta^{(t)} - \sin \Theta^{(t)}) \cdot I \\ &\stackrel{(ii)}{\succeq} N(1 - 2\epsilon) \cdot \left(\pi - 2 \arcsin\left(\frac{\|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2}{\|\mathbf{w}^*\|_2}\right) \right) \cdot I \\ &\stackrel{(iii)}{\succeq} N(1 - 2\epsilon) \cdot \pi \cdot \left(1 - \frac{\|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2}{\|\mathbf{w}^*\|_2} \right) \cdot I \end{aligned}$$

In the above, (i) follows from Lemma 29, the inequality in (ii) follows when $\|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2 \leq \|\mathbf{w}^*\|_2$, and (iii) follows from noting $\arcsin(x) \leq \frac{\pi}{2} \cdot x$ for $x \in [0, 1]$. We now bound the second-moment matrix approximation. Let $dP(\mathbf{x})$ be a Dirac-measure for $\mathbf{x} \in \text{TP}$. We will now upper bound the third term in Equation (B.10).

$$\begin{aligned} \eta^2 \cdot \|\nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP})\|_2^2 &= \frac{4\eta^2}{[(1 - \epsilon)N]^2} \cdot \left\| \sum_{i \in \text{TP}} (\psi(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - \psi(\mathbf{x}_i^\top \mathbf{w}^*) - \xi_i) \cdot \mathbf{x}_i \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^{(t)} \geq 0\} \right\|_2^2 \\ &\leq \frac{8\eta^2}{[(1 - \epsilon)N]^2} \cdot \left(\left\| \sum_{i \in \text{TP}} (\psi(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - \psi(\mathbf{x}_i^\top \mathbf{w}^*)) \cdot \mathbf{x}_i \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^{(t)} \geq 0\} \right\|_2^2 + \left\| \sum_{i \in \text{TP}} \xi_i \mathbf{x}_i \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^{(t)} \geq 0\} \right\|_2^2 \right) \end{aligned} \quad (\text{B.11})$$

Recall we have from Lemma 28, we have an upper bound on the second term of Equation (B.11). We give a bound on the first term of Equation (B.11).

$$\begin{aligned} &\frac{8\eta^2}{[(1 - \epsilon)N]^2} \cdot \left\| \sum_{i \in \text{TP}} (\psi(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - \psi(\mathbf{x}_i^\top \mathbf{w}^*)) \cdot \mathbf{x}_i \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^{(t)} \geq 0\} \right\|_2^2 \\ &\leq \frac{8\eta^2}{[(1 - \epsilon)N]^2} \cdot \|X_{\text{TP}} X_{\text{TP}}^\top\|_2 \sum_{i \in \text{TP}} (\psi(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - \psi(\mathbf{x}_i^\top \mathbf{w}^*))^2 \\ &\leq \frac{8\eta^2}{[(1 - \epsilon)N]^2} \cdot \|X_{\text{TP}} X_{\text{TP}}^\top\|_2^2 \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2 \end{aligned} \quad (\text{B.12})$$

Then, by choosing $\eta \leq \frac{\lambda_{\min}(\Sigma)}{80\lambda_{\max}^2(\Sigma)}$, we have that the RHS in Equation (B.12) will be less than $\frac{\lambda_{\min}(\Sigma)}{8}$.

Step 3: Upper bounding the corrupted gradient. We now upper bound the third term in Equation (B.10).

$$\begin{aligned} \eta^2 \cdot \|\nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{FP})\|_F^2 &= \frac{4\eta^2}{[(1 - \epsilon)N]^2} \cdot \left\| \sum_{i \in \text{FP}} (\psi(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - \psi(\mathbf{x}_i^\top \mathbf{w}^*)) \cdot \mathbf{x}_i \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^{(t)} \geq 0\} \right\|_2^2 \\ &\leq \frac{4\eta^2}{[(1 - \epsilon)N]^2} \cdot \|\Sigma_{\text{FP}}(\mathbf{w}^{(t)}, \mathbf{w}^{(t)})\|_2 \sum_{i \in \text{FP}} (\psi(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - \psi(\mathbf{x}_i^\top \mathbf{w}^*))^2 \\ &\leq \frac{4\eta^2}{[(1 - \epsilon)N]^2} \cdot \|X_{\text{FP}} X_{\text{FP}}^\top\|_2 \sum_{i \in \text{FN}} (\psi(\mathbf{x}_i^\top \mathbf{w}^{(t)}) - \psi(\mathbf{x}_i^\top \mathbf{w}^*))^2 \end{aligned}$$

$$\leq \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot \|X_{\text{FP}}X_{\text{FP}}^\top\|_2 \|X_{\text{FN}}X_{\text{FN}}^\top\|_2 \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2$$

In the above, the first inequality follows from the same argument as Lemma 32, the second inequality follows from the optimality of the Subquantile set, and the final inequality follows from noting that ψ is 1-Lipschitz. We now conclude Steps 1-3 with our linear convergence result.

$$\begin{aligned} \|\mathbf{w}^{(t+1)} - \mathbf{w}^*\|_2 &\leq \|\mathbf{w}^{(t)} - \mathbf{w}^*\| \left(1 - \frac{\eta}{32} \cdot \lambda_{\min}(\Sigma) + \frac{\sqrt{8}\eta}{(1-\epsilon)N} \cdot \|X_{\text{FP}}\|_2 \|X_{\text{FN}}\| \right) \\ &\quad + \frac{16}{3} KC_\sigma \cdot \left(\frac{2d}{N} \log(5) + \frac{2}{N} \log(1/\delta) + 6\epsilon \log \epsilon^{-1} \right)^{1/2} \end{aligned}$$

Step 4: Concentration Inequalities. From our previous theorems, we have $\|X_{\text{FP}}\|_2 \|X_{\text{FN}}\|_2 \leq N\epsilon \sqrt{10C_3 \log \epsilon^{-1}}$ with probability exceeding $1 - \delta$ and $N = \Omega\left(\frac{d + \log(2/\delta)}{\epsilon}\right)$. If $N = \Omega\left(\frac{2d \log 5 + \log(1/\delta)}{3\epsilon \log \epsilon^{-1}}\right)$ and $\epsilon \leq \frac{1}{64\sqrt{80 \log(2)}}$, we obtain after $T = O\left(\kappa^2(\Sigma) \log\left(\frac{\|\mathbf{w}^*\|_2}{\epsilon}\right)\right)$ gradient descent iterations,

$$\|\mathbf{w}^{(t+1)} - \mathbf{w}^*\|_2 \leq \epsilon + 1024 \sqrt{\epsilon \log \epsilon^{-1}} = O\left(\sqrt{\epsilon \log \epsilon^{-1}}\right)$$

when we choose $\epsilon = O\left(\sqrt{\epsilon \log \epsilon^{-1}}\right)$. Our proof is complete. \blacksquare

C Probability Theory

In this section we will give various concentration inequalities on functions of the good data.

Lemma 24 (Upper Bound on Sum of Chi-Squared Variables [LM00]). *Suppose $\xi_i \sim \mathcal{N}(0, \sigma^2)$ for $i \in [n]$, then*

$$\Pr\{\|\xi\|_2^2 \geq \sigma(n + 2\sqrt{nx} + 2x)\} \leq e^{-x}$$

Proposition 25 (Probabilistic Upper Bound on Sum of Chi-Squared Variables). *Suppose $\xi_i \sim \mathcal{N}(0, \sigma^2)$ for $i \in [n]$. Let $S \subset [n]$ such that $|S| = \epsilon n$ for $\epsilon \in (0, 0.5)$ and let \mathcal{C} represent all such subsets. Given a failure probability $\delta \in (0, 1)$, when $n \geq \log(1/\delta)$, with probability exceeding $1 - \delta$,*

$$\max_{S \in \mathcal{C}} \|\xi_S\|_2^2 \leq \sigma(30n\epsilon \log \epsilon^{-1})$$

Proof. Directly from Lemma 24, we have with probability exceeding $1 - \delta$.

$$\|\xi\|_2^2 \leq \sigma\left(n + 2\sqrt{n \log(1/\delta)} + 2 \log(1/\delta)\right)$$

We now can prove the claimed bound using the layer-cake representation,

$$\Pr\left\{\max_{S \in \mathcal{C}} \|\xi\|_2^2 \geq \sigma(\epsilon n + 2\sqrt{\epsilon n x} + 2x)\right\} \leq \left(\frac{e}{\epsilon}\right)^{\epsilon n} \Pr\{\|\xi\|_2^2 \geq \sigma(\epsilon n + 2\sqrt{\epsilon n x} + 2x)\} \leq \left(\frac{e}{\epsilon}\right)^{\epsilon n} e^{-x}$$

In the first inequality we apply a union bound over \mathcal{C} with Lemma 36, and in the second inequality we use Lemma 24. We then obtain with probability exceeding $1 - \delta$,

$$\begin{aligned} \max_{S \in \mathcal{C}} \|\xi_S\|_2^2 &\leq \sigma\left(\epsilon n + 2\sqrt{n\epsilon \log(1/\delta)} + 3n^2\epsilon^2 \log \epsilon^{-1} + 2 \log(1/\delta) + 6n\epsilon \log \epsilon^{-1}\right) \\ &\leq \sigma\left(9n\epsilon \log \epsilon^{-1} + 2\sqrt{n\epsilon \log(1/\delta)} + 2\sqrt{3}n\epsilon \sqrt{\log \epsilon^{-1}} + 2 \log(1/\delta)\right) \\ &\leq \sigma\left(15n\epsilon \log \epsilon^{-1} + 2\sqrt{n\epsilon \log(1/\delta)} + 2 \log(1/\delta)\right) \\ &\leq \sigma(30n\epsilon \log \epsilon^{-1}) \end{aligned}$$

In the above, in the first inequality, we note that $\log \binom{n}{\epsilon n} \leq 3n\epsilon \log \epsilon^{-1}$ as $\epsilon < 0.5$, in the second inequality we note that $\sqrt{\log \epsilon^{-1}} \leq (\log(2))^{-1/2} \log \epsilon^{-1} \leq \sqrt{3} \log \epsilon^{-1}$ when $\epsilon < 0.5$, the final inequality holds when $n \geq \log(1/\delta)$ by solving for the quadratic equation. The proof is complete. \blacksquare

Lemma 26 (Sub-Gaussian Covariance Matrix Estimation [Ver10] Theorem 5.40). *Let $X \in \mathbb{R}^{d \times n}$ have columns sampled from a sub-Gaussian distribution with sub-Gaussian norm K and second-moment matrix Σ , then there exists positive constants c_k, C_Σ , dependent on the sub-Gaussian norm such that with probability at least $1 - 2e^{-c_k t^2}$,*

$$\lambda_{\max}(XX^\top) \leq n \cdot \lambda_{\max}(\Sigma) + \lambda_{\max}(\Sigma) \cdot \left(C_\Sigma \cdot \sqrt{dn} + t \cdot \sqrt{n} \right)$$

Lemma 27. *Let $X \in \mathbb{R}^{d \times n}$ have columns sampled from a sub-Gaussian distribution with sub-Gaussian norm K and second-moment matrix Σ . Let $S \subset [n]$ such that $|S| = \epsilon n$ for $\epsilon \in (0, 0.5)$ and let \mathcal{C} represent all such subsets. Then with probability at least $1 - \delta$,*

$$\begin{aligned} \max_{S \in \mathcal{C}} \lambda_{\max}(X_S X_S^\top) &\leq \lambda_{\max}(\Sigma) \cdot (10n\epsilon \log \epsilon^{-1}) \\ \min_{S \in \mathcal{C}} \lambda_{\min}(X_{X \setminus S} X_{X \setminus S}^\top) &\geq \frac{n}{4} \cdot \lambda_{\min}(\Sigma) \end{aligned}$$

when

$$n \geq \frac{2}{\epsilon} \cdot \left(C_\Sigma^2 \cdot d + \frac{\log(2/\delta)}{c_K} \right)$$

and $\epsilon \leq \frac{1}{60} \cdot \kappa^{-1}(\Sigma)$.

Proof. We will use the layer-cake representation to obtain our claimed error bound.

$$\begin{aligned} \Pr \left\{ \max_{S \in \mathcal{C}} \lambda_{\max}(X_S X_S^\top) \geq n\epsilon \cdot \lambda_{\max}(\Sigma) + \lambda_{\max}(\Sigma) \cdot \left(C_\Sigma \cdot \sqrt{dn\epsilon} + t\sqrt{n\epsilon} \right) \right\} \\ \leq \left(\frac{e}{\epsilon} \right)^{\epsilon n} \Pr \left\{ \lambda_{\max}(X_S X_S^\top) \geq n\epsilon \cdot \lambda_{\max}(\Sigma) + \lambda_{\max}(\Sigma) \cdot \left(C_\Sigma \cdot \sqrt{dn\epsilon} + t\sqrt{n\epsilon} \right) \right\} \leq 2 \cdot \left(\frac{e}{\epsilon} \right)^{\epsilon n} e^{-c_K t^2} \end{aligned}$$

In the above, the first inequality follows from a union bound over \mathcal{C} and Lemma 36, the second inequality follows from Lemma 26. Then from elementary inequalities, we obtain with probability $1 - \delta$,

$$\begin{aligned} \max_{S \in \mathcal{C}} \lambda_{\max}(X_S X_S^\top) &\leq n\epsilon \cdot \lambda_{\max}(\Sigma) + \lambda_{\max}(\Sigma) \cdot \left(C_\Sigma \cdot \sqrt{dn\epsilon} + \sqrt{\frac{1}{c_K} (n\epsilon \cdot \log(2/\delta) + 3n^2 \epsilon^2 \log \epsilon^{-1})} \right) \\ &\leq n \cdot \lambda_{\max}(\Sigma) \cdot (\epsilon + 3^{3/4} \epsilon \log \epsilon^{-1}) + \lambda_{\max}(\Sigma) \cdot \left(C_\Sigma \cdot \sqrt{dn\epsilon} + \sqrt{\frac{1}{c_K} n\epsilon \cdot \log(2/\delta)} \right) \\ &\leq \lambda_{\max}(\Sigma) \cdot (6n\epsilon \log \epsilon^{-1}) + \lambda_{\max}(\Sigma) \cdot \left(C_\Sigma \cdot \sqrt{dn\epsilon} + \sqrt{\frac{1}{c_K} n\epsilon \cdot \log(2/\delta)} \right) \\ &\leq \lambda_{\max}(\Sigma) \cdot (10n\epsilon \log \epsilon^{-1}) \end{aligned}$$

In the above, the last inequality holds when

$$n \geq \frac{2}{\epsilon} \cdot \left(C_\Sigma^2 \cdot d + \frac{\log(2/\delta)}{c_K} \right)$$

and our proof of the upper bound for the maximal eigenvalue is complete. We have from Weyl's Inequality for any $S \in \mathcal{C}$,

$$\lambda_{\min}(X_{X \setminus S} X_{X \setminus S}^\top) = \lambda_{\min}(XX^\top - X_S X_S^\top) \geq \lambda_{\min}(XX^\top) - \lambda_{\max}(X_S X_S^\top)$$

We then have with probability at least $1 - \delta$,

$$\lambda_{\min}(X_{X \setminus S} X_{X \setminus S}^\top) \geq n \cdot \lambda_{\min}(\Sigma) - C_\Sigma \cdot \sqrt{dn} - \sqrt{\frac{1}{c_K} \cdot n \cdot \log(2/\delta)} - \lambda_{\max}(\Sigma) \cdot (10n\epsilon \log \epsilon^{-1})$$

$$\geq \frac{n}{2} \cdot \lambda_{\min}(\Sigma) - \lambda_{\max}(\Sigma) \cdot (10n\epsilon \log \epsilon^{-1}) \geq \frac{n}{4} \cdot \lambda_{\min}(\Sigma)$$

In the above, the first inequality follows when $n \geq \frac{1}{\lambda_{\min}^2(\Sigma)} \left(8C_\Sigma \cdot d + \frac{2}{c_K} \cdot \log(2/\delta) \right)$, and the last inequality follows when $\epsilon \leq \frac{1}{60} \cdot \kappa^{-1}(\Sigma)$. The proof is complete. \blacksquare

Lemma 28 (Noise Bound). *Let $X = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ be the data matrix such that for $i \in [N]$, \mathbf{x}_i are sampled from a sub-Gaussian distribution with second-moment matrix Σ with sub-Gaussian norm C_Σ and ξ_i are sampled from sub-Gaussian distribution with sub-Gaussian norm C_σ . Assume $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is bounded over \mathbb{R} and Lipschitz. Let \mathcal{S} represent all subsets of $[N]$ of size empty to $(1-\epsilon)N$. Set a failure probability $\delta \in (0, 1)$, then with probability at least $1 - \delta$,*

$$\max_{S \in \mathcal{S}} \sup_{\mathbf{w} \in \mathcal{B}(\mathbf{0}, R)} \|(\xi_S \circ f(X_S^\top \mathbf{w}))^\top X_S\|_2 \leq 2NC_\Sigma C_\sigma \|f\|_\infty \left(\frac{2d}{N} \log(6) + \frac{2Rd}{N} \log(3N) + \frac{2}{N} \log(2/\delta) + 6\epsilon \log \epsilon^{-1} \right)^{1/2}$$

Proof. We will use the following characterization of the spectral norm.

$$\left\| \sum_{i \in \text{TP}} f(\mathbf{x}_i^\top \mathbf{w}) \xi_i \mathbf{x}_i \right\|_2 = \max_{\mathbf{v} \in \mathbb{S}^{d-1}} \left| \sum_{i \in \text{TP}} f(\mathbf{x}_i^\top \mathbf{w}) \xi_i \mathbf{x}_i^\top \mathbf{v} \right|$$

Our proof uses standard ideas in High-Dimensional Probability and will follow similarly to [ZYWG19]. We will first show that $f(\mathbf{x}_i^\top \mathbf{w}) \mathbf{x}_i$ is sub-Gaussian. We first note for any $\mathbf{v} \in \mathbb{S}^{d-1}$, the random variable, $\mathbf{x}_i^\top \mathbf{v}$ is sub-Gaussian by definition. We then have,

$$(\mathbf{E} |f(\mathbf{x}_i^\top \mathbf{w}) \mathbf{x}_i^\top \mathbf{v}|^p)^{1/p} \stackrel{(i)}{\leq} (\mathbf{E} |f(\mathbf{x}_i^\top \mathbf{w})|^{2p} \mathbf{E} |\mathbf{x}_i^\top \mathbf{v}|^{2p})^{1/2p} \stackrel{(ii)}{\leq} (\|f\|_\infty C_\Sigma \sqrt{2}) \sqrt{p}$$

In the above, (i) follows from Hölder's Inequality (see Proposition 35), (ii) follows from noting from letting $q = 2p$ and noting from Definition 7 that $\|\mathbf{x}_i^\top \mathbf{v}\|_{L_q}$ is upper bounded by $C_\Sigma \sqrt{q}$. We thus have $f(\mathbf{x}_i^\top \mathbf{w}) \mathbf{x}_i^\top \mathbf{v}$ is sub-Gaussian for any $\mathbf{w} \in \mathbb{R}^d$ and $\|f(\mathbf{x}_i^\top \mathbf{w}) \mathbf{x}_i^\top \mathbf{v}\|_{\psi_2} \leq \sqrt{2} C_\Sigma \|f\|_\infty$. Recall $C_\sigma \triangleq \|\xi_i\|_{\psi_2}$, then from Lemma 8, the random variable $\xi_i f(\mathbf{x}_i^\top \mathbf{w}) \mathbf{x}_i^\top \mathbf{v}$ is sub-exponential s.t. $\|\xi_i f(\mathbf{x}_i^\top \mathbf{w}) \mathbf{x}_i^\top \mathbf{v}\|_{\psi_1} \leq \sqrt{2} C_\Sigma C_\sigma \|f\|_\infty$. Let $\tilde{\mathbf{w}} \in \mathcal{M}$ such that $\tilde{\mathbf{w}} = \arg \min_{\mathbf{u} \in \mathcal{M}} \|\mathbf{w} - \mathbf{u}\|_2$, where \mathcal{M} is a ε -cover of $\mathcal{B}(\mathbf{0}, R)$, then

$$\begin{aligned} \Pr \left\{ \max_{S \in \mathcal{S}} \sup_{\mathbf{w} \in \mathcal{B}(\mathbf{0}, R)} \frac{1}{(1-\epsilon)N} \cdot \left\| \sum_{i \in S} \xi_i f(\mathbf{x}_i^\top \mathbf{w}) \mathbf{x}_i \right\|_2 \geq t \right\} &\leq \Pr \left\{ \max_{S \in \mathcal{S}} \max_{\mathbf{u} \in \mathcal{M}} \frac{1}{(1-\epsilon)N} \cdot \left\| \sum_{i \in S} \xi_i f(\mathbf{x}_i^\top \mathbf{u}) \mathbf{x}_i \right\|_2 \geq \frac{t}{2} \right\} \\ &+ \Pr \left\{ \max_{S \in \mathcal{S}} \sup_{\mathbf{w} \in \mathcal{B}(\mathbf{0}, R)} \frac{1}{(1-\epsilon)N} \cdot \left\| \sum_{i \in S} \xi_i f(\mathbf{x}_i^\top \mathbf{w}) \mathbf{x}_i - \sum_{i \in S} \xi_i f(\mathbf{x}_i^\top \tilde{\mathbf{w}}) \mathbf{x}_i \right\|_2 \geq \frac{t}{2} \right\} \quad (\text{C.1}) \end{aligned}$$

We will now use a ε -covering argument to bound the first term of Equation (C.1). Let \mathcal{C} be a ε -net of \mathbb{S}^{d-1} such that for any $\mathbf{v} \in \mathbb{S}^{d-1}$, there exists $\mathbf{u} \in \mathcal{C}$ such that $\|\mathbf{u} - \mathbf{v}\|_2 \leq \varepsilon$. Let $\mathbf{u}^* = \arg \max_{\mathbf{u} \in \mathcal{C}} |(\xi \circ f(X^\top \tilde{\mathbf{w}}))^\top X \mathbf{u}|$ and $\mathbf{v}^* = \arg \max_{\mathbf{v} \in \mathbb{S}^{d-1}} |(\xi \circ f(X^\top \tilde{\mathbf{w}}))^\top X \mathbf{v}|$. We then have from the triangle inequality,

$$|(\xi \circ f(X^\top \tilde{\mathbf{w}}))^\top X \mathbf{v}^* - (\xi \circ f(X^\top \tilde{\mathbf{w}}))^\top X \mathbf{u}^*| \leq |(\xi \circ f(X^\top \tilde{\mathbf{w}}))^\top X| \|\mathbf{u}^* - \mathbf{v}^*\|_2 \leq \varepsilon \cdot |(\xi \circ f(X^\top \tilde{\mathbf{w}}))^\top X|$$

where in the final inequality we use the definition of a ε -net. We then have from the reverse triangle inequality,

$$\begin{aligned} |(\xi \circ f(X^\top \tilde{\mathbf{w}}))^\top X \mathbf{u}^*| &\geq |(\xi \circ f(X^\top \tilde{\mathbf{w}}))^\top X \mathbf{v}^*| - |(\xi \circ f(X^\top \tilde{\mathbf{w}}))^\top X \mathbf{u}^* - (\xi \circ f(X^\top \tilde{\mathbf{w}}))^\top X \mathbf{v}^*| \\ &\geq (1 - \varepsilon) |(\xi \circ f(X^\top \tilde{\mathbf{w}}))^\top X \mathbf{v}^*| \end{aligned}$$

From rearranging, we have

$$|(\xi \circ f(X^\top \tilde{\mathbf{w}}))^\top X \mathbf{v}^*| \leq \frac{1}{1 - \varepsilon} \cdot |(\xi \circ f(X^\top \tilde{\mathbf{w}}))^\top X \mathbf{u}^*|$$

With this result we are ready to make the probabilistic bounds. Suppose \mathcal{S} represents all subsets of $[N]$ of size $(1 - 2\epsilon)N$ to $(1 - \epsilon)N$. Suppose \mathcal{C} is a ε_c net of \mathbb{S}^{d-1} and \mathcal{M} is a ε_m net of $\mathcal{B}(\mathbf{0}, R)$. Then,

$$\begin{aligned}
& \Pr\left\{\frac{1}{(1-\epsilon)N} \cdot \max_{S \in \mathcal{S}} \max_{\mathbf{w} \in \mathcal{M}} \|(\boldsymbol{\xi}_S \circ f(X_S^\top \mathbf{w}))^\top X_S\| \geq \frac{t}{2}\right\} \\
& \leq \Pr\left\{\frac{1}{1-\varepsilon_c} \cdot \frac{1}{(1-\epsilon)N} \cdot \max_{S \in \mathcal{S}} \max_{\mathbf{v} \in \mathcal{C}} \max_{\mathbf{w} \in \mathcal{M}} |(\boldsymbol{\xi}_S \circ f(X_S^\top \mathbf{w}))^\top X_S \mathbf{v}| \geq \frac{t}{2}\right\} \\
& \leq \sum_{j \in [|\mathcal{S}|]} \sum_{i \in [|\mathcal{M}|]} \sum_{k \in [|\mathcal{C}|]} \Pr\left\{\frac{1}{1-\varepsilon_c} \cdot \frac{1}{(1-\epsilon)N} \cdot |(\boldsymbol{\xi}_{S_j} \circ f(X_{S_j}^\top \mathbf{w}_i))^\top X_{S_j} \mathbf{v}_k| \geq \frac{t}{2}\right\} \\
& \stackrel{(iii)}{\leq} \left(\frac{3}{\varepsilon_m}\right)^{Rd} \left(\frac{3}{\varepsilon_c}\right)^d \cdot \left(\frac{e}{\epsilon}\right)^{N\epsilon} \exp\left(-\frac{3}{16} \cdot (1-\epsilon)N \left(\frac{t^2}{2C_\Sigma^2 C_\sigma^2 \|f\|_\infty^2} \wedge \frac{t}{\sqrt{2} C_\Sigma C_\sigma \|f\|_\infty}\right)\right) \leq \frac{\delta}{2}
\end{aligned}$$

In the above, (iii) follows from Bernstein's Inequality (see Lemma 9). We can now note that $\log\left(\frac{N}{(1-\epsilon)N}\right) = \log\left(\frac{N}{\epsilon N}\right)$. Then to satisfy the above probabilistic condition, it must hold that

$$t \geq 2C_\Sigma C_\sigma \|f\|_\infty \left(\frac{2d}{N} \log\left(\frac{3}{\varepsilon_c}\right) + \frac{2Rd}{N} \log\left(\frac{3}{\varepsilon_m}\right) + \frac{2}{N} \log(2/\delta) + 6\epsilon \log \epsilon^{-1}\right)^{1/2}$$

We now bound the second term of Equation (C.1). For any \mathbf{w} , recall $\tilde{\mathbf{w}} = \arg \min_{\mathbf{u} \in \mathcal{M}} \|\mathbf{w} - \mathbf{u}\|_2$ and thus for any \mathbf{w} , we have $\|\mathbf{w} - \tilde{\mathbf{w}}\|_2 \leq \varepsilon_m$. It then follows,

$$\begin{aligned}
\sup_{\mathbf{w} \in \mathcal{B}(\mathbf{0}, R)} \frac{1}{(1-\epsilon)N} \left\| \sum_{i \in \mathcal{S}} \xi_i f(\mathbf{x}_i^\top \mathbf{w}) \mathbf{x}_i - \sum_{i \in \mathcal{S}} \xi_i f(\mathbf{x}_i^\top \tilde{\mathbf{w}}) \mathbf{x}_i \right\|_2 & \leq \sup_{\mathbf{w} \in \mathcal{B}(\mathbf{0}, R)} \max_{i \in [N]} \|\xi_i \mathbf{x}_i\|_2 \|f(\mathbf{x}_i^\top \mathbf{w}) - f(\mathbf{x}_i^\top \tilde{\mathbf{w}})\|_2 \\
& \leq \varepsilon_m \|f'\|_\infty \max_{i \in [N]} \|\xi_i \mathbf{x}_i\|_2 \max_{i \in [N]} \|\mathbf{x}_i\|_2 \\
& \stackrel{(iv)}{\leq} \frac{\varepsilon_m \|f'\|_\infty}{2} \left(\max_{i \in [N]} \|\xi_i \mathbf{x}_i\|_2^2 + \max_{i \in [N]} \|\mathbf{x}_i\|_2^2 \right)
\end{aligned}$$

In the above, (iv) follows from Young's Inequality (see Proposition 34). Note that $\mathbf{x}_i^\top \mathbf{v}$ for any $\mathbf{v} \in \mathbb{S}^{d-1}$ and $i \in [N]$ is sub-Gaussian, we thus have,

$$\begin{aligned}
\Pr\left\{\frac{\varepsilon_m \|f'\|_\infty}{2} \cdot \max_{j \in [N]} \|\mathbf{x}_j\|_2^2 \geq \frac{t}{4}\right\} & \leq \frac{1}{1-\varepsilon_c} \cdot \sum_{i \in \mathcal{C}} \sum_{j \in [N]} \Pr\left\{|\mathbf{x}_j^\top \mathbf{v}_i|^2 \geq \frac{t}{2\varepsilon_m \|f'\|_\infty}\right\} \\
& \leq \frac{N}{1-\varepsilon_c} \cdot \left(\frac{3}{\varepsilon_c}\right)^d \cdot \exp\left(-\frac{1}{2} \left(\frac{t}{2LC_\Sigma \varepsilon_m}\right)^2 \wedge \left(\frac{t}{2LC_\Sigma \varepsilon_m}\right)\right) \leq \frac{\delta}{4}
\end{aligned}$$

where in the above, the final inequality follows when

$$t \geq 2LC_\Sigma \varepsilon_m \left(2 \log\left(\frac{1}{1-\varepsilon_c}\right) + 2 \log N + d \log\left(\frac{3}{\varepsilon_c}\right) + 2 \log(4/\delta)\right)$$

To bound $\max_{i \in \mathcal{S}} \|\xi_i \mathbf{x}_i\|_2$, we note that for any $\mathbf{v} \in \mathbb{R}^d$ that $\xi_i \mathbf{x}_i^\top \mathbf{v}$ is sub-exponential with norm $\|\xi_i \mathbf{x}_i^\top \mathbf{v}\|_{\psi_1} \leq C_\sigma C_\Sigma$. Similarly,

$$\Pr\left\{\frac{\varepsilon_m \|f'\|_\infty}{2} \max_{i \in \mathcal{S}} \|\xi_i \mathbf{x}_i\|_2^2 \geq \frac{t}{4}\right\} \leq \frac{\delta}{4}$$

The final inequality holds when,

$$t \geq 2C_\sigma C_\Sigma \|f'\|_\infty \varepsilon_m \left(2 \log\left(\frac{1}{1-\varepsilon_c}\right) + 2 \log(N) + d \log\left(\frac{3}{\varepsilon_c}\right) + 2 \log(4/\delta)\right)$$

We now choose $\varepsilon_c \triangleq \frac{1}{2}$ and $\varepsilon_m = \frac{1}{N}$. Combining our estimates, we obtain,

$$\Pr\left\{\frac{1}{(1-\epsilon)N} \cdot \max_{S \in \mathcal{S}} \|(\boldsymbol{\xi}_S \circ f(X_S^\top \mathbf{w}))^\top X_S\| \geq t\right\} \leq \delta$$

when

$$t \geq 2C_\Sigma C_\sigma \|f\|_\infty \left(\frac{2d}{N} \log(6) + \frac{2Rd}{N} \log(3N) + \frac{2}{N} \log(2/\delta) + 6\epsilon \log \epsilon^{-1} \right)^{1/2}$$

Our proof is complete. \blacksquare

Lemma 29. Let $\mu(\mathbf{x})$ represent the PDF of $\mathcal{N}(\mathbf{0}, I)$. Fix $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^d$ and define the sets $\Omega_1 = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{x}^\top \mathbf{w}_1 \geq 0\}$ and $\Omega_2 = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{x}^\top \mathbf{w}_2 \geq 0\}$, let $\Theta = \arccos\left(\frac{\mathbf{w}_1^\top \mathbf{w}_2}{\|\mathbf{w}_1\| \|\mathbf{w}_2\|}\right)$, then

$$\lambda_{\min}\left(\int_{\Omega_1, \Omega_2} \mathbf{x} \mathbf{x}^\top d\mu(\mathbf{x})\right) \succeq (1 \wedge (\pi - \Theta - \sin \Theta)) \cdot I$$

Proof. From Lemma 1 in [DLT⁺18], we have for $U = \left[\bar{\mathbf{w}}_1, \frac{\bar{\mathbf{w}}_2 - \cos \Theta (\bar{\mathbf{w}}_1, \bar{\mathbf{w}}_2)}{\sin \Theta (\bar{\mathbf{w}}_1, \bar{\mathbf{w}}_2)}, \mathbf{u}_3, \dots, \mathbf{u}_n\right]^\top$.

$$\begin{aligned} \mathbf{E}[\mathbf{x} \mathbf{x}^\top \mid \mathbf{x}^\top \mathbf{w}_1 \geq 0, \mathbf{x}^\top \mathbf{w}_2 \geq 0] &= \mathbf{E}[U^\top U \mathbf{x} \mathbf{x}^\top U^\top U \mid \mathbf{x}^\top \mathbf{w}_1 \geq 0, \mathbf{x}^\top \mathbf{w}_2 \geq 0] \\ &= U^\top \cdot \mathbf{E}[U \mathbf{x} \mathbf{x}^\top U^\top \mid \mathbf{x}^\top \mathbf{w}_1 \geq 0, \mathbf{x}^\top \mathbf{w}_2 \geq 0] \cdot U \\ &= U^\top \cdot \begin{bmatrix} \frac{1}{2\pi}(\pi - \Theta + \sin \Theta \cos \Theta) & \frac{1}{2\pi} \sin^2 \Theta & 0 & \cdots & 0 \\ \frac{1}{2\pi} \sin^2 \Theta & \frac{1}{2\pi}(\pi - \Theta + \sin \Theta \cos \Theta) & & & 0 \\ 0 & & 1 & & \vdots \\ \vdots & & & \ddots & \\ 0 & & & & 1 \end{bmatrix} \cdot U \end{aligned}$$

The matrix has $n - 2$ unitary eigenvalues and the final two eigenvalues can be calculated to be $\pi - \Theta \pm \sin \Theta$. We then have

$$\lambda_{\min}(\mathbf{E}[\mathbf{x} \mathbf{x}^\top \mid \mathbf{x}^\top \mathbf{w}_1 \geq 0, \mathbf{x}^\top \mathbf{w}_2 \geq 0]) \geq 1 \wedge (\pi - \Theta - \sin \Theta)$$

Our proof is complete. \blacksquare

Lemma 30. Fix $\mathbf{w}^* \in \mathbb{R}^{d-1}$ and suppose $\mathbf{w} \in \mathcal{B}(\mathbf{w}^*, R)$ for a constant $R < \|\mathbf{w}^*\|$. Sample $\mathbf{x}_1, \dots, \mathbf{x}_N$ i.i.d from a sub-Gaussian distribution with second-moment matrix Σ and sub-Gaussian norm C_Σ . Suppose S represents a subset of $[N]$ s.t. $|S| \leq (1 - \epsilon)N$. Then with probability at least $1 - \delta$,

$$\left\| \sum_{i \in S} \mathbf{x}_i \mathbf{x}_i^\top \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^* \geq 0\} \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^{(t)} \geq 0\} - \mathbf{E}[\mathbf{x} \mathbf{x}^\top \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^* \geq 0\} \cdot \mathbb{I}\{\mathbf{x}_i^\top \mathbf{w}^{(t)} \geq 0\}] \right\|_2 \leq \Xi$$

Proof. We will use the decomposition given in [ZYWG19]. Let $\tilde{\mathbf{w}} = \arg \min_{\mathbf{v} \in \mathcal{C}} \|\mathbf{w} - \mathbf{v}\|_2$ where \mathcal{C} is a ϵ_c -cover of $\mathcal{B}(\mathbf{w}^*, R)$ and let \mathcal{S} be defined as in the Lemma statement.

$$\begin{aligned} &\Pr\left\{ \max_{S \in \mathcal{S}} \sup_{\mathbf{w} \in \mathcal{B}(\mathbf{w}^*; R)} \|\Sigma_S(\mathbf{w}, \mathbf{w}^*) - \mathbf{E}[\Sigma_S(\mathbf{w}, \mathbf{w}^*)]\|_2 \geq t \right\} \\ &\leq \Pr\left\{ \max_{S \in \mathcal{S}} \sup_{\mathbf{w} \in \mathcal{B}(\mathbf{w}^*; R)} \|\Sigma_S(\mathbf{w}, \mathbf{w}^*) - \Sigma_S(\tilde{\mathbf{w}}, \mathbf{w}^*)\|_2 \geq \frac{t}{3} \right\} \\ &+ \Pr\left\{ \max_{S \in \mathcal{S}} \sup_{\tilde{\mathbf{w}} \in \mathcal{C}} \|\Sigma_S(\tilde{\mathbf{w}}, \mathbf{w}^*) - \mathbf{E}[\Sigma_S(\tilde{\mathbf{w}}, \mathbf{w}^*)]\|_2 \geq \frac{t}{3} \right\} \\ &+ \Pr\left\{ \sup_{\mathbf{w} \in \mathcal{B}(\mathbf{w}^*; R)} \|\mathbf{E}[\Sigma_S(\tilde{\mathbf{w}}, \mathbf{w}^*)] - \mathbf{E}[\Sigma_S(\mathbf{w}, \mathbf{w}^*)]\|_2 \geq \frac{t}{3} \right\} \end{aligned}$$

We bound all terms separately.

$$\sup_{\mathbf{w} \in \mathcal{B}(\mathbf{w}^*; R)} \|\Sigma_S(\mathbf{w}, \mathbf{w}^*) - \Sigma_S(\tilde{\mathbf{w}}, \mathbf{w}^*)\|_2 \leq \sup_{\mathbf{w} \in \mathcal{B}(\mathbf{w}^*; R)} \|\Sigma_S(\mathbf{w}, -\tilde{\mathbf{w}})\|_2 + \sup_{\mathbf{w} \in \mathcal{B}(\mathbf{w}^*; R)} \|\Sigma_S(-\mathbf{w}, \tilde{\mathbf{w}})\|_2$$

\blacksquare

Lemma 31. Fix $S \in \mathbb{R}^{K \times N\epsilon}$, $T \in \mathbb{R}^{N\epsilon \times L}$, then sample a matrix $G \in \mathbb{R}^{N\epsilon \times N\epsilon}$ such that each column of G represents a ϵ -subset of a n -dimensional vector sampled from $\mathcal{N}(\mathbf{0}, \sigma^2 \cdot I)$, then with probability exceeding $1 - \delta$,

$$\|SGT\|_{\mathbb{F}}^2 \leq \|S\|_{\mathbb{F}}^2 \|T\|_{\mathbb{F}}^2 \cdot \sigma \sqrt{2 \log \left(\frac{2N^2}{\delta} \right)}$$

Proof. The proof will be a calculation.

$$\|SGT\|_{\mathbb{F}}^2 = \sum_{i \in [K]} \sum_{j \in [L]} \sum_{k_1, k_2 \in [N\epsilon] \times [N\epsilon]} (S_{i,k_1} G_{k_1, k_2} T_{k_2, j})^2 \leq \|S\|_{\mathbb{F}}^2 \|T\|_{\mathbb{F}}^2 \max_{i, j \in [N\epsilon] \times [N\epsilon]} (G_{i, j})^2$$

It then suffices to bound the maximum value of a Gaussian squared over N^2 samples. From Lemma 37, we have

$$\Pr \left\{ \max_{i, j} G_{i, j}^2 \geq t \right\} = \Pr \left\{ \max_{i, j} |G_{i, j}| \geq \sqrt{t} \right\} \leq 2Ne^{-\frac{t}{2\sigma^2}}$$

We thus obtain from elementary inequalities, with probability at least $1 - \delta$,

$$\max_{i, j} G_{i, j}^2 \leq \sigma \left(2 \log \left(\frac{2N^2}{\delta} \right) \right)$$

Our proof is complete. ■

D Mathematical Tools

In this section, we state additional lemmas referenced throughout the text for completeness.

Lemma 32. Let $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ and $X \in \mathbb{R}^{p \times n}$, then

$$\left\| \sum_{i=1}^n a_i b_i \mathbf{x}_i \right\|_2^2 \leq \|\mathbf{a}\|_{\infty}^2 \|\mathbf{b}\|_2^2 \|X X^{\top}\|_2$$

Proof. The proof is a simple calculation. Expanding out the LHS, we have

$$\left\| \sum_{i=1}^n a_i b_i \mathbf{x}_i \right\|_2^2 = \sum_{i=1}^n \sum_{j=1}^n a_i a_j b_i b_j \mathbf{x}_i^{\top} \mathbf{x}_j = (\mathbf{a} \circ \mathbf{b})^{\top} X^{\top} X (\mathbf{a} \circ \mathbf{b}) \leq \|\mathbf{a} \circ \mathbf{b}\|_2^2 \|X^{\top} X\|_2 \leq \|\mathbf{a}\|_{\infty}^2 \|\mathbf{b}\|_2^2 \|X^{\top} X\|_2$$

where the final inequality comes from noting

$$\|\mathbf{a} \circ \mathbf{b}\|^2 = \sum_{i=1}^n a_i^2 b_i^2 \leq \max_{i \in [n]} a_i^2 \cdot \sum_{i=1}^n b_i^2$$

Our proof is complete. ■

Lemma 33 (Lemma 3.11 [B⁺15]). Let f be β -smooth and α -strongly convex over \mathbb{R}^n , then for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \frac{\alpha\beta}{\alpha + \beta} \|\mathbf{x} - \mathbf{y}\|^2 + \frac{1}{\alpha + \beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2$$

Proposition 34 (Young's Inequality [You12]). Suppose $a, b \in \mathbb{R}_+$, then for $p, q > 1$ such that $\frac{1}{p} + \frac{1}{q} = 1$, then

$$ab \leq \frac{a^p}{p} + \frac{b^q}{q}$$

Proposition 35 (Hölder's Inequality [H89]). *Consider the probability space $(\Omega, \mathcal{F}, \mathbf{Pr})$. Suppose the random variables X and Y are defined over Ω , for $p, q > 1$ such that $\frac{1}{p} + \frac{1}{q} = 1$, it holds*

$$\mathbf{E}|XY| \leq (\mathbf{E}|X|^p)^{1/p} (\mathbf{E}|Y|^q)^{1/q}$$

Lemma 36 (Sum of Binomial Coefficients [CLRS22]). *Let $k, n \in \mathbb{N}$ such that $k \leq n$, then*

$$\sum_{i=0}^k \binom{n}{i} \leq \left(\frac{en}{k}\right)^k$$

Lemma 37 (Max Gaussian [RH23]). *Let x_1, \dots, x_n be sampled i.i.d from $\mathcal{N}(0, \sigma^2)$. Then,*

$$\mathbf{Pr}\{\|\mathbf{x}\|_\infty > t\} \leq N \exp\left(-\frac{t^2}{\frac{16}{3}\sigma^2}\right)$$

Lemma 38 (Corollary 4.2.13 in [Ver20]). *The covering number of the ℓ_2 -norm ball $\mathcal{B}(\mathbf{0}; 1)$ for $\varepsilon < 0$, satisfies,*

$$\mathcal{N}(\mathcal{B}_{\ell_2}^d(\mathbf{0}, 1), \varepsilon) \leq \left(\frac{3}{\varepsilon}\right)^d$$