

Non-Linear Learning in the Huber ϵ -Contamination Model

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

Alex Gittens
RPI CS, gittea@rpi.edu

Abstract

In this paper we study Subquantile Minimization for learning the Adversarial Huber- ϵ Contamination Problem for Kernel Learning. We first reduce the Subquantile Minimization Algorithm to Iterative Thresholding using ideas from convex optimization. Let the target data be distributed as $y = \sum_{k=1}^K \sigma(\mathbf{x}^T \mathbf{w}_k^*) + \xi_i$ for a non-linear function, σ , s.t. $\text{range}(\sigma') \in [C_1, C_2]$ for positive constants C_1, C_2 and noise $\xi_i \sim \mathcal{N}(0, \sigma^2)$ and \mathbf{x}_i be sampled from a centered sub-Gaussian distribution with multivariate proxy, $\mathbf{\Gamma}$. Our main result is for sufficiently large n , there exists an algorithm that returns $\mathbf{W}^{(T)}$ with high probability such that $\|\mathbf{W}^* - \mathbf{W}^{(T)}\|_F \leq \epsilon$ in time,

$$O(\text{polylog}(n, \|\mathbf{W}^*\|_F, \kappa(\mathbf{\Gamma}), K, (1/C_1), C_2, (1/\epsilon), 1 - \epsilon))$$

for a sufficiently small ϵ dependent on $\mathbf{\Gamma}$, N , and $\|\mathbf{X}_Q\|$. Furthermore, we consider the noisy Kernelized Generalized Linear Model (GLM) where $y = \omega(f^*(\mathbf{x})) + \xi$ where $\xi \sim \mathcal{N}(0, \sigma^2)$ and prove the same algorithm returns $f^{(T)}$ with high probability such that $\|f^{(T)} - f^*\| \leq \epsilon + O(\sigma)$ after T iterations. The iterative thresholding algorithm has been used in large neural network models in prior research [SS19]. Our work provides the first steps for theoretical guarantees for neural networks and non-linear models for iterative thresholding algorithms in the Huber- ϵ contamination model.

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 (Huber ϵ -Contamination Model [HR09]). *Given a corruption parameter $0 < \epsilon < 0.5$, a data matrix, \mathbf{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 \mathbf{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 \mathbf{X} be a sub-Gaussian data matrix, and $\mathbf{y} = \mathbf{X}^T \mathbf{w}^* + \mathbf{e}$ where \mathbf{e} is 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.*

We will now give our results for the Kernelized GLM problem.

Theorem 4 (Informal of Theorem 12). *Let the dataset be given as $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ such that for $i \in P$, $y_i = \omega(f^*(\mathbf{x}_i)) + \xi_i$. Then there exists an algorithm such $\|f^{(t)} - f^*\|_{\mathcal{H}} \leq \epsilon + O(\|\Gamma\|\sigma) + O(\sigma)$ after $t = O\left(n \log\left(\frac{\|f^*\|_{\mathcal{H}}}{\epsilon}\right)\right)$ iterations.*

1.1 Contributions

Our main contribution is the approximation bounds for Subquantile Minimization for non-linear learning problems from the iterative thresholding algorithm developed in Algorithm 1. Our proof techniques extend [BJK15, ADKS22] as we do not assume the covariates follow the spherical Gaussian property, as such a property will not hold for any infinite-dimensional Hilbert Space.

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 $\mathbf{P}, \mathbf{Q}, \mathbf{R}$ as matrices. Throughout the paper, we will use the following matrix norms for a matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$,

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

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

2.1 Reproducing Kernel Hilbert Spaces

Let the function $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$ represent the Hilbert Space Representation or ‘feature transform’ from a vector in the original covariate space to the RKHS. We define $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ as $k(\mathbf{x}, \mathbf{x}) \triangleq \langle \phi(\mathbf{x}), \phi(\mathbf{x}) \rangle_{\mathcal{H}}$. For a function in a RKHS, $f \in \mathcal{H}$, it follows for a function f parameterized by weights $\mathbf{w} \in \mathbb{R}^n$, that the point evaluation function is given as $f : \mathbb{R}^d \rightarrow \mathbb{R}$ and defined $f(\cdot) \triangleq \sum_{i \in [n]} w_i k(\mathbf{x}_i, \cdot)$.

Definition 5 (Reproducing Property). *Let $\mathbf{x} \in \mathcal{X}$, then for any $f \in \mathcal{H}$,*

$$f(\mathbf{x}) = \langle f, k(\mathbf{x}, \cdot) \rangle_{\mathcal{H}} = \langle f, \phi(\mathbf{x}) \rangle_{\mathcal{H}}$$

Definition 6 (Pushforward Measure). *Let $\phi : \mathcal{X} \rightarrow \mathcal{H}$ represent the mapping from the input dimension to the Hilbert Space, and let \mathcal{P} be the Probability Measure of the uncorrupted data over \mathcal{X} . Then $\mathcal{P}_\# \phi(X) = \mathcal{P}(\phi^{-1}(X))$ represents the measure over the Hilbert Space \mathcal{H} using the measure of the good data defined over the original data space \mathcal{X} .*

The norm of a function $f \in \mathcal{H}$ is given as $\|f\|_{\mathcal{H}} = \sqrt{\langle f, f \rangle_{\mathcal{H}}}$.

2.2 Tensor Products

Let \mathcal{H}, \mathcal{K} be Hilbert Spaces, then $\mathcal{H} \otimes \mathcal{K}$ is the tensor product space and is also a Hilbert Space [RaR02]. For $\phi_1, \psi_1 \in \mathcal{H}$ and $\phi_2, \psi_2 \in \mathcal{K}$, the inner product is defined as $\langle \phi_1 \otimes \phi_2, \psi_1 \otimes \psi_2 \rangle_{\mathcal{H} \otimes \mathcal{K}} = \langle \phi_1, \psi_1 \rangle_{\mathcal{H}} \langle \phi_2, \psi_2 \rangle_{\mathcal{K}}$. We will utilize tensor products when we discuss infinite dimensional covariance estimation.

2.3 Sub-Gaussian Random Functions in the Hilbert Space

In this paper we sample the target covariates $\mathbf{x} \sim \mathcal{X}$ such that $\phi(\mathbf{x}) \triangleq X \sim \mathcal{P}_\# \phi$ is sub-Gaussian in the Hilbert Space where $\mathbf{E}[X] = \mathbf{0}$ and covariance $\mathbf{E}[X \otimes X] = \mathbf{\Sigma}$ with proxy $\mathbf{\Gamma}$, where $\mathbf{\Sigma} \preceq 4\|X\|_{\psi_2}^2 \mathbf{\Gamma}$, where we denote \preceq as the Löwner order. We have X is a centered Hilbert Space sub-Gaussian random function if for all $\theta > 0$,

$$\mathbf{E}_{X \sim \mathcal{P}} [\exp(\theta \langle X, v \rangle_{\mathcal{H}})] \leq \exp\left(\frac{\alpha^2 \theta^2 \langle v, \mathbf{\Gamma} v \rangle_{\mathcal{H}}}{2}\right) \quad (2.1)$$

where the sub-Gaussian Norm for a centered Hilbert Space Function is given as

$$\|X\|_{\psi_2} \triangleq \inf \left\{ \alpha \geq 0 : \mathbf{E} \left[e^{\langle v, X \rangle_{\mathcal{H}}} \right] \leq e^{\alpha^2 \langle v, \mathbf{\Gamma} v \rangle_{\mathcal{H}} / 2} : \forall v \in \mathcal{H} \right\}$$

Then we say $X \sim \mathcal{SG}(\mathbf{\Gamma}, \alpha)$, where if $\alpha = 1$, we will say $X \sim \mathcal{SG}(\mathbf{\Gamma})$. The Gaussian Design for the Feature Space has gained popularity in the study of kernel learning [CLKZ21]. The sub-Gaussian design is the standard assumed distribution in the robust statistics literature, [JLT20, ADKS22], and has been studied extensively in the context of iterative thresholding algorithms for linear regression.

2.4 Assumptions

We will first give our assumptions for robust kernelized regression.

Assumption 7 (Sub-Gaussian Design). *We assume for $\mathbf{x}_i \sim \mathcal{X}$, then it follows for the function to the Hilbert Space, $\phi(\cdot) : \mathcal{X} \rightarrow \mathcal{H}$,*

$$\phi(\mathbf{x}) \triangleq X \sim \mathcal{P}_{\#} \phi \triangleq \mathcal{SG}(\mathbf{\Gamma}, 1/2)$$

where $\mathbf{\Gamma}$ is a possibly infinite dimensional covariance operator.

Assumption 8 (Bounded Functions). *We assume for $\mathbf{x}_i \sim \mathcal{P} \in \mathcal{X}$, then it follows for the feature map, $\phi(\cdot) : \mathcal{X} \rightarrow \mathcal{H}$,*

$$\sup_{\mathbf{x} \in \mathcal{X}} \|\phi(\mathbf{x})\|_{\mathcal{H}}^2 \leq P_k < \infty$$

where \mathcal{H} is a Reproducing Kernel Hilbert Space.

Assumption 9 (Normal Residuals). *Let $\inf_{f \in \mathcal{H}} \mathbf{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\mathcal{R}(f; \mathbf{x}, y)]$. The residual is defined as $\xi_i \triangleq f^*(\mathbf{x}_i) - y_i$. Then we assume for some $\sigma > 0$, it follows*

$$\xi_i \sim \mathcal{N}(0, \sigma^2)$$

2.5 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]. Therefore, we will dedicate this subsection to reviewing such works and to make clear our contributions to the iterative thresholding literature.

[BJK15] study iterative thresholding for least squares regression / sparse recovery. In particular, one part of their study is of a gradient descent algorithm when the data $\mathcal{P} = \mathcal{Q} = \mathcal{N}(\mathbf{0}, \mathbf{I})$ or multivariate sub-Gaussian with proxy \mathbf{I} . Their approximation bounds relies on the fact that $\lambda_{\min}(\mathbf{\Sigma}) = \lambda_{\max}(\mathbf{\Sigma})$ and with sufficiently large data and sufficiently small ϵ , $\lambda_{\max}(\mathbf{X})/\lambda_{\min}(\mathbf{X}) \searrow 1$. This is similar to the study by [ADKS22], where the iterative trimmed maximum likelihood estimator is studied for General Linear Models. The algorithm studied by [ADKS22] utilizes a filtering algorithm with the sketching matrix $\mathbf{\Sigma}^{-1/2}$ so the columns of \mathbf{X} are sampled from a multivariate sub-Gaussian Distribution with proxy \mathbf{I} before running the iterative thresholding procedure. This ‘whitening’ procedure to decrease the conditioning number of the covariates is also done in recent work, [SBRJ19, BJKK17].

Conditioning covariates does not generalize to kernel learning where we are given a matrix \mathbf{K} which is equivalent to inner product of the quasimatrix¹, Φ , with itself. In the infinite dimensional case, it is not possible to sketch the kernel matrix [W⁺14] in order to have the original covariates be well-conditioned. In the finite dimensional case, the feature maps can be quite large and it is very difficult to obtain in practice. Thus, we are left with Φ where the columns are sampled from a sub-Gaussian Distribution with proxy $\mathbf{\Gamma}$ is a trace-class operator, which implies the eigenvalues tend to zero, i.e. $\lambda_{\inf}(\mathbf{\Gamma}) = 0$, and there is no longer a notion of $\lambda_{\min}(\mathbf{\Gamma})$.

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

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

¹ A quasimatrix is an infinite-dimensional analogue of a tall-skinny matrix that represents an ordered set of functions in ℓ_2 (see e.g. [TT15]).

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

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

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]. Let $\text{Proj}_{\mathcal{K}}$ be the projection of a function on to the convex set $\mathcal{K} \triangleq \{f \in \mathcal{H} : \|f\|_{\mathcal{H}} \leq R\}$, then our update steps are

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

$$f^{(k+1)} = \text{Proj}_{\mathcal{K}} \left[f^{(k)} - \eta \nabla_f g(f^{(k)}, t^{(k+1)}) \right] \quad (3.2)$$

The proof of convergence for the above algorithm was given in [JNJ20][Theorem 35]. The sufficient condition for convergence is $g(f, t)$ is concave with respect to t , which for the subquantile objective is simple to show.

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 $\hat{\nu}_i \triangleq \mathcal{R}(f; \mathbf{x}_{[i]}, y_{[i]})$, it then follows $\hat{\nu}_{(1-\epsilon)N} \in \arg \max_{t \in \mathbb{R}} g(t, f)$.*

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 .

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

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

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

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

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

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

This concludes the proof. \blacksquare

We have therefore shown that the two-step optimization of Subquantile Minimization gives the iterative thresholding algorithm.

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.

4.1 Algorithm

Algorithm 1 Subquantile Minimization for One Hidden-Layer Neural Networks

input: Possibly corrupted $\mathbf{X} \in \mathbb{R}^{d \times N}$ with outputs $\mathbf{y} \in \mathbb{R}^N$ and corruption parameter $\epsilon = O(\text{poly}(C_3, \Gamma))$.

output: Approximate solution $\mathbf{W} \in \mathbb{R}^{k \times d}$ to $\min \|\mathbf{W} - \mathbf{W}^*\|_F$.

1: $\mathbf{W}^{(0)} \leftarrow \mathbf{0}$

2: **for** $t \in [T]$ **do**

3: $S^{(t)} \leftarrow \{i \in [n] : \mathcal{L}(\mathbf{W}^{(t)}; \mathbf{x}_i, y_i) \leq \mathcal{L}(\mathbf{W}^{(t)}; \mathbf{x}_{\lfloor (1-\epsilon)N \rfloor}; y_{\lfloor (1-\epsilon)N \rfloor})\}$

4: $\nabla \mathcal{R}(\mathbf{W}; S^{(t)}) \leftarrow \frac{2}{(1-\epsilon)N} \cdot \sum_{i \in S^{(t)}} (\sum_{k \in [K]} \sigma(\mathbf{x}_i^T \mathbf{w}_k) - y_i) \cdot (\sigma \circ (\mathbf{W}^{(t)} \mathbf{x}_i)) \cdot \mathbf{x}_i^T$

5: $\mathbf{W}^{(t+1)} \leftarrow \mathbf{W}^{(t)} - \eta \nabla \mathcal{R}(\mathbf{W}^{(t)}; S^{(t)})$

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

4.2 Kernelized GLMs

The error function for the the Kernelized GLM problem is given by the following equation for a single training pair $(\mathbf{x}_i, y_i) \sim \mathcal{D}$.

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

Theorem 12 (Subquantile Minimization for Generalized Linear Models is Good with High Probability). *Let Algorithm 1 be run on a dataset $\mathcal{D} \sim \hat{\mathcal{P}}$ with learning rate $\eta \triangleq \Omega(\ell^{-1})$ and link function $\omega : \mathbb{R} \rightarrow \mathbb{R}$, s.t. $C_1 \leq \omega'(x) \leq C_2$ for absolute constants $C_1, C_2 > 0$. Then after $O\left(n \log\left(\frac{\|f^*\|_{\mathcal{H}}}{\epsilon}\right)\right)$ gradient descent iterations, with probability exceeding $1 - \delta$ and a positive constant C ,*

$$\|f^{(T)} - f^*\|_{\mathcal{H}} \leq \epsilon$$

for $n \geq (1 - \epsilon)^{-1} (16\|\Gamma\|_{\text{op}}^2 + 2P_k^2 \log(2/\delta))$.

Proof. The proof is deferred to § A.1. \blacksquare

4.3 Neural Networks

In this section we will consider Iterative Thresholding for a linear one-layer neural network and then a general two-layer neural network.

4.3.1 One Layer Linear Network

We start with the simple case of a linear one-layer neural network for multivariate regression. The error function for the linear one-layer Neural Network problem is given by the following equation for $\mathbf{X} \in \mathbb{R}^{d \times n}$ and $\mathbf{Y} \in \mathbb{R}^{k \times n}$.

$$\mathcal{L}(\mathbf{W}; \mathbf{X}, \mathbf{Y}) = \|\mathbf{WX} - \mathbf{Y}\|_{\text{F}}^2$$

Theorem 13 (Subquantile Minimization for a One-layer Linear Network is Good with High Probability). *Let Algorithm 1 be run on a dataset $\mathcal{D} \sim \hat{\mathcal{P}}$ such that $\mathbf{X} \in \mathbb{R}^{d \times n}$ and $\mathbf{Y} \in \mathbb{R}^{k \times n}$ with learning rate $\eta \triangleq \Omega(\|\mathbf{\Gamma}\|^{-1})$. Then after $O\left(n \log\left(\frac{\|\mathbf{W}^*\|_{\text{F}}}{\varepsilon}\right)\right)$ gradient descent iterations, with probability exceeding $1 - \delta$ and a positive constant C ,*

$$\|\mathbf{W}^{(T)} - \mathbf{W}^*\|_{\text{F}} \leq \varepsilon + O(k\sigma) + O(k\sigma\|\mathbf{\Gamma}\|)$$

for $n = \Omega(\Xi)$.

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

We will now consider the problem of learning neurons.

4.3.2 Single Neuron

Theorem 14. *Let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$ be the data matrix and $\mathbf{y} = [y_1, \dots, y_n]^T$ be the output, such that for $i \in P$, $\mathbf{x}_i \sim \mathcal{P}$ with multivariate proxy $\mathbf{\Gamma}$, and moments, $\mathbf{E}[\mathbf{x}] = \mathbf{0}$ and $\mathbf{E}[\mathbf{x}\mathbf{x}^T] = \mathbf{\Sigma}$. Suppose the output is given as $y_i = \sigma(\mathbf{x}^T \mathbf{w}^*)$ for $\xi_i \sim \mathcal{N}(0, \sigma^2)$ and $\sigma = \max\{0, x\}$ 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.2 ■

4.3.3 Two Layer Neural Network

We now give our assumption of the data.

Definition 15 (One Hidden Layer Network). *The true data is given by the following relation for a $\mathbf{W}^* \in \mathbb{R}^{k \times d}$ and $\xi \sim \mathcal{N}(0, \sigma^2)$.*

$$y = \sum_{k=1}^K \sigma(\mathbf{x}^T \mathbf{w}_k^*) + \xi = f^*(\mathbf{x}) + \xi$$

where $\sigma : \mathbb{R} \rightarrow \mathbb{R}$.

Theorem 16. *Let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{N \times d}$ be the data matrix and $\mathbf{y} = [y_1, \dots, y_n]^T$ be the output, such that for $i \in P$, $\mathbf{x}_i \sim \mathcal{P}$ are sub-Gaussian with multivariate proxy $\mathbf{\Gamma}$ with moments, $\mathbf{E}[\mathbf{x}] = \mathbf{0}$, $\mathbf{E}[\mathbf{x}_i \mathbf{x}_i^T] = \mathbf{\Sigma}$ and $y_i = \sum_{k=1}^K \sigma(\mathbf{x}^T \mathbf{w}_k^*) + \xi_i$ for $\xi_i \sim \mathcal{N}(0, \sigma^2)$. Then after $O\left(C_1^{-2} C_2^2 \kappa(\mathbf{\Gamma}) \log\left(\frac{\|\mathbf{W}^*\|_{\text{F}}}{\varepsilon}\right)\right)$ gradient descent iterations, with probability exceeding $1 - \delta$ and a positive constant C , Algorithm 1 with learning rate $\eta = O(C_2^{-2} \|\mathbf{\Gamma}\|^{-1})$ returns $\mathbf{W}^{(T)}$ such that $\|\mathbf{W}^{(T)} - \mathbf{W}^*\|_{\text{F}} \leq \varepsilon$.*

Proof. The proof is deferred to § B.3. ■

5 Discussion

The main contribution of this paper is the study of a nonconvex-concave formulation of Subquantile minimization for the robust learning problem for kernel ridge regression and kernel classification. We present an algorithm to solve the nonconvex-concave formulation and prove rigorous error bounds which show that the more good data that is given decreases the error bounds.

Extension to Infinite Dimensional Kernels.

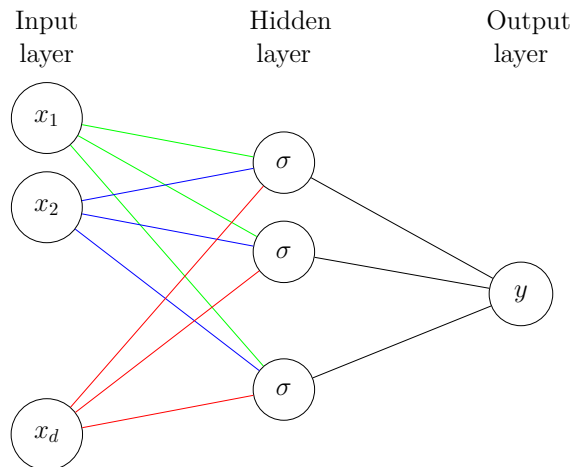


Figure 1: One Hidden Layer Neural Network.

Theory. We develop strong theoretical bounds on the normed difference between the function returned by Subquantile Minimization and the optimal function for data in the target distribution, \mathcal{P} , in the sub-Gaussian Design. We are able to show if the number of inliers is sufficiently small, then the kernelized binary classification problem with binary cross-entropy loss is consistent.

Future Work. The analysis of Subquantile Minimization can be extended to neural networks as kernel learning can be seen as a one-layer network. This generalization will appear in subsequent work. Another interesting direction work in optimization is for accelerated methods for optimizing non-convex concave min-max problems with a maximization oracle. The current theory analyzes standard gradient descent for the minimization. Ideas such as Momentum and Nesterov Acceleration in conjunction with the maximum oracle are interesting and can be analyzed in future work.

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. [1](#), [1.1](#), [2.3](#), [2.5](#)
- [AMMIL12] Yaser S Abu-Mostafa, Malik Magdon-Ismail, and Hsuan-Tien Lin. *Learning from data*, volume 4. AMLBook New York, 2012. [1](#)
- [B⁺15] Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015. [25](#)
- [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. [3](#), [1.1](#), [2.5](#)
- [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. [1](#), [2.5](#)
- [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. [1](#)
- [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. [1](#)
- [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. [1](#), [2.3](#)
- [CLRS22] Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, and Clifford Stein. *Introduction to algorithms*. MIT press, 2022. [29](#)
- [Dic16] Lee H Dicker. Ridge regression and asymptotic minimax estimation over spheres of growing dimension. 2016. [1](#)
- [DK23] Ilias Diakonikolas and Daniel M Kane. *Algorithmic high-dimensional robust statistics*. Cambridge University Press, 2023. [1](#)
- [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. [1](#)
- [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. [1](#)
- [FWZ18] Jianqing Fan, Weichen Wang, and Yiqiao Zhong. An l eigenvector perturbation bound and its application to robust covariance estimation. *Journal of Machine Learning Research*, 18(207):1–42, 2018. [1](#)
- [HR09] Peter J. Huber and Elvezio. Ronchetti. *Robust statistics*. Wiley series in probability and statistics. Wiley, Hoboken, N.J., 2nd ed. edition, 2009. [1](#)
- [HSS08] Thomas Hofmann, Bernhard Schölkopf, and Alexander J. Smola. Kernel methods in machine learning. *The Annals of Statistics*, 36(3):1171 – 1220, 2008. [1](#)
- [HYW⁺23] Shu Hu, Zhenhuan Yang, Xin Wang, Yiming Ying, and Siwei Lyu. Outlier robust adversarial training. *arXiv preprint arXiv:2309.05145*, 2023. [2.5](#)
- [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. [1](#)
- [Jen06] Johan Ludwig William Valdemar Jensen. Sur les fonctions convexes et les inégalités entre les valeurs moyennes. *Acta mathematica*, 30(1):175–193, 1906. [19](#)
- [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. [2.3](#)
- [JNJ20] Chi Jin, Praneeth Netrapalli, and Michael Jordan. What is local optimality in nonconvex-nonconcave minimax optimization? 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 4880–4889. PMLR, 13–18 Jul 2020. [3](#)

- [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. [1](#)
- [KLA18] Ashish Khetan, Zachary C. Lipton, and Anima Anandkumar. Learning from noisy singly-labeled data. In *International Conference on Learning Representations*, 2018. [1](#)
- [LBSS21] Tian Li, Ahmad Beirami, Maziar Sanjabi, and Virginia Smith. Tilted empirical risk minimization. In *International Conference on Learning Representations*, 2021. [1](#), [1](#)
- [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. [2.5](#)
- [LM00] Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of statistics*, pages 1302–1338, 2000. [17](#)
- [M⁺89] Colin McDiarmid et al. On the method of bounded differences. *Surveys in combinatorics*, 141(1):148–188, 1989. [20](#)
- [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. [1](#), [2.5](#)
- [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. [1](#)
- [PBR19] Adarsh Prasad, Sivaraman Balakrishnan, and Pradeep Ravikumar. A unified approach to robust mean estimation. *arXiv preprint arXiv:1907.00927*, 2019. [1](#)
- [PP⁺08] Kaare Brandt Petersen, Michael Syskind Pedersen, et al. The matrix cookbook. *Technical University of Denmark*, 7(15):510, 2008. [B.1](#)
- [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. [1](#)
- [RaR02] Raymond A Ryan and R a Ryan. *Introduction to tensor products of Banach spaces*, volume 73. Springer, 2002. [2.2](#)
- [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. [3](#)
- [SBRJ19] Arun Sai Suggala, Kush Bhatia, Pradeep Ravikumar, and Prateek Jain. Adaptive hard thresholding for near-optimal consistent robust regression. 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 2892–2897. PMLR, 25–28 Jun 2019. [2.5](#)
- [SS19] Yanyao Shen and Sujay Sanghavi. Learning with bad training data via iterative trimmed loss minimization. In *International Conference on Machine Learning*, pages 5739–5748. PMLR, 2019. [\(document\)](#)
- [TSM⁺17] Ilya Tolstikhin, Bharath K Sriperumbudur, Krikamol Mu, et al. Minimax estimation of kernel mean embeddings. *Journal of Machine Learning Research*, 18(86):1–47, 2017. [21](#), [C.2](#)
- [TT15] Alex Townsend and Lloyd N Trefethen. Continuous analogues of matrix factorizations. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 471(2173):20140585, 2015. [1](#)

- [W⁺14] David P Woodruff et al. Sketching as a tool for numerical linear algebra. *Foundations and Trends® in Theoretical Computer Science*, 10(1–2):1–157, 2014. [1](#)
- [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. [26](#)

A Proofs for Kernelized GLMs

In this section, we will prove error bounds for Subquantile Minimization in the Kernelized GLM Problem.

A.1 Proof of Theorem 12

Proof. From Algorithm 1, we have the gradient update for the generalized linear model.

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

We will now use the standard arguments for convergence of gradient descent methods,

$$\begin{aligned} \|f^{(t+1)} - f^*\|_{\mathcal{H}} &= \|f^{(t)} - \eta \nabla_f \mathcal{R}_{S^{(t)}}(f^{(t)}) - f^*\|_{\mathcal{H}} \\ &= \|f^{(t)} - f^* - \eta \nabla \mathcal{R}_{\text{TP}}(f^{(t)}) - \eta \nabla \mathcal{R}_{\text{FP}}(f^{(t)})\|_{\mathcal{H}} \\ &\leq \|f^{(t)} - f^* - \eta \nabla \mathcal{R}_{\text{TP}}(f^{(t)})\|_{\mathcal{H}} + \|\eta \nabla \mathcal{R}_{\text{FP}}(f^{(t)})\|_{\mathcal{H}} \end{aligned} \quad (\text{A.1})$$

We will now analyze the first term of Equation (A.1) through its square,

$$\|f^{(t)} - f^* - \eta \nabla \mathcal{R}_{\text{TP}}(f^{(t)})\|_{\mathcal{H}}^2 = \|f^{(t)} - f^*\|_{\mathcal{H}}^2 - 2\eta \cdot \langle f^{(t)} - f^*, \nabla \mathcal{R}_{\text{TP}}(f^{(t)}) \rangle_{\mathcal{H}} + \eta^2 \cdot \|\nabla \mathcal{R}_{\text{TP}}(f^{(t)})\|_{\mathcal{H}}^2 \quad (\text{A.2})$$

We then have,

$$\begin{aligned} 2\eta \cdot \langle f^{(t)} - f^*, \nabla \mathcal{R}_{\text{TP}}(f^{(t)}) \rangle_{\mathcal{H}} &\stackrel{\text{def}}{=} \frac{4\eta}{(1-\epsilon)N} \cdot \langle f^{(t)} - f^*, \sum_{i \in \text{TP}} ((\omega(f^{(t)}(\mathbf{x}_i)) - \omega(f^*(\mathbf{x}_i)) - \xi_i) \cdot \omega'(f^{(t)}(\mathbf{x}_i)) \cdot X_i) \rangle_{\mathcal{H}} \\ &\geq \frac{4C_1\eta}{(1-\epsilon)N} \cdot \langle f^{(t)} - f^*, \sum_{i \in \text{TP}} ((\omega(f^{(t)}(\mathbf{x}_i)) - \omega(f^*(\mathbf{x}_i))) \cdot X_i) \rangle_{\mathcal{H}} - \frac{4\eta}{(1-\epsilon)N} \cdot \langle f^{(t)} - f^*, \sum_{i \in \text{TP}} \xi_i \omega'(f^{(t)}(\mathbf{x}_i)) \cdot X_i \rangle_{\mathcal{H}} \end{aligned} \quad (\text{A.3})$$

In the above, in the first inequality we can note with the reproducing property that as the link function is monotonic increasing, we have $(f^{(t)}(\mathbf{x}) - f^*(\mathbf{x}))(\omega(f^{(t)}(\mathbf{x})) - \omega(f^*(\mathbf{x}))) \geq 0$ for any $\mathbf{x} \in \mathcal{X}$. We will first lower bound the first term in Equation (A.3). Then, consider the function, $h : \mathcal{H} \rightarrow \mathbb{R}$, for a Dirac measure $P(X) = \frac{1}{|\text{TP}|} \delta_{\text{TP}}(X)$.

$$h(f) \triangleq \int_{\mathcal{H}} \left(\int \omega(y) dy \right) (f(X)) dP(X)$$

The first derivative gives us the following,

$$\nabla h(f) \triangleq \int_{\mathcal{H}} \omega(f(X)) \cdot X dP(X)$$

From which we have the following,

$$\begin{aligned} \nabla^2 h(f) &\triangleq \int_{\mathcal{H}} \omega'(f(X)) \cdot X \otimes X dP(X) \succeq \min_{y \in \mathbb{R}} \omega'(y) \cdot \lambda_{\min} \left(\int_{\mathcal{H}} X \otimes X dP(X) \right) \cdot \mathbf{I} \\ &\stackrel{\text{def}}{=} C_1 \lambda_{\min} \left(\int_{\mathcal{H}} X \otimes X dP(X) \right) \cdot \mathbf{I} \triangleq \alpha \cdot \mathbf{I} \end{aligned}$$

Finally, we can note that

$$\nabla^2 h(f) \preceq \max_{y \in \mathbb{R}} \omega'(y) \cdot \lambda_{\max} \left(\int_{\mathcal{H}} X \otimes X dP(X) \right) \cdot \mathbf{I} \stackrel{\text{def}}{=} C_2 \cdot \lambda_{\max} \left(\int_{\mathcal{H}} X \otimes X dP(X) \right) \cdot \mathbf{I} \triangleq \beta \cdot \mathbf{I}$$

Then, with Lemma 25, we have

$$2\eta \cdot \langle f^{(t)} - f^*, \nabla \mathcal{R}_{\text{TP}}(f^{(t)}) \rangle_{\mathcal{H}}$$

$$\begin{aligned}
&\geq \frac{2C_1\eta}{(1-\epsilon)N} \cdot \lambda_{\min}\left(\sum_{i \in \text{TP}} X_i \otimes X_i\right) \|f^{(t)} - f^*\|_{\mathcal{H}}^2 + \frac{2C_2^{-1}\eta}{(1-\epsilon)N} \lambda_{\max}^{-1}\left(\sum_{i \in \text{TP}} X_i \otimes X_i\right) \left\| \sum_{i \in \text{TP}} (\omega(f^{(t)}(\mathbf{x}_i)) - \omega(f^*(\mathbf{x}_i))) \cdot X_i \right\|_{\mathcal{H}}^2 \\
&\geq \frac{2C_1\eta}{(1-\epsilon)N} \cdot \lambda_{\min}\left(\sum_{i \in \text{TP}} X_i \otimes X_i\right) \|f^{(t)} - f^*\|_{\mathcal{H}}^2 + \frac{2C_2^{-1}\eta}{(1-\epsilon)N} \cdot \kappa^{-1}\left(\sum_{i \in \text{TP}} X_i \otimes X_i\right) \cdot \sum_{i \in \text{TP}} (\omega(f^{(t)}(\mathbf{x}_i)) - \omega(f^*(\mathbf{x}_i)))^2
\end{aligned} \tag{A.4}$$

We will now upper bound the second term in Equation (A.3).

$$\begin{aligned}
&\frac{4\eta}{(1-\epsilon)N} \cdot \langle f^{(t)} - f^*, \sum_{i \in \text{TP}} \xi_i \omega'(f^{(t)}(\mathbf{x}_i)) \cdot X_i \rangle_{\mathcal{H}} \leq \frac{4\eta}{(1-\epsilon)N} \cdot \|f^{(t)} - f^*\|_{\mathcal{H}} \left\| \sum_{i \in \text{TP}} \xi_i \omega'(f^{(t)}(\mathbf{x}_i)) \cdot X_i \right\|_{\mathcal{H}} \\
&\leq \frac{C_1\eta}{(1-\epsilon)N} \cdot \lambda_{\min}\left(\sum_{i \in \text{TP}} X_i \otimes X_i\right) \|f^{(t)} - f^*\|_{\mathcal{H}}^2 + \frac{4C_1^{-1}\eta}{(1-\epsilon)N} \cdot \lambda_{\min}^{-1}\left(\sum_{i \in \text{TP}} X_i \otimes X_i\right) \left\| \sum_{i \in \text{TP}} \xi_i \omega'(f^{(t)}(\mathbf{x}_i)) \cdot X_i \right\|_{\mathcal{H}}^2 \\
&\leq \frac{C_1\eta}{(1-\epsilon)N} \cdot \lambda_{\min}\left(\sum_{i \in \text{TP}} X_i \otimes X_i\right) \|f^{(t)} - f^*\|_{\mathcal{H}}^2 + \frac{4C_1^{-1}C_2^2\eta}{(1-\epsilon)N} \cdot \kappa\left(\sum_{i \in \text{TP}} X_i \otimes X_i\right) \|\xi_{\text{TP}}\|_2^2
\end{aligned} \tag{A.5}$$

In the above, in the first inequality we applied the Cauchy-Schwarz Inequality, in the second inequality we utilize Young's Inequality (see Proposition 26), in the third inequality we applied Lemma 23. Then for the third term of Equation (A.2), we have

$$\begin{aligned}
&\|\eta \nabla \mathcal{R}_{\text{TP}}(f^{(t)})\|_{\mathcal{H}}^2 \stackrel{\text{def}}{=} \frac{\eta^2}{[(1-\epsilon)N]^2} \cdot \left\| \sum_{i \in \text{TP}} (\omega(f^{(t)}(\mathbf{x}_i)) - y_i) \cdot \omega'(f^{(t)}(\mathbf{x}_i)) \cdot X_i \right\|_{\mathcal{H}}^2 \\
&= \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot \left\| \sum_{i \in \text{TP}} (\omega(f^{(t)}(\mathbf{x}_i)) - \omega(f^*(\mathbf{x}_i)) + \xi_i) \cdot \omega'(f^{(t)}(\mathbf{x}_i)) \cdot X_i \right\|_{\mathcal{H}}^2 \\
&\leq \frac{8\eta^2}{[(1-\epsilon)N]^2} \cdot \left(\left\| \sum_{i \in \text{TP}} (\omega(f^{(t)}(\mathbf{x}_i)) - \omega(f^*(\mathbf{x}_i)) \cdot \omega'(f^{(t)}(\mathbf{x}_i)) \cdot X_i \right\|_{\mathcal{H}}^2 + \left\| \sum_{i \in \text{TP}} \xi_i \omega'(f^{(t)}(\mathbf{x}_i)) \cdot X_i \right\|_{\mathcal{H}}^2 \right) \\
&\leq \frac{8C_2^2\eta^2}{[(1-\epsilon)N]^2} \cdot \left\| \sum_{i \in \text{TP}} X_i \otimes X_i \right\|_{\text{op}} \cdot \sum_{i \in \text{TP}} (\omega(f^{(t)}(\mathbf{x}_i)) - \omega(f^*(\mathbf{x}_i)))^2 + \frac{8C_2^2\eta^2}{[(1-\epsilon)N]^2} \cdot \|\xi_{\text{TP}}\|_2^2 \left\| \sum_{i \in \text{TP}} X_i \otimes X_i \right\|_{\text{op}}
\end{aligned} \tag{A.6}$$

In the above, in the first inequality we utilize the elementary inequality $(a+b)^2 \leq 2a^2 + 2b^2$, in the second inequality we utilize Lemma 23. We then see the second term of Equation (A.4) is greater than the first term of Equation (A.6) when $\eta \leq (1-\epsilon)N(2C_2^3\kappa(\Phi_{\text{TP}}\Phi_{\text{TP}}^*)\|\Phi_{\text{TP}}\Phi_{\text{TP}}^*\|_{\text{op}})^{-1}$. We now will upper bound second term of Equation (A.1) through its square.

$$\begin{aligned}
&\|\eta \nabla \mathcal{R}_{\text{FP}}(f^{(t)})\|_{\mathcal{H}}^2 \stackrel{\text{def}}{=} \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot \left\| \sum_{i \in \text{FP}} (\omega(f^{(t)}(\mathbf{x}_i)) - y_i) \cdot \omega'(f^{(t)}(\mathbf{x}_i)) \cdot X_i \right\|_{\mathcal{H}}^2 \\
&\leq \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot \left\| \sum_{i \in \text{FP}} X_i \otimes X_i \right\|_{\text{op}} \sum_{i \in \text{FP}} [(\omega(f^{(t)}(\mathbf{x}_i)) - y_i) \cdot \omega'(f^{(t)}(\mathbf{x}_i))]^2 \\
&\leq \frac{4C_2^2\eta^2}{[(1-\epsilon)N]^2} \cdot \left\| \sum_{i \in \text{FP}} X_i \otimes X_i \right\|_{\text{op}} \sum_{i \in \text{FN}} (\omega(f^{(t)}(\mathbf{x}_i)) - y_i)^2 \\
&= \frac{4C_2^2\eta^2}{[(1-\epsilon)N]^2} \cdot \left\| \sum_{i \in \text{FP}} X_i \otimes X_i \right\|_{\text{op}} \sum_{i \in \text{FN}} (\omega(f^{(t)}(\mathbf{x}_i)) - \omega(f^*(\mathbf{x}_i)) + \xi_i)^2 \\
&\leq \frac{4C_2^4\eta^2}{[(1-\epsilon)N]^2} \cdot \left\| \sum_{i \in \text{FP}} X_i \otimes X_i \right\|_{\text{op}} \left(\sum_{i \in \text{FN}} (f^{(t)}(\mathbf{x}_i) - f^*(\mathbf{x}_i))^2 + \xi_i^2 \right) \\
&\leq \frac{8C_2^4\eta^2}{[(1-\epsilon)N]^2} \cdot \left\| \sum_{i \in \text{FP}} X_i \otimes X_i \right\|_{\text{op}} \left(\left\| \sum_{i \in \text{FN}} X_i \otimes X_i \right\|_{\text{op}} \|f^{(t)} - f^*\|_{\mathcal{H}}^2 + \|\xi_{\text{FN}}\|^2 \right)
\end{aligned} \tag{A.7}$$

In the above, the first inequality follows from Lemma 24, the second inequality follows from the optimality of the Subquantile set and the bounded link function gradient, the third inequality follows from noting for any $x, y \in \mathbb{R}$, the bounded gradient implies Lipschitzness, i.e. $|\omega(x) - \omega(y)| \leq C_2|x - y|$. The final equality follows from the following,

$$\begin{aligned} \sum_{i \in \text{FN}} (f^{(t)}(\mathbf{x}_i) - f^*(\mathbf{x}_i))^2 &= \langle f^{(t)} - f^*, [\sum_{i \in \text{FN}} X_i \otimes X_i] (f^{(t)} - f^*) \rangle_{\mathcal{H}} \stackrel{\text{def}}{=} \|f^{(t)} - f^*\|_{\Sigma_{\text{FN}, \mathcal{H}}}^2 \\ &\leq \|f^{(t)} - f^*\|_{\mathcal{H}}^2 \left\| \sum_{i \in \text{FN}} X_i \otimes X_i \right\|_{\text{op}} \end{aligned}$$

In the above, the first equality follows from the reproducing property of the RKHS. Then from Equations (A.4), (A.5), and (A.7), we have

$$\begin{aligned} \|f^{(t+1)} - f^*\|_{\mathcal{H}}^2 &\leq \|f^{(t)} - f^*\|_{\mathcal{H}} \left(1 - \frac{C_1 \eta}{(1 - \epsilon)N} \cdot \lambda_{\min} \left(\sum_{i \in \text{TP}} X_i \otimes X_i \right) + \frac{8C_2^4 C_3 \eta^2}{(1 - \epsilon)N} \cdot \left\| \sum_{i \in \text{FN}} X_i \otimes X_i \right\|_{\text{op}} \right) \\ &\quad + \frac{8C_2^4 \eta^2}{[(1 - \epsilon)N]^2} \cdot \left\| \sum_{i \in \text{FP}} X_i \otimes X_i \right\|_{\text{op}} \|\xi_{\text{FN}}\|_2^2 + \frac{8C_2^2 \eta^2}{[(1 - \epsilon)N]^2} \cdot \left\| \sum_{i \in \text{TP}} X_i \otimes X_i \right\|_{\text{op}} \|\xi_{\text{TP}}\|_2^2 \\ &\quad + \frac{4C_1^{-1} C_2^2 \eta}{(1 - \epsilon)N} \cdot \kappa \left(\sum_{i \in \text{TP}} X_i \otimes X_i \right) \|\xi_{\text{TP}}\|_2^2 \end{aligned}$$

It is clear there exists ϵ such that the multiplicative term is on the order of,

$$1 - \Omega \left(\frac{C_1 \eta}{(1 - \epsilon)N} \cdot \lambda_{\min} \left(\sum_{i \in \text{TP}} X_i \otimes X_i \right) \right) = 1 - \Omega \left(C_1 C_2^{-3} \kappa^{-2} \left(\sum_{i \in \text{TP}} X_i \otimes X_i \right) \right)$$

Solving the induction, we obtain ... ■

B Proofs for Neural Networks

B.1 Proof of Theorem 13

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

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

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

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

Our proof will begin similarly to the proof of Theorem 12. We then have,

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

We first will upper bound the first term in Equation (B.1).

$$\begin{aligned} &\|\mathbf{W}^{(t)} - \mathbf{W}^* - \eta \nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{TP})\|_{\text{F}}^2 \\ &= \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_{\text{F}}^2 - \eta \cdot \text{Tr}((\mathbf{W}^{(t)} - \mathbf{W}^*)^T (\nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{TP}))) + \|\eta \nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{TP})\|_{\text{F}}^2 \end{aligned} \quad (\text{B.2})$$

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

$$\begin{aligned}
& 2\eta \cdot \text{Tr}((\mathbf{W}^{(t)} - \mathbf{W}^*)^T (\nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{TP}))) \stackrel{\text{def}}{=} \frac{4\eta}{(1-\epsilon)N} \cdot \text{Tr}((\mathbf{W}^{(t)} - \mathbf{W}^*)^T (\mathbf{W}^{(t)} \mathbf{X}_{\text{TP}} - \mathbf{Y}_{\text{TP}}) \mathbf{X}_{\text{TP}}^T) \\
&= \frac{4\eta}{(1-\epsilon)N} \cdot \text{Tr}((\mathbf{W}^{(t)} - \mathbf{W}^*)^T (\mathbf{W}^{(t)} - \mathbf{W}^*) \mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T) - \frac{4\eta}{(1-\epsilon)N} \cdot \text{Tr}((\mathbf{W}^{(t)} - \mathbf{W}^*)^T \mathbf{E}_{\text{TP}} \mathbf{X}_{\text{TP}}^T) \quad (\text{B.3})
\end{aligned}$$

We will first lower bound the first term of Equation B.3.

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

In the above, in the first equality we use the cyclic property of the trace, in the first inequality we apply the Cauchy-Schwarz Inequality, in the first inequality we apply Lemma 25 to the first term, in the final inequality we apply Young's inequality (see Proposition 26), and in the third equality we have \otimes as the Kronecker product which gives the following equality,

$$(\mathbf{I} \otimes \mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T) \text{vec}((\mathbf{W}^{(t)} - \mathbf{W}^*)^T) = \begin{bmatrix} \mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T & & \\ & \ddots & \\ & & \mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T \end{bmatrix} \begin{bmatrix} \mathbf{w}_1 - \mathbf{w}_1^* \\ \vdots \\ \mathbf{w}_K - \mathbf{w}_K^* \end{bmatrix} = \sum_{k \in [K]} \mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T (\mathbf{w}_k - \mathbf{w}_k^*)$$

We now upper bound the second term in Equation (B.3).

$$\begin{aligned}
& \frac{4\eta}{(1-\epsilon)N} \cdot \text{Tr}((\mathbf{W}^{(t)} - \mathbf{W}^*)^T \mathbf{E}_{\text{TP}} \mathbf{X}_{\text{TP}}^T) \leq \frac{4\eta}{(1-\epsilon)N} \cdot \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F \|\mathbf{E}_{\text{TP}} \mathbf{X}_{\text{TP}}^T\|_F \\
&\leq \frac{2\eta}{(1-\epsilon)N} \cdot \lambda_{\min}(\mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T) \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F^2 + \frac{2\eta}{(1-\epsilon)N} \cdot \kappa(\mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T) \|\mathbf{E}_{\text{TP}}\|_F^2
\end{aligned}$$

In the above, in the first inequality we apply the Cauchy-Schwarz Inequality and in the second inequality we use Young's Inequality (see Proposition 26). We now upper bound the second term in Equation (B.1),

$$\begin{aligned}
& \|\eta \nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{FP})\|_F \stackrel{\text{def}}{=} \frac{2\eta}{(1-\epsilon)N} \cdot \|(\mathbf{W}^{(t)} \mathbf{X}_{\text{FP}} - \mathbf{Y}_{\text{FP}}) \mathbf{X}_{\text{FP}}^T\|_F \\
&\leq \frac{2\eta}{(1-\epsilon)N} \cdot \|\mathbf{X}_{\text{FP}}\|_2 \|\mathbf{W}^{(t)} \mathbf{X}_{\text{FP}} - \mathbf{Y}_{\text{FP}}\|_F \\
&\leq \frac{2\eta}{(1-\epsilon)N} \cdot \|\mathbf{X}_{\text{FP}}\|_2 \|\mathbf{W}^{(t)} \mathbf{X}_{\text{FN}} - \mathbf{Y}_{\text{FN}}\|_F \\
&\leq \frac{2\eta}{(1-\epsilon)N} \cdot \|\mathbf{X}_{\text{FP}}\|_2 (\|\mathbf{W}^{(t)} \mathbf{X}_{\text{FN}} - \mathbf{W}^* \mathbf{X}_{\text{FN}}\|_F^2 + \|\mathbf{E}_{\text{FN}}\|_F^2) \\
&\leq \frac{2\eta}{(1-\epsilon)N} \cdot \|\mathbf{X}_{\text{FP}}\|_2 \|\mathbf{X}_{\text{FN}}\|_2 \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F + \frac{2\eta}{(1-\epsilon)N} \cdot \|\mathbf{X}_{\text{FP}}\|_2 \|\mathbf{E}_{\text{FN}}\|_F
\end{aligned}$$

In the above, the first and fourth inequalities from the fact that for any two size compatible matrices, \mathbf{A}, \mathbf{B} , it holds that $\|\mathbf{AB}\|_F \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_2$, the second inequality follows from the optimality of the Subquantile set,

the third inequality follows from the sub-additivity of the Frobenius norm. We will now upper bound the third term in Equation (B.2),

$$\begin{aligned}\|\eta \nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{TP})\|_F &\stackrel{\text{def}}{=} \frac{2\eta}{(1-\epsilon)N} \cdot \|(\mathbf{W}^{(t)} \mathbf{X}_{\text{TP}} - \mathbf{Y}_{\text{FP}}) \mathbf{X}_{\text{TP}}^T\|_F \\ &\leq \frac{2\eta}{(1-\epsilon)N} \cdot \|(\mathbf{W}^{(t)} - \mathbf{W}^*) \mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T\|_F + \frac{2\eta}{(1-\epsilon)N} \cdot \|\mathbf{X}_{\text{TP}}\|_2 \|\mathbf{E}_{\text{TP}}\|_F\end{aligned}\quad (\text{B.5})$$

In the above, we use the elementary inequality $(a+b)^2 \leq 2a^2 + 2b^2$. We then have from choosing $\eta \leq (1-\epsilon)N(4\|\mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T\|_2^{-1})$, the first term in Equation (B.5) will be less than the second term in Equation (B.4). We thus obtain from noting that $\sqrt{1-2x} \leq 1-x$ for any $x \leq 1/2$,

$$\begin{aligned}\|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_F &\leq \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F \left(1 - \frac{2\eta}{(1-\epsilon)N} \cdot \lambda_{\min}(\mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T) + \frac{2\eta}{(1-\epsilon)N} \cdot \|\mathbf{X}_{\text{FP}}\|_2 \|\mathbf{X}_{\text{FN}}\|_2\right) \\ &\quad + \frac{2\eta}{(1-\epsilon)N} \cdot \|\mathbf{E}_{\text{FN}}\|_F \|\mathbf{X}_{\text{FP}}\|_2 + \sqrt{\frac{2\eta}{(1-\epsilon)N}} \cdot \kappa(\mathbf{X}_{\text{TP}}) \|\mathbf{E}_{\text{TP}}\|_F\end{aligned}$$

From Proposition 18, we obtain with probability exceeding $1-\delta$ that $\|\mathbf{E}_{\text{FN}}\|_F \leq \sqrt{30NK\epsilon \log \epsilon^{-1}}$ and from noting that for $\epsilon < 0.5$ that $(1-\epsilon) \log((1-\epsilon)^{-1}) \leq \epsilon \log \epsilon^{-1}$ we also have $\|\mathbf{E}_{\text{TP}}\|_F \leq \sqrt{30NK\epsilon \log \epsilon^{-1}}$ with the same probability guarantees. ■

B.2 Proof of Theorem 14

Proof. The neuron function is given as follows,

$$f_{\mathbf{w}}(\mathbf{x}) = \sigma(\mathbf{x}^T \mathbf{w}) = \mathbf{x}^T \mathbf{w} \cdot \mathbb{I}\{\mathbf{x}^T \mathbf{w} \geq 0\}$$

The empirical risk function is given as follows,

$$\mathcal{R}(\mathbf{w}; \mathbf{S}) = \frac{1}{(1-\epsilon)N} \cdot \sum_{i \in \mathbf{S}} (\sigma(\mathbf{x}_i^T \mathbf{w}) - y_i)^2$$

The gradient is then given as follows,

$$\frac{\partial \mathcal{R}(\mathbf{w}; \mathbf{S})}{\partial \mathbf{w}} = \frac{2}{(1-\epsilon)N} \cdot \sum_{i \in \mathbf{S}} (\sigma(\mathbf{x}_i^T \mathbf{w}) - y_i) \cdot \mathbf{x}_i \cdot \mathbb{I}\{\mathbf{x}_i^T \mathbf{w} \geq 0\}$$

Then, from noting that the ReLU function is non-differentiable at one point, then almost surely the Hessian is given as

$$\frac{\partial^2 \mathcal{R}(\mathbf{w}; \mathbf{S})}{\partial \mathbf{w}^2} = \frac{2}{(1-\epsilon)N} \cdot \sum_{i \in \mathbf{S}} \mathbf{x}_i \mathbf{x}_i^T \cdot \mathbb{I}\{\mathbf{x}_i^T \mathbf{w} \geq 0\}$$

We will now begin our standard analysis.

$$\begin{aligned}\|\mathbf{w}^{(t+1)} - \mathbf{w}^*\|_2 &= \|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}; \mathbf{S}^{(t)})\|_2 \\ &= \|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) + \eta \nabla \mathcal{R}(\mathbf{w}^*; \text{TP}) - \eta \nabla \mathcal{R}(\mathbf{w}^*; \text{TP}) - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{FP})\|_2 \\ &\leq \|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) + \eta \nabla \mathcal{R}(\mathbf{w}^*; \text{TP})\|_2 + \|\eta \mathcal{R}(\mathbf{w}^{(t)}; \text{TP})\|_2 + \|\eta \mathcal{R}(\mathbf{w}^{(t)}; \text{FP})\|_2\end{aligned}\quad (\text{B.6})$$

We will now upper bound the first term of Equation B.6 through its square,

$$\begin{aligned}\|\mathbf{w}^{(t)} - \mathbf{w}^* - \eta \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) + \eta \nabla \mathcal{R}(\mathbf{w}^*; \text{TP})\|_2^2 &= \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2^2 \\ &\quad - 2\eta \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) - \nabla \mathcal{R}(\mathbf{w}^*; \text{TP}) \rangle + \eta^2 \cdot \|\nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) - \nabla \mathcal{R}(\mathbf{w}^*; \text{TP})\|_2^2\end{aligned}\quad (\text{B.7})$$

We will lower bound the second term of Equation B.7.

$$2\eta \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \nabla \mathcal{R}(\mathbf{w}^{(t)}; \text{TP}) - \nabla \mathcal{R}(\mathbf{w}^*; \text{TP}) \rangle$$

$$\begin{aligned}
&\stackrel{\text{def}}{=} \frac{4\eta}{(1-\epsilon)N} \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \sum_{i \in \text{TP}} (\sigma(\mathbf{x}_i^T \mathbf{w}^{(t)}) - y_i) \cdot \mathbf{x}_i \cdot \mathbb{I}\{\mathbf{x}_i^T \mathbf{w}^{(t)} \geq 0\} - \xi_i \mathbf{x}_i \cdot \mathbb{I}\{\mathbf{x}_i^T \mathbf{w}^* \geq 0\} \rangle \\
&= \frac{4\eta}{(1-\epsilon)N} \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \sum_{i \in \text{TP}} (\sigma(\mathbf{x}_i^T \mathbf{w}^{(t)}) - \sigma(\mathbf{x}_i^T \mathbf{w}^*)) \cdot \mathbf{x}_i \cdot \mathbb{I}\{\mathbf{x}_i^T \mathbf{w}^{(t)} \geq 0\} \rangle \\
&= \frac{4\eta}{(1-\epsilon)N} \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \sum_{i \in \text{TP}} \mathbf{x}_i \mathbf{x}_i^T (\mathbf{w}^{(t)} - \mathbf{w}^*) \cdot \mathbb{I}\{\mathbf{x}_i^T \mathbf{w}^{(t)} \geq 0\} \cdot \mathbb{I}\{\mathbf{x}_i^T \mathbf{w}^* \geq 0\} \rangle \\
&\quad + \frac{4\eta}{(1-\epsilon)N} \cdot \langle \mathbf{w}^{(t)} - \mathbf{w}^*, \sum_{i \in \text{TP}} \mathbf{x}_i \mathbf{x}_i^T \mathbf{w}^{(t)} \cdot \mathbb{I}\{\mathbf{x}_i^T \mathbf{w}^{(t)} \geq 0\} \cdot \mathbb{I}\{\mathbf{x}_i^T \mathbf{w}^* < 0\} \rangle
\end{aligned}$$

■

B.3 Proof of Theorem 16

Proof. Recall the neural network function is given as follows,

$$f_{\mathbf{W}}(\mathbf{x}) = \sum_{k \in [K]} \sigma(\mathbf{x}^T \mathbf{w}_k)$$

Then we have the empirical risk as,

$$\mathcal{R}(\mathbf{W}; \mathbf{S}^{(t)}) = \frac{1}{(1-\epsilon)N} \cdot \sum_{i \in \mathbf{S}^{(t)}} \left(\sum_{k \in [K]} \sigma(\mathbf{x}_i^T \mathbf{w}_k) - y_i \right)^2$$

The gradient is then given as follows for a column \mathbf{w}_k of \mathbf{W} for $k \in [K]$,

$$\frac{\partial \mathcal{R}(\mathbf{W}; \mathbf{S}^{(t)})}{\partial \mathbf{w}_k} = \frac{2}{(1-\epsilon)N} \cdot \sum_{i \in \mathbf{S}^{(t)}} \sum_{j \in [K]} (\sigma(\mathbf{x}_i^T \mathbf{w}_j) - y_i) \cdot \sigma'(\mathbf{x}_i^T \mathbf{w}_k) \cdot \mathbf{x}_i$$

We can then compute the second partial derivatives, first for $k \neq \ell$ and $k, \ell \in [K]$, we have

$$\frac{\partial^2 \mathcal{R}(\mathbf{W}; \mathbf{S}^{(t)})}{\partial \mathbf{w}_k \partial \mathbf{w}_\ell} = \frac{2}{(1-\epsilon)N} \cdot \sum_{i \in \mathbf{S}^{(t)}} \sigma'(\mathbf{x}_i^T \mathbf{w}_k) \sigma'(\mathbf{x}_i^T \mathbf{w}_\ell) \cdot \mathbf{x}_i \mathbf{x}_i^T$$

then on the diagonal, we have

$$\frac{\partial^2 \mathcal{R}(\mathbf{W}; \mathbf{S}^{(t)})}{\partial \mathbf{w}_k^2} = \frac{2}{(1-\epsilon)N} \cdot \left(\sum_{i \in \mathbf{S}^{(t)}} \sum_{j \in [K]} (\sigma(\mathbf{x}_i^T \mathbf{w}_j) - y_i) \cdot \sigma''(\mathbf{x}_i^T \mathbf{w}_k) \cdot \mathbf{x}_i \mathbf{x}_i^T + \sum_{i \in \mathbf{S}^{(t)}} [\sigma'(\mathbf{x}_i^T \mathbf{w}_k)]^2 \cdot \mathbf{x}_i \mathbf{x}_i^T \right)$$

Then, noting that $\sigma''(x) = 0$ a.s. we have,

$$\frac{\partial^2 \mathcal{R}(\mathbf{W}; \mathbf{S}^{(t)})}{\partial \mathbf{w}_k^2} = \frac{2}{(1-\epsilon)N} \cdot \sum_{i \in \mathbf{S}^{(t)}} [\sigma'(\mathbf{x}_i^T \mathbf{w}_k)]^2 \cdot \mathbf{x}_i \mathbf{x}_i^T$$

Step 1: Upper bounding the Frobenius norm distance between $\mathbf{W}^{(t+1)}$ and \mathbf{W}^* .

$$\begin{aligned}
\|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_F &= \|\mathbf{W}^{(t)} - \mathbf{W}^* - \eta \nabla \mathcal{R}(\mathbf{W}^{(t)}; \mathbf{S}^{(t)})\|_F \\
&= \|\mathbf{W}^{(t)} - \mathbf{W}^* - \eta \nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{TP}) - \nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{FP})\|_F \\
&\leq \|\mathbf{W}^{(t)} - \mathbf{W}^* - \eta \nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{TP})\|_F + \|\eta \nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{FP})\|_F
\end{aligned} \tag{B.8}$$

We will first upper bound the first term of Equation (B.8) through its square.

$$\|\mathbf{W}^{(t)} - \mathbf{W}^* - \eta \nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{TP})\|_F^2$$

$$\begin{aligned}
&= \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F^2 - 2\eta \cdot \langle \mathbf{W}^{(t)} - \mathbf{W}^*, \nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{TP}) \rangle_{\text{Tr}} + \eta^2 \cdot \|\nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{TP})\|^2 \\
&= \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F^2 - 2\eta \cdot \langle \mathbf{W}^{(t)} - \mathbf{W}^*, \nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{TP}) - \nabla \mathcal{R}(\mathbf{W}^*, \text{TP}) \rangle_{\text{Tr}} \\
&\quad + \eta^2 \cdot \|\nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{TP}) - \nabla \mathcal{R}(\mathbf{W}^*, \text{TP})\|_F^2
\end{aligned} \tag{B.9}$$

In the above, in the final equality we use the fact that $\nabla \mathcal{R}(\mathbf{W}^*, \text{TP}) = \mathbf{0}$. Then from Equation (B.3), we can calculate the smoothness constant of $\mathcal{R}(\mathbf{W}^{(t)}; \text{TP})$ over \mathbf{w}_j for $j \in [K]$ and fixing \mathbf{w}_ℓ for $\ell \in [K]$ and $\ell \neq j$ as constant vectors,

$$\frac{2}{(1-\epsilon)N} \cdot \sum_{i \in \text{TP}} [(\sigma'(\mathbf{x}_i^T \mathbf{w}_j))^2 \cdot \mathbf{x}_i \mathbf{x}_i^T] \preceq \frac{2C_2^2}{(1-\epsilon)N} \cdot \|\mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T\|_2 \cdot \mathbf{I} \triangleq \frac{\beta}{(1-\epsilon)N} \cdot \mathbf{I} \tag{B.10}$$

We can also similarly obtain the strong convexity constant,

$$\frac{2}{(1-\epsilon)N} \cdot \sum_{i \in \text{TP}} [(\sigma'(\mathbf{x}_i^T \mathbf{w}_j))^2 \cdot \mathbf{x}_i \mathbf{x}_i^T] \succeq \frac{2C_1^2}{(1-\epsilon)N} \cdot \lambda_{\min}(\mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T) \cdot \mathbf{I} \triangleq \frac{\alpha}{(1-\epsilon)N} \cdot \mathbf{I} \tag{B.11}$$

We can now note

$$\begin{aligned}
&2\eta \cdot \langle \mathbf{W}^{(t)} - \mathbf{W}^*, \nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{TP}) - \nabla \mathcal{R}(\mathbf{W}^*, \text{TP}) \rangle_{\text{Tr}} \\
&= 2\eta \cdot \langle \text{vec}(\mathbf{W}^{(t)}) - \text{vec}(\mathbf{W}^*), \text{vec}(\nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{TP})) - \text{vec}(\nabla \mathcal{R}(\mathbf{W}^*, \text{TP})) \rangle \\
&= 2\eta \cdot \sum_{k \in [K]} \left\langle \mathbf{w}_k^{(t)} - \mathbf{w}_k^*, \frac{\partial \mathcal{R}(\mathbf{W}^{(t)}; \text{TP})}{\partial \mathbf{w}_k} - \frac{\partial \mathcal{R}(\mathbf{W}^*; \text{TP})}{\partial \mathbf{w}_k} \right\rangle \\
&\geq \frac{2C_1^2 \eta}{(1-\epsilon)N} \cdot \lambda_{\min}(\mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T) \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F^2 + \frac{2C_2^2 \eta}{(1-\epsilon)N} \cdot \|\mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}\|^{-1} \|\nabla \mathcal{L}(\mathbf{W}^{(t)}; \text{TP}) - \nabla \mathcal{L}(\mathbf{W}^*; \text{TP})\|_F^2
\end{aligned}$$

In the above, the first equality follows from Lemma 28, in the first inequality we apply Lemma 25 with α and β defined in Equations (B.10) and (B.11), respectively, and f defined as $\mathcal{R}(\mathbf{W}^{(t)}; \text{TP})$ where \mathbf{w}_j for $j \in [K]$ is the input variable and \mathbf{w}_ℓ for $\ell \in [K]$ and $\ell \neq j$ as fixed vectors, and in the second equality we simply use the definition of the Frobenius norm. Then, solving the quadratic equation, we set $\eta \leq 2C_2^{-2}(1-\epsilon)N\|\mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T\|_2^{-1}$, and obtain

$$\|\mathbf{W}^{(t+1)} - \mathbf{W}^* - \eta \nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{TP})\|_F^2 \leq \left(1 - \frac{2C_1^2 \eta}{(1-\epsilon)N} \cdot \lambda_{\min}(\mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T)\right) \cdot \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F^2$$

Step 2: Upper bounding the corrupted gradient. We now upper bound the third term in Equation (B.9).

$$\begin{aligned}
\eta^2 \cdot \|\nabla \mathcal{R}(\mathbf{W}^{(t)}; \text{FP})\|_F^2 &= \frac{4\eta^2}{[(1-\epsilon)N]^2} \cdot \sum_{k \in [K]} \left\| \sum_{i \in \text{FP}} \sum_{j \in [K]} (\sigma(\mathbf{x}_i^T \mathbf{w}_j) - \sigma(\mathbf{x}_i^T \mathbf{w}_j^*)) \cdot \sigma'(\mathbf{x}_i^T \mathbf{w}_k^{(t)}) \cdot \mathbf{x}_i \right\|_2^2 \\
&\leq \frac{4KC_2^2 \eta^2}{[(1-\epsilon)N]^2} \cdot \|\mathbf{X}_{\text{FP}} \mathbf{X}_{\text{FP}}^T\|_2 \sum_{i \in \text{FP}} \left(\sum_{j \in [K]} \sigma(\mathbf{x}_i^T \mathbf{w}_j) - \sigma(\mathbf{x}_i^T \mathbf{w}_j^*) \right)^2 \\
&\leq \frac{4KC_2^2 \eta^2}{[(1-\epsilon)N]^2} \cdot \|\mathbf{X}_{\text{FP}} \mathbf{X}_{\text{FP}}^T\|_2 \sum_{i \in \text{FN}} \left(\sum_{j \in [K]} \sigma(\mathbf{x}_i^T \mathbf{w}_j) - \sigma(\mathbf{x}_i^T \mathbf{w}_j^*) \right)^2 \\
&\leq \frac{4K^2 C_2^2 \eta^2}{[(1-\epsilon)N]^2} \cdot \|\mathbf{X}_{\text{FP}} \mathbf{X}_{\text{FP}}^T\|_2 \sum_{i \in \text{FN}} \sum_{j \in [K]} (\sigma(\mathbf{x}_i^T \mathbf{w}_j) - \sigma(\mathbf{x}_i^T \mathbf{w}_j^*))^2 \\
&\leq \frac{4K^2 C_2^4 \eta^2}{[(1-\epsilon)N]^2} \cdot \|\mathbf{X}_{\text{FP}} \mathbf{X}_{\text{FP}}^T\|_2 \|\mathbf{X}_{\text{FN}} \mathbf{X}_{\text{FN}}^T\|_2 \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F^2
\end{aligned}$$

In the above, the first inequality follows from Lemma 23, the second inequality follows from the optimality of the Subquantile set, the third inequality follows from the AM-QM inequality (see Lemma 27), the final inequality follows from the C_2 -Lipschitzness of σ .

Step 3: Combining Step 1 and Step 2. We obtain,

$$\|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_F^2 \leq \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F^2 \left(1 - \frac{C_1^2 \eta}{(1-\epsilon)N} \cdot \lambda_{\min}(\mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T) + \frac{4K^2 C_2^4 C_3 \eta^2}{(1-\epsilon)N} \cdot \|\mathbf{X}_{\text{FN}} \mathbf{X}_{\text{FN}}^T\|_2 \right)$$

From which we obtain clearly that there exists sufficiently small ϵ such that

$$\epsilon = O(\text{poly}(K, C_1, C_2, C_3, (1-\epsilon)N))$$

That gives us

$$\|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_F^2 \leq \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F^2 (1 - \Theta(C_1^2 C_2^{-2} \cdot \kappa^{-1}(\mathbf{X}_{\text{TP}} \mathbf{X}_{\text{TP}}^T)))$$

We solve the induction to obtain the iteration complexity and complete the proof. \blacksquare

C Probability Theory

In this section we will give various concentration inequalities on the inlier data for functions in the Reproducing Kernel Hilbert Space.

C.1 Finite Dimensional Concentrations of Measure

Lemma 17 (Upper Bound on Sum of Chi-Squared Variables [LM00]). *Suppose $\xi_i \sim \mathcal{N}(0, \sigma^2)$ for $i \in [n]$, then*

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

Proposition 18 (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 $\epsilon \in (0, 1]$ and let \mathcal{C} represent all such subsets. Then, with probability exceeding $1 - \delta$,*

$$\max_{S \in \mathcal{C}} \|\boldsymbol{\xi}_S\|_2^2 \leq 30\sigma n \epsilon \log \epsilon^{-1}$$

Proof. Directly from Lemma 17, we have with probability exceeding $1 - \delta$.

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

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

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

In the first inequality we apply a union bound over \mathcal{C} with Lemma 29, and in the second inequality we use Lemma 17. We then obtain,

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

In the above, in the first inequality, we note that $\log\left(\frac{n}{\epsilon n}\right) \leq 3n \epsilon \log \epsilon^{-1}$ as $\epsilon < 0.5$, in the second inequality we note that $\sqrt{\log \epsilon^{-1}} \leq C \log \epsilon^{-1}$ for a $C \leq (\log(2))^{-1/2} \leq \sqrt{3}$ 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

C.2 Hilbert Space Concentrations of Measure

Proposition 19 (Jensen's Inequality [Jen06]). *Suppose φ is a convex function, then for a random variable X , it holds*

$$\varphi(\mathbf{E}[X]) \leq \mathbf{E}[\varphi(X)]$$

The inequality is reversed for φ concave.

We will now study the covariance approximation problem. Our main probabilistic tool will be McDiarmid's Inequality.

Proposition 20 (McDiarmid's Inequality [M⁺89]). *Suppose $f : \mathcal{X}_1 \times \dots \times \mathcal{X}_n \rightarrow \mathbb{R}$. Consider i.i.d X_1, \dots, X_n where $X_i \in \mathcal{X}_i$ for all $i \in [n]$. If there exists constants c_1, \dots, c_n , such that for all $x_i \in \mathcal{X}_i$ for all $i \in [n]$, it holds*

$$\sup_{\tilde{X}_i \in \mathcal{X}_i} |f(X_1, \dots, X_{i-1}, X_i, X_{i+1}, \dots, X_n) - f(X_1, \dots, X_{i-1}, \tilde{X}_i, X_{i+1}, \dots, X_n)| \leq c_i$$

Then for any $t > 0$, it holds

$$\Pr\{f(X_1, \dots, X_n) - \mathbf{E}[f(X_1, \dots, X_n)] \geq t\} \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n c_i^2}\right)$$

Theorem 21 (Mean Estimation in the Hilbert Space [TSM⁺17]). *Define $P_n \triangleq \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$ and P be the distribution of the covariates in \mathcal{X} . Suppose $r : \mathcal{X} \rightarrow \mathcal{H}$ is a continuous function such that $\sup_{X \in \mathcal{X}} \|r(X)\|_{\mathcal{H}}^2 \leq C_k < \infty$. Then with probability at least $1 - \delta$,*

$$\left\| \int_{\mathcal{X}} r(x) dP_n(x) - \int_{\mathcal{X}} r(x) dP(x) \right\| \leq \sqrt{\frac{C_k}{n}} + \sqrt{\frac{2C_k \log(1/\delta)}{n}}$$

We will strengthen upon the result by [TSM⁺17] by using knowledge of the distribution to first derive the expectation.

Proposition 22 (Probabilistic Bound on Infinite Dimensional Covariance Estimation). *Let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be i.i.d sampled from \mathcal{P} such that $\phi(\mathbf{x}_i) \triangleq X_i \sim \mathcal{P}_{\sharp} \phi$ (Assumption 7). Denote \mathcal{S} as all subsets of $[n]$ with size from $(1 - 2\epsilon)N$ to $(1 - \epsilon)N$. We then have simultaneously with probability exceeding $1 - \delta$,*

$$\begin{aligned} \left\| \frac{1}{n} \sum_{i=1}^n X_i \otimes X_i - \Sigma \right\|_{\text{HS}} &\leq \sqrt{\frac{8}{n}} \|\Gamma\|_{\text{op}} + \sqrt{\frac{2 \log(2/\delta)}{n}} P_k \\ \max_{A \in \mathcal{S}} \left\| \frac{1}{(1 - \epsilon)N} \sum_{i \in A} X_i \otimes X_i - \Sigma \right\|_{\text{HS}} &\leq \sqrt{\frac{8}{(1 - \epsilon)N}} \|\Gamma\|_{\text{op}} + \sqrt{\frac{2P_k^2 \log(2/\delta)}{(1 - \epsilon)N}} + P_k \sqrt{\frac{\epsilon \log \epsilon^{-1}}{(1 - \epsilon)}} \end{aligned}$$

Proof. We will calculate the mean operator in the Hilbert Space $\mathcal{H} \otimes \mathcal{H}$ and use the \sqrt{n} -consistency of estimating the mean-element in a Hilbert Space to obtain the probability bounds.

$$\begin{aligned} \mathbf{E}_{X_i \sim \mathcal{P}_{\sharp} \phi} \left\| \frac{1}{(1 - \epsilon)N} \sum_{i=1}^{(1 - \epsilon)N} X_i \otimes X_i - \Sigma \right\|_{\text{HS}} &\stackrel{(ii)}{\leq} \mathbf{E}_{X_i \sim \mathcal{P}_{\sharp} \phi} \mathbf{E}_{\tilde{X}_i \sim \mathcal{P}_{\sharp} \phi} \left\| \frac{1}{(1 - \epsilon)N} \sum_{i=1}^{(1 - \epsilon)N} X_i \otimes X_i - \tilde{X}_i \otimes \tilde{X}_i \right\|_{\text{HS}} \\ &\stackrel{(iii)}{=} \mathbf{E}_{X_i \sim \mathcal{P}_{\sharp} \phi} \mathbf{E}_{\tilde{X}_i \sim \mathcal{P}_{\sharp} \phi} \mathbf{E}_{\xi_i \sim \mathcal{R}} \left\| \frac{1}{(1 - \epsilon)N} \sum_{i=1}^{(1 - \epsilon)N} \xi_i (X_i \otimes X_i - \tilde{X}_i \otimes \tilde{X}_i) \right\|_{\text{HS}} \\ &\leq \mathbf{E}_{X_i \sim \mathcal{P}_{\sharp} \phi} \mathbf{E}_{\xi_i \sim \mathcal{R}} \left\| \frac{2}{(1 - \epsilon)N} \sum_{i=1}^{(1 - \epsilon)N} \xi_i (X_i \otimes X_i) \right\|_{\text{HS}} \\ &\leq \frac{2}{(1 - \epsilon)N} \mathbf{E}_{X_i \sim \mathcal{P}_{\sharp} \phi} \left(\mathbf{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^{(1 - \epsilon)N} \xi_i (X_i \otimes X_i) \right\|_{\text{HS}}^2 \right)^{1/2} \end{aligned}$$

In (ii) we note that $X_i \otimes X_i - \mathbf{\Gamma}$ is a mean $\mathbf{0}$ operator in the tensor product space $\mathcal{H} \otimes \mathcal{H}$. Then for $X, Y \in \mathcal{H} \otimes \mathcal{H}$ s.t. $\mathbf{E}[Y] = \mathbf{0}$ it follows $\|X\|_{\text{HS}} = \|X - \mathbf{E}[Y]\|_{\text{HS}} = \|\mathbf{E}[X - Y]\|_{\text{HS}}$ and finally we apply Jensen's Inequality. Let e_k for $k \in [p]$ (p possibly infinite) represent a complete orthonormal basis for the image of $\mathbf{\Gamma}$. By expanding out the Hilbert-Schmidt Norm, we then have

$$\begin{aligned}
& \frac{2}{(1-\epsilon)N} \left(\mathbf{E}_{X_i \sim \mathcal{P}_\# \phi} \mathbf{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^{(1-\epsilon)N} \xi_i (X_i \otimes X_i) \right\|_{\text{HS}}^2 \right)^{1/2} \\
&= \frac{2}{(1-\epsilon)N} \left(\mathbf{E}_{X_i \sim \mathcal{P}_\# \phi} \mathbf{E}_{\xi_i \sim \mathcal{R}} \sum_{k=1}^p \left\langle \sum_{i=1}^{(1-\epsilon)N} \xi_i (X_i \otimes X_i) e_k, \sum_{j=1}^{(1-\epsilon)N} \xi_j (X_j \otimes X_j) e_k \right\rangle_{\mathcal{H}} \right)^{1/2} \\
&= \frac{2}{(1-\epsilon)N} \left(\mathbf{E}_{X_i \sim \mathcal{P}_\# \phi} \mathbf{E}_{\xi_i \sim \mathcal{R}} \sum_{k=1}^p \sum_{i=1}^{(1-\epsilon)N} \sum_{j=1}^{(1-\epsilon)N} \xi_i \xi_j \langle (X_i \otimes X_i) e_k, (X_j \otimes X_j) e_k \rangle_{\mathcal{H}} \right)^{1/2} \\
&\stackrel{(iv)}{=} \frac{2}{(1-\epsilon)N} \left(\mathbf{E}_{X_i \sim \mathcal{P}_\# \phi} \sum_{k=1}^p \sum_{i=1}^{(1-\epsilon)N} \langle (X_i \otimes X_i) e_k, (X_i \otimes X_i) e_k \rangle_{\mathcal{H}} \right)^{1/2} \\
&= \frac{2}{(1-\epsilon)N} \left(\sum_{i=1}^{(1-\epsilon)N} \mathbf{E}_{X_i \sim \mathcal{P}_\# \phi} \|X_i \otimes X_i\|_{\text{HS}}^2 \right)^{1/2} \\
&\stackrel{(v)}{=} \frac{2}{\sqrt{(1-\epsilon)N}} \left(\mathbf{E}_{X_i \sim \mathcal{P}_\# \phi} \|X_i\|_{\mathcal{H}}^4 \right)^{1/2}
\end{aligned}$$

(iv) follows from noticing $\mathbf{E}_{\xi_i, \xi_j \sim \mathcal{R}} [\xi_i \xi_j] = \delta_{ij}$. (v) follows from expanding the Hilbert-Schmidt Norm and applying Parseval's Identity. We will now calculate the fourth moment of a norm of sub-Gaussian function in the Hilbert Space.

$$\begin{aligned}
& \mathbf{E}_{X \sim \mathcal{P}_\# \phi} [\|X\|_{\mathcal{H}}^4] = \int_0^\infty \mathbf{Pr}_{X \sim \mathcal{P}_\# \phi} \{ \|X\|_{\mathcal{H}}^4 \geq t \} dt = \int_0^\infty \mathbf{Pr}_{X \sim \mathcal{P}_\# \phi} \{ \|X\|_{\mathcal{H}} \geq t^{1/4} \} dt \\
&\stackrel{(vi)}{\leq} \int_0^\infty \inf_{\theta > 0} \mathbf{E}_{X \sim \mathcal{P}_\# \phi} [\exp(\theta \|X\|_{\mathcal{H}})] \exp(-\theta t^{1/4}) dt \leq \int_0^\infty \inf_{\theta > 0} \exp\left(\frac{\theta^2 \|\mathbf{\Gamma}\|_{\text{op}}}{2} - \theta t^{1/4}\right) dt \\
&= \int_0^\infty \exp\left(-\frac{\sqrt{t}}{\|\mathbf{\Gamma}\|_{\text{op}}}\right) dt = 2\|\mathbf{\Gamma}\|_{\text{op}}^2
\end{aligned}$$

In (vi) we apply Markov's Inequality. From which we obtain,

$$\mathbf{E}_{X_i \sim \mathcal{P}_\# \phi} \left\| \frac{1}{(1-\epsilon)N} \sum_{i=1}^{(1-\epsilon)N} X_i \otimes X_i - \mathbf{\Sigma} \right\|_{\text{HS}} \leq \sqrt{\frac{8}{(1-\epsilon)N}} \|\mathbf{\Gamma}\|_{\text{op}}$$

Then, define the function $r(\mathbf{x}) : \mathcal{X} \rightarrow \mathcal{H} \otimes \mathcal{H}$, $\mathbf{x} \rightarrow \phi(\mathbf{x}) \otimes \phi(\mathbf{x})$. From Assumption 8, we have $r(\mathbf{x}) = \|\phi(\mathbf{x}) \otimes \phi(\mathbf{x})\|_{\text{HS}} \leq \|\phi(\mathbf{x})\|_{\mathcal{H}}^2 \leq P_k$. We will use McDiarmid's Inequality, consider $\tilde{P} \triangleq \delta_{X_i}$ with one modified element. Then consider the equation $f(x_1, \dots, x_n) : \mathcal{X} \times \dots \times \mathcal{X} \rightarrow \mathcal{H} \otimes \mathcal{H} \times \dots \times \mathcal{H} \otimes \mathcal{H}$, $x_1, \dots, x_n \rightarrow \int_{\mathcal{X}} r(x) dP_B(x) - \int_{\mathcal{X}} r(x) dP(x) \|_{\text{HS}}$.

$$\begin{aligned}
& \left\| \int_{\mathcal{X}} r(x) dP_B(x) - \int_{\mathcal{X}} r(x) dP(x) \right\|_{\text{HS}} - \left\| \int_{\mathcal{X}} r(x) d\tilde{P}(x) - \int_{\mathcal{X}} r(x) dP(x) \right\|_{\text{HS}} \\
&\leq \frac{1}{(1-\epsilon)N} (\|r(x_i)\|_{\text{HS}} + \|r(\tilde{x}_i)\|_{\text{HS}}) \leq \frac{2P_k}{(1-\epsilon)N}
\end{aligned}$$

Then, we have from McDiarmid's inequality (Proposition 20),

$$\mathbf{Pr} \left\{ \left\| \int_{\mathcal{X}} r(x) dP_B(x) - \int_{\mathcal{X}} r(x) dP(x) \right\|_{\text{HS}} - \sqrt{\frac{8}{(1-\epsilon)N}} \|\mathbf{\Gamma}\|_{\text{op}} \geq t \right\} \leq \exp\left(-\frac{t^2(1-\epsilon)N}{P_k^2}\right)$$

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

$$\left\| \int_{\mathcal{X}} r(x) dP_B(x) - \int_{\mathcal{X}} r(x) dP(x) \right\|_{\text{HS}} \leq \sqrt{\frac{8}{(1-\epsilon)N}} \|\mathbf{\Gamma}\|_{\text{op}} + \sqrt{\frac{P_k^2 \log(2/\delta)}{(1-\epsilon)N}}$$

Next, applying a union bound over \mathcal{S} with lemma 29, we have

$$\max_{B \in \mathcal{S}} \left\| \int_{\mathcal{X}} r(x) dP_B(x) - \int_{\mathcal{X}} r(x) dP(x) \right\|_{\text{HS}} \leq \sqrt{\frac{8}{(1-\epsilon)N}} \|\mathbf{\Gamma}\|_{\text{op}} + \sqrt{\frac{P_k^2 \log(2/\delta)}{(1-\epsilon)N} + \frac{P_k^2 \epsilon \log \epsilon^{-1}}{(1-\epsilon)}}$$

Simplifying the resultant bound completes the proof. ■

D Additional Lemmas

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

Lemma 23. *Let $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ and $\mathbf{X} \in \mathbb{R}^{p \times n}$, then*

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

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

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

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 24. *Consider a determinate set of numbers $(a_i)_{i=1}^n$, and determinate set of functions in the Hilbert Space, $(X_i)_{i=1}^n$, then*

$$\left\| \sum_{i=1}^n a_i X_i \right\|_{\mathcal{H}}^2 \leq \|\mathbf{a}\|_2^2 \|\mathbf{K}\|$$

Proof. The proof is a calculation.

$$\left\| \sum_{i=1}^n a_i X_i \right\|_{\mathcal{H}}^2 = \left\langle \sum_{i=1}^n a_i X_i, \sum_{j=1}^n a_j X_j \right\rangle_{\mathcal{H}} = \sum_{i=1}^n \sum_{j=1}^n a_i a_j k(x_i, x_j) = \mathbf{a}^T \mathbf{K} \mathbf{a} \leq \|\mathbf{a}\|_2^2 \|\mathbf{K}\|$$

where $[\mathbf{K}]_{i,j} = k(x_i, x_j)$ is the kernel Gram matrix. ■

Lemma 25 (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 26 (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*

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

Lemma 27 (AM-QM Inequality). *Let x_1, \dots, x_n be scalars in \mathbb{R} , then*

$$\left(\sum_{i \in [n]} x_i\right)^2 \leq n \sum_{i \in [n]} x_i^2$$

Lemma 28. *Suppose $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$, then*

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

Lemma 29 (Sum of Binomial Coefficients [CLRS22]). *Let $k, n \in \mathbb{N}$ such that $k \leq n$, then*

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