

Subquantile Minimization for Kernel Learning in the Huber ϵ -Contamination Model*

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

Alex Gittens
RPI CS, gittaa@rpi.edu

Abstract

In this paper we propose Subquantile Minimization for learning with adversarial corruption in the training set. Superquantile objectives have been formed in the past in the context of fairness where one wants to learn an underrepresented distribution equally [19, 30]. Our intuition is to learn a more favorable representation of the *majority* class, thus we propose to optimize over the p -subquantile of the loss in the dataset. In particular, we study the Huber Contamination Problem for Kernel Learning where the distribution is formed as, $\hat{\mathbb{P}} = (1 - \epsilon)\mathbb{P} + \epsilon\mathbb{Q}$, and we want to find the function $\inf_f \mathbb{E}_{x \in \mathbb{P}} [\ell_f(x)]$, from the noisy distribution, $\hat{\mathbb{P}}$. We assume the adversary has knowledge of the true distribution of \mathbb{P} , and is able to corrupt the covariates and the labels of ϵ samples. To our knowledge, we are the first to study the problem of general kernel learning in the Huber Contamination Model. In our theoretical analysis, we analyze our non-convex concave objective function with the Moreau Envelope. We show (i) a stationary point with respect to the Moreau Envelope is a good point and (ii) we can reach a stationary point with gradient descent methods. Further, we analyze accelerated gradient methods for the non-convex concave minimax optimization problem. We empirically test Kernel Ridge Regression and Kernel Classification on various state of the art datasets and show Subquantile Minimization gives strong results. Furthermore, we run experiments on various datasets and compare with the state-of-the-art algorithms to show the superior performance of Subquantile Minimization.

*Preliminary Work

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), [8, 1, 21, 25, 11] as well as for linear regression [2, 24]. Robust Statistics has been studied extensively [9] for problems such as high-dimensional mean estimation [26, 4] and Robust Covariance Estimation [5, 10]. Recently, there has been an interest in solving robust machine learning problems by gradient descent [27, 8]. Subquantile minimization aims to address the shortcomings of standard ERM in applications of noisy/corrupted data [18, 17]. In many real-world applications, linear models are insufficient to model the data. Therefore, we consider the problem of Robust Learning for Kernel Learning.

Definition 1. (Huber ϵ -Contamination Model [16]). Given a corruption parameter $0 < \epsilon < 0.5$, a data matrix, X and labels y . An adversary is allowed to inspect all samples and modify $n\epsilon$ samples arbitrarily. The algorithm is then given the ϵ -corrupted data matrix X and y as training data.

Current approaches for robust learning across various machine learning tasks often use gradient descent over a robust objective, [21]. 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 has also been proposed recently in [15] where there has been an analysis in the Binary Classification Task. We show Subquantile Minimization reduces to the same objective in [15]. We use theory from the weakly-convex concave optimization literature for our error bounds. We are able to leverage this theory by analyzing the asymptotic distribution of a softplus approximation of the Subquantile objective.

Theorem 2. (Informal). Let the dataset be given as $\{(x_i, y_i)\}_{i=1}^n$ such that the labels and features of ϵn samples are arbitrarily corrupted by an adversary. Assume Subquantile Minimization returns $f_{\hat{w}}$ for $n \geq \frac{(1-2\epsilon)(C_k \|\Sigma\|_{\text{op}} + \beta)}{(1-c_1)\lambda_{\min}(\Sigma)} + \sqrt{\beta}$ for a constant $c_1 \in (0, 1)$ such that for Kernelized Regression:

$$\mathbb{E}_{\mathcal{D} \sim \mathbb{P}} \|f_{\hat{w}} - f_w^*\|_{\mathcal{H}} \leq O\left(\frac{\gamma\sigma}{\sqrt{\lambda_{\min}(\Sigma)}}\right)$$

where $\epsilon \rightarrow 0$ as number of gradient descent iterations goes to ∞ and $\Sigma = \mathbb{E}[\phi(x) \otimes \phi(x)]$.

Kernel Binary Classification:

$$\mathbb{E}_{\mathcal{D} \sim \mathbb{P}} \|f_{\hat{w}} - f_w^*\|_{\mathcal{H}} \leq O\left(\frac{\sqrt{\text{Tr}(\Sigma)} + \sqrt{Q_k}}{\sqrt{n(1-2\epsilon)\lambda_{\min}(\Sigma)}}\right)$$

Kernel Multi-Class Classification:

$$\mathbb{E}_{\mathcal{D} \sim \mathbb{P}} \|f_{\hat{W}} - f_W^*\|_{\mathcal{H}} \leq O(\gamma)$$

1.1 Related Work

The idea of iterative thresholding algorithms for robust learning tasks dates back to 1806 by Legendre [20]. From the popularity of Machine Learning, numerous algorithms have been developed in this ideology. Therefore, we will dedicate this section to reviewing such works and to make clear our contributions to the iterative thresholding literature.

1.1.1 Robust Regression via Hard Thresholding [3]

Bhatia et al. consider robust linear regression by considering an active set S , which contains the points with the lowest error. This set is updated each iteration in conjunction with either a full solve (TORRENT-FC) or a gradient iteration (TORRENT-GD). TORRENT-GD is an unconstrained variant of our algorithm. The main limitation of this work is that only the case of label corruption is considered. We pick up the result of Theorem 9 and Theorem 11 in [3] (up to constants) for linear regression with and without feature corruption, which is one of our key contributions.

1.1.2 Learning with bad training data via iterative trimmed loss minimization [33]

This work considers optimizing over the bottom- k errors by choosing the αn points with smallest error and then updating the model from these αn . This general model is the same as ours. Theoretically, this work considers only general linear models. Experimentally, this work considers more general machine learning models such as GANS.

1.1.3 Trimmed Maximum Likelihood Estimation for Robust Generalized Linear Model [1]

This work studies a different class of generalized linear models. Interestingly, they show for Gaussian Regression the iterative trimmed maximum likelihood estimator is able to achieve near minimax optimal error. This work does not consider feature corruption and primarily focuses on the covariates sampled with Gaussian Design from Identity covariance.

1.1.4 Sum of Ranked Range Loss for Supervised Learning [15]

Hu et al. proposed learning over the bottom k losses, this is an alternative formulation of our algorithm. They solve their optimization problem with difference of sums convex solvers. This work considers only the classification task and does not give rigorous error bounds.

1.2 Contributions

We will now state our main contributions clearly.

1. We provide a novel theoretical framework using the Moreau Envelope for analyzing the iterative trimmed estimator for machine learning tasks.
2. We provide rigorous error bounds for subquantile minimization in the kernel regression, kernel binary classification, and kernel multi-class classification. Furthermore, we provide our bounds for both label and feature comparison with a general Gaussian Design.
3. We perform experiments on state-of-the-art matrices in kernel learning and show the effectiveness of our algorithm compared to other robust meta-algorithms.

2 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 .

$$\mathcal{Q}_p(U) \leq t \iff \mathbb{P}\{U \leq t\} \geq p$$

The p -subquantile of the risk is then given by

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

Given an objective function, ℓ , the kernelized learning problem becomes:

$$f_{\hat{w}} = \arg \min_{f_w \in \mathcal{K}} \max_{t \in \mathbb{R}} \left\{ g(t, f_w) \triangleq t - \sum_{i=1}^n (t - (f_w(x_i) - y_i)^2)^+ \right\}$$

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 w . Thus, to solve this problem we use the iterations from equation 11 in [28]. Let $\Pi_{\mathcal{K}}$ be the projection of a vector on to the convex set $\mathcal{K} \triangleq \{f \in \mathcal{H} : \|f\|_{\mathcal{H}} \leq R\}$, then our update steps are

$$\begin{aligned} t^{(k+1)} &= \arg \max_{t \in \mathbb{R}} g(f_w^{(k)}, t) \\ f_w^{(k+1)} &= \Pi_{\mathcal{K}} \left(f_w^{(k)} - \alpha \nabla_f g(f_w^{(k)}, t^{(k+1)}) \right) \end{aligned}$$

We provide an algorithm for Subquantile Minimization of the ridge regression and classification kernel learning algorithm. Algorithm 1 is applicable to both kernel ridge regression and kernel classification.

Algorithm 1: SUBQ-GRADIENT

Input: Iterations: T ; Quantile: p ; Data Matrix: $X, (n \times d), n \gg d$; Learning schedule: $\alpha_1, \dots, \alpha_T$; Ridge parameter: λ

Output: Trained Parameters, $w_{(T)}$

```

1:  $w_{(0)} \leftarrow \mathcal{N}_d(0, \sigma)$ 
2: for  $k \in 1, 2, \dots, T$  do
3:    $S_{(k)} \leftarrow \text{SUBQUANTILE}(w^{(k)}, X)$ 
4:    $w^{(k+1)} \leftarrow w^{(k)} - \alpha_{(k)} \nabla_w g(t^{(k+1)}, w^{(k)})$ 
5: end
6: return  $w_{(T)}$ 

```

Algorithm 2: SUBQUANTILE

Input: Parameters w , Data Matrix: $X, (n \times d)$, Convex Loss Function f

Output: Subquantile Matrix S

```

1:  $\hat{\nu}_i \leftarrow \ell(x_i; f_w, y_i)$  s.t.  $\hat{\nu}_{i-1} \leq \hat{\nu}_i \leq \hat{\nu}_{i+1}$ 
2:  $t \leftarrow \hat{\nu}_{np}$ 
3: Let  $x_1, \dots, x_{np}$  be  $np$  points such that
    $\ell(x_i; f_w, y_i) \leq t$ 
4:  $S \leftarrow (x_1^\top \dots x_{np}^\top)^\top$ 
5: return  $S$ 

```

3 Structural Results

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 Section 2.

Lemma 3. *Let $f(x; w)$ be a convex loss function. Let x_1, x_2, \dots, x_n denote the n data points ordered such that $f(x_1; w, y_1) \leq f(x_2; w, y_2) \leq \dots \leq f(x_n; w, y_n)$. If we denote $\hat{\nu}_i \triangleq f(x_i; w, y_i)$, it then follows $\arg \max_{t \in \mathbb{R}} g(t, w) = \hat{\nu}_{np}$.*

Proof is given in ???. It therefore follows,

$$\sum_{i=1}^n \mathbb{I} \left\{ \hat{\nu}_{np} \geq \left(f_w^{(k)}(x_i) - y_i \right)^2 \right\} \left(f_w^{(k)}(x_i) - y_i \right)^2 \in \max_{t \in \mathbb{R}} g(t, f_w^{(k)})$$

Interpretation 4. From Lemma 3, we see the t will be greater than or equal to the errors of exactly np points. Thus, we are continuously updating over the np minimum errors.

Lemma 5. *Let $\hat{\nu}_i \triangleq f(x_i; w, y_i)$ s.t. $\hat{\nu}_{i-1} \leq \hat{\nu}_i \leq \hat{\nu}_{i+1}$, if we choose $t^{(k+1)} = \hat{\nu}_{np}$ as by Lemma 3, it then follows $\nabla_w g(t^{(k)}, f_w^{(k)}) = \frac{1}{np} \sum_{i=1}^{np} \nabla f(x_i; f_w^{(k)}, y_i)$*

Proof is given in Appendix B.2.

3.1 On the Softplus Approximation

It is clear our objective function is non-smooth. Thus we propose to use the Softplus approximation to smooth the function. The main idea is to *first* approximate ReLU, consider the theory with respect to the approximation, and then take the limit as the approximation goes to the ReLU. The softplus approximation is given as follows,

$$\zeta_\lambda(x) = \frac{1}{\lambda} \log(1 + e^{\lambda x})$$

We then have the approximation of g as

$$\begin{aligned} \tilde{g}_\lambda(t, f_w) &\triangleq t - \sum_{i=1}^n \zeta_\lambda(t - \ell(f_w; x_i, y_i)) \\ &= t - \frac{1}{np} \sum_{i=1}^n \frac{1}{\lambda} \log(1 + \exp(\lambda(t - \ell(f_w; x_i, y_i)))) \end{aligned}$$

Now we compute the derivatives w.r.t to the softplus approximation, and then we consider the limit of the derivative as $\lambda \rightarrow \infty$.

$$\begin{aligned}\nabla_t \tilde{g}_\lambda(t, f_w) &= \nabla_t \left(t - \frac{1}{np} \sum_{i=1}^n \frac{1}{\lambda} \ln(1 + \exp(\lambda(t - \ell(f_w; x_i, y_i)))) \right) \\ &= 1 - \frac{1}{np} \sum_{i=1}^n \sigma(\lambda(t - \ell(f_w; x_i, y_i)))\end{aligned}$$

where $\sigma(\cdot)$ is the sigmoid function. It therefore follows,

$$\lim_{\lambda \rightarrow \infty} \nabla_t \tilde{g}_\lambda(t, f_w) = 1 - \frac{1}{np} \sum_{i=1}^n \mathbb{I}\{t - \ell(f_w; x_i, y_i)\}$$

$$\begin{aligned}\nabla_f \tilde{g}_\lambda(t, f_w) &= \nabla_f \left(t - \frac{1}{np} \sum_{i=1}^n \frac{1}{\lambda} \ln(1 + \exp(\lambda(t - \ell(f_w; x_i, y_i)))) \right) \\ &= \frac{1}{np} \sum_{i=1}^n \nabla_f \ell(f_w; x_i, y_i) \sigma(\lambda(t - \ell(f_w; x_i, y_i)))\end{aligned}$$

We therefore similarly have,

$$\lim_{\lambda \rightarrow \infty} \nabla_f \tilde{g}_\lambda(t, f_w) = \frac{1}{np} \sum_{i=1}^n \mathbb{I}\{t - \ell(f_w; x_i, y_i)\} \nabla_f \ell(f_w; x_i, y_i)$$

We can then calculate the Lipschitz constant of the approximation function with respect to f_w .

Lemma 6 (Lipschitz continuous gradient). *Let $f_w, f_{\tilde{w}} \in \mathcal{K}$, then we have*

$$\lim_{\lambda \rightarrow \infty} |\nabla_f \tilde{g}_\lambda(t, f_w) - \nabla_f \tilde{g}_\lambda(t, f_{\tilde{w}})| \leq \beta \|f_w - f_{\tilde{w}}\|_{\mathcal{H}}$$

where

$$\beta = \frac{1}{np} \sum_{i=1}^n \|\nabla_f^2 \ell(f_w; x_i, y_i)\|_{\text{op}}$$

and β has no dependence on λ .

Proof is in Appendix B.3.

3.2 Weakly Convex Concave Optimization Theory

With our smoothed function, we are now able to use the weakly-convex concave minimization literature to analyze g . The Moreau Envelope can be interpreted as an infimal convolution of the function f . When f is ρ -weakly convex, if $\lambda \leq \rho^{-1}$, then the Moreau Envelope is smooth.

Definition 7. (Moreau Envelope on closed, convex set, [23]). Let f be proper lower semi-continuous convex function $\ell : \mathcal{K} \rightarrow \mathbb{R}$, where $\mathcal{K} \subset \mathcal{X}$ is a closed and convex set, then the Moreau Envelope is defined as:

$$\mathcal{M}_{\lambda\ell}(f_w) \triangleq \inf_{f_{\tilde{w}} \in \mathcal{K}} \left\{ \ell(f_{\tilde{w}}) + \frac{1}{2\rho} \|f_w - f_{\tilde{w}}\|_{\mathcal{H}}^2 \right\}$$

Definition 8. Define the function $\Phi(f_w) \triangleq \max_{t \in \mathbb{R}} g(t, f_w)$. This function is a L -weakly convex function in \mathcal{K} , i.e., $\Phi(f_w) + \frac{L}{2} \|f_w\|_{\mathcal{H}}^2$ is a convex function over w in the convex and compact set \mathcal{K} .

Definition 9 (First Order Stationary Point). Let $f_{\tilde{w}}$ be a first-order stationary point, then for any $f_w \in \mathcal{K}$, it follows

$$\nabla_f g(f_{\tilde{w}}) (\overline{f_w - f_{\tilde{w}}}) \geq 0 \quad \forall f_w \in \mathcal{K}$$

Definition 10 (Stationary Point of Moreau Envelope). A point $f_{\hat{w}}$ is a stationary point of the Moreau Envelope defined in Definition 7 of Φ defined in Definition 8 if

$$f_{\hat{w}} = \arg \inf_{f_w \in \mathcal{K}} \left\{ \Phi_{\lambda}(f_w) + \frac{1}{2\rho} \|f_w - f_{\hat{w}}\|_{\mathcal{H}}^2 \right\}$$

We will show that if a point f_w is a stationary point then this point is close to the optimal point for the uncorrupted distribution, i.e. $\|f_{\hat{w}} - f_w^*\|_{\mathcal{H}}$ is small.

Lemma 11 (Lower bound on distance from stationary point and optimal point). *Let Φ_{λ} be defined as in Definition 8, then if $f_{\hat{w}}$ is a stationary point as defined in Definition 10, then*

$$\lim_{\lambda \rightarrow \infty} (\Phi_{\lambda}(f_{\hat{w}}) - \Phi_{\lambda}(f_w^*)) \leq \beta \|f_{\hat{w}} - f_w^*\|_{\mathcal{H}}^2$$

We can now upper bound $\|f_{\hat{w}} - f_w^*\|_{\mathcal{H}}$. We proceed by contradiction, i.e. if a stationary point is sufficiently far from the optimal point, then this will break the stationary property proved in Lemma 11. This bound is different for each of the loss functions, so we must upper bound $\|f_{\hat{w}} - f_w^*\|_{\mathcal{H}}$ separately for each loss function with the same high level overview.

3.3 Kernelized Regression

The loss for the Kernel Ridge Regression problem for a single training pair (x_i, y_i) is given by the following equation

$$\ell(f_w; x_i, y_i) = (f_w(x_i) - y_i)^2$$

For our theory, we need the L -lipschitz constant and β -smoothness constant.

Lemma 12. (L -Lipschitz of $g(t, f_w)$ w.r.t f_w). *Let x_1, x_2, \dots, x_n , represent the data vectors. It then follows for any $f_w, f_{\hat{w}} \in \mathcal{K}$:*

$$|g(t, f_w) - g(t, f_{\hat{w}})| \leq L \|f_w - f_{\hat{w}}\|_{\mathcal{H}}$$

where

$$L = \frac{2R}{np} \left(\sum_{i=1}^n \sqrt{k(x_i, x_i)} \right)^2 + \frac{2\|y\|}{np} \left(\sum_{i=1}^n \sqrt{k(x_i, x_i)} \right)$$

Proof is given in Appendix B.5.

Lemma 13. (β -Smoothness of $g(t, w)$ w.r.t w). *Let x_1, x_2, \dots, x_n represent the rows of the data matrix X . It then follows:*

$$\|\nabla_w g(t, f_w) - \nabla_w g(t, f_{\hat{w}})\| \leq \beta \|f_w - f_{\hat{w}}\|_{\mathcal{H}}$$

where

$$\beta = \frac{2}{np} \sum_{i \in X} k(x_i, x_i) = \frac{2}{np} \text{Tr}(K)$$

Proof. W.L.O.G, let S be the set of points such that if $x \in S$, then $t \geq (f_w(x) - y)^2$. Since g is twice differentiable, we will analyze the second derivative.

$$\|\nabla_f^2 g(t, f_w)\|_{\mathcal{H}} = \left\| \frac{2}{np} \sum_{i=1}^n \mathbb{I}\{t \geq (f_w(x_i) - y_i)^2\} k(x_i, x_i) \right\| \leq \frac{2}{np} \sum_{i=1}^n k(x_i, x_i) = \frac{2}{np} \text{Tr}(K)$$

This concludes the proof. ■

Similar results for the Lipschitz Constant for non-kernelized learning algorithms can be seen in [35]. It is important to note that β is upper bounded as

$$\beta = \frac{2}{np} \text{Tr}(K) \leq \frac{2}{np} (n(1 - \varepsilon) \max_{i \in P} k(x_i, x_i) + n\varepsilon \max_{j \in Q} k(x_j, x_j)) = 2p^{-1}((1 - \varepsilon)P_k + \varepsilon Q_k)$$

which is independent of n .

Lemma 14. If $\|f_w - f_{w^*}\| \geq \eta$, then it follows

$$\lim_{\lambda \rightarrow \infty} (\Phi_\lambda(f_w) - \Phi_\lambda(f_{w^*})) \geq \eta^2 n(1 - 2\varepsilon) \lambda_{\min} \left(\mathbb{E}_{x \sim \mathbb{P}} [\phi(x) \otimes \phi(x)] \right) - O \left(\sigma \sqrt{n(1 - 2\varepsilon) \log(n(1 - 2\varepsilon)) \|\Sigma\|_{\text{HS}}} \right) - 2\eta \left\| \sum_{i \in S \cap P} \eta_i \phi(x_i) \right\| - \sum_{j \in P \setminus S} \eta_j^2$$

The proof is deferred to Appendix B.6.

Theorem 15 (Stationary Point for Kernelized Regression is Good). *Let $f_{\hat{w}}$ be a stationary point defined in Definition 10 for the function Φ defined in Definition 8. Then for a constant $c_1 \in (0, 1)$, if*

$$n \geq \frac{8 \text{Tr}(\Sigma)^2}{\lambda_{\min}(\Sigma)(1 - c_1)^2(1 - 2\varepsilon)} + \frac{8\beta}{(1 - c_1)^2(1 - 2\varepsilon)},$$

$$\mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} \|f_{\hat{w}} - f_{w^*}\|_{\mathcal{H}} \leq \left(\frac{\sigma \varepsilon n}{c_1 n(1 - 2\varepsilon) \lambda_{\min}(\Sigma)} \right)^{1/2} + \frac{O \left(\sigma \sqrt{\log(n(1 - 2\varepsilon)) \text{Tr}(\Sigma)} \right)}{c_1 \sqrt{n(1 - 2\varepsilon) \lambda_{\min}(\Sigma)}}$$

where β is the Lipschitz Gradient Constant given in Lemma 13.

The proof is given in Appendix B.7. In Theorem 15, we have an upper bound on the expected distance from a stationary point to the optimal point over the distance of the dataset. The numerator of the second term grows in $O(\sqrt{\log(n)})$ and the denominator grows in $O(\sqrt{n})$ as can be shown by choosing sufficiently large n . Asymptotically the second term will then go to 0. In the first term, we have both the numerator and denominator scale in $O(n)$. Furthermore, when we consider the case of feature noise, e.g. a large multiplicative term on the features, we simply require more data to obtain the same bounds. Such a result is corroborated in [32]. For the linear and polynomial kernel, we then have β increases, therefore to obtain the same bound on η as with no feature noise, we simply need more data.

Corollary 16 (¹Linear Regression Expected Error Bound). *Consider Subquantile Minimization for Linear Regression on the data X with optimal parameters w^* . Assume $x_i \sim \mathcal{N}(0, \Sigma)$ for $i \in [n]$. Then after T iterations of Algorithm 1, we have the following error bounds for robust kernelized linear regression. Given sufficient data*

$$\mathbb{E} \|w^{(T)} - w^*\|_2 \leq O \left(\sqrt{\frac{\varepsilon}{1 - 2\varepsilon}} \frac{\sigma}{\sqrt{\lambda_{\min}(\Sigma)}} \right)$$

Proof given in Appendix B.8. It is important to note in all our bounds, the $\sqrt{\frac{\varepsilon}{1 - 2\varepsilon}}$ is a theoretical worst case bound when the Subquantile contains the minimum possible number of uncorrupted points. In other words, we have $\gamma \triangleq \frac{|P \setminus S|}{|S \cap P|} \leq \frac{n\varepsilon}{n(1 - 2\varepsilon)} = \frac{\varepsilon}{1 - 2\varepsilon}$. So, as $|S \cap P|$ increases, we have a better error bound as $|P \setminus S|$ decreases. As is typical in the robust statistics literature, we make no assumptions on the distribution of the corrupted data so we cannot say anything about $|S \cap P|$. We will have γ decreases if stationary points give high error for corrupt points as our optimization procedure moves toward a stationary point.

3.4 Kernel Binary Classification

The Negative Log Likelihood for the the Kernel Classification problem is given by the following equation for a single training pair (x_i, y_i)

$$\ell(x_i, y_i; f_w) = -(y_i \log(\sigma(f_w(x_i))) - (1 - y_i) \log(1 - \sigma(f_w(x_i))))$$

Similar to Section 3.3, we require the L -Lipschitz constant and β -smoothness constant.

Lemma 17. (L -Lipschitz of $g(t, w)$ w.r.t w). *Let x_1, x_2, \dots, x_n , represent the data vectors. It then follows:*

$$|g(t, f_w) - g(t, f_{\hat{w}})| \leq L \|f_w - f_{\hat{w}}\|_{\mathcal{H}}$$

where

$$L = \frac{1}{np} \sum_{i \in X} \sqrt{k(x_i, x_i)} = \frac{1}{np} \text{Tr}(K)$$

¹In Progress

Proof is given in Appendix B.9.

Lemma 18. (β -Smoothness of $g(t, w)$ w.r.t w). Let x_1, x_2, \dots, x_n represent the rows of the data matrix X . It then follows:

$$\|\nabla_f g(t, f_w) - \nabla_f g(t, f_{\hat{w}})\| \leq \beta \|f_w - f_{\hat{w}}\|_{\mathcal{H}}$$

where

$$\beta = \frac{1}{4p} \sum_{i=1}^n k(x_i, x_i) = \frac{1}{4p} \text{Tr}(K)$$

Proof is given in Appendix B.10.

Lemma 19. ² Assume $f_{\hat{w}}$ is a first-order stationary point as defined in Definition 9. If $\|f_{\hat{w}} - f_{w^*}\|_{\mathcal{H}} \geq \eta$, then it follows

$$\mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} \|f_{\hat{w}} - f_w^*\|_{\mathcal{H}} \leq O \left(\frac{\sqrt{n(1-2\varepsilon) \text{Tr}(\Sigma)} + \sqrt{n\varepsilon Q_k}}{n(1-2\varepsilon)c_4 \lambda_{\min}(\Sigma)} \right)$$

Proof is given in Appendix B.11.

Theorem 20. [A stationary point is good for kernel binary classification] Let $f_{\hat{w}}$ be a stationary point defined in Definition 9 for the function Φ defined in Definition 8. Then for a constant $c_4 \in (0, 1)$, if $n \geq \frac{4 \text{Tr}(\Sigma)}{\lambda_{\min}(\Sigma)(1-2\varepsilon)(1-c_4)}$, then in expectation over the dataset distribution,

$$\mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} \|f_{\hat{w}} - f_w^*\|_{\mathcal{H}} \leq O \left(\frac{\sqrt{\text{Tr}(\Sigma)} + \sqrt{Q_k}}{\sqrt{n(1-2\varepsilon) \lambda_{\min}(\Sigma)}} \right)$$

Proof is given in Appendix B.12. This result although shows consistency, i.e. when $n \rightarrow \infty$, then we have in expectation $\|f_w - f_w^*\| \rightarrow 0$, however it does crucially rely on the fact that Q_k is bounded, and in general when n is not large, a large Q_k does affect the error bounds.

3.5 Kernel Multi-Class Classification

The Negative Log-Likelihood Loss for the the Kernel Multi-Class Classification problem is given by the following equation for a single training pair (x_i, y_i) , note W is now a matrix

$$\ell(x_i, y_i; W) = - \sum_{j=1}^{|\mathcal{Y}|} \mathbb{I}\{j = y_i\} \log \left(\frac{\exp(f_{W_j}(x_i))}{\sum_{k=1}^{|\mathcal{Y}|} \exp(f_{W_k}(x_i))} \right)$$

Lemma 21 (L -Lipschitz of $g(t, w)$ w.r.t w). ³ Let x_1, x_2, \dots, x_n , represent the data vectors. It then follows:

$$|g(t, f_w) - g(t, f_{\hat{w}})| \leq L \|f_w - f_{\hat{w}}\|_{\mathcal{H}}$$

where

$$L = \frac{1}{np} \text{Tr}(K)$$

Lemma 22. (β -Smoothness of $g(t, w)$ w.r.t w). Let x_1, x_2, \dots, x_n represent the rows of the data matrix X . It then follows:

$$\|\nabla_f g(t, f_w) - \nabla_f g(t, f_{\hat{w}})\| \leq \beta \|f_w - f_{\hat{w}}\|_{\mathcal{H}}$$

where

$$\beta = \frac{1}{4p} \sum_{i=1}^n k(x_i, x_i) = \frac{1}{4p} \text{Tr}(K)$$

Proof is given in ??.

²In Progress

³In Progress

Lemma 23. ⁴ Assume $f_{\hat{w}}$ is a first-order stationary point as defined in Definition 9. If $\|f_{\hat{w}} - f_{w^*}\|_{\mathcal{H}} \geq \eta$, then it follows

$$\mathbb{E}_{\mathcal{D} \sim \mathbb{P}} \|f_{\hat{w}} - f_w^*\|_{\mathcal{H}} \leq$$

Proof is given in Appendix B.15.

In practice, however, it is important to note that solving for $\|\nabla \Phi_{\lambda}\|_{\mathcal{H}} = 0$ is NP-Hard. Thus, we will analyze the approximate stationary point.

Lemma 24 ([29, 7]). Assume the function Φ is β -weakly convex. Let $\lambda < \frac{1}{\beta}$, and let $f_{\hat{w}} = \arg \min_{f_w \in \mathcal{K}} (\Phi(f_w) + \frac{1}{2\lambda} \|f_w - f_{\hat{w}}\|_{\mathcal{H}}^2)$, then $\|\nabla \Phi_{\lambda}(f_w)\|_{\mathcal{H}} \leq \epsilon$ implies:

$$\|f_{\hat{w}} - f_w\|_{\mathcal{H}} = \lambda \epsilon \quad \text{and} \quad \min_{g \in \partial \Phi(f_{\hat{w}}) + \partial \mathcal{I}_{\mathcal{K}}(f_{\hat{w}})} \|g\|_{\mathcal{H}} \leq \epsilon$$

4 Experiments

We perform numerical experiments on state of the art datasets comparing with other state of the art methods. We initialize the weights parameterizing f_w with the Glorot Initialization Scheme [12].

Algorithm 3: SUBQUANTILE-KERNEL

Input: Iterations: T ; Quantile: p ; Data Matrix: $X \in \mathbb{R}^{n \times d}$, $n \gg d$; Labels: $y \in \mathbb{R}^{n \times 1}$; Learning Rate schedule: $\alpha_1, \dots, \alpha_T$; Ridge parameter: λ

Output: Trained Parameters: $f_w^{(T)}$

- 1: $w_i^{(0)} \leftarrow \text{Unif} \left[-\sqrt{\frac{6}{n}}, \sqrt{\frac{6}{n}} \right]$, $\forall i \in [n]$ ▷ Initialize Weights Randomly
- 2: **for** $k = 1, 2, \dots, T$ **do**
- 3: $S^{(k)} \leftarrow \text{SUBQUANTILE}(f_w^{(k)}, X)$ ▷ Algorithm 2
- 4: $\nabla_{fg}(t^{(k+1)}, f_w^{(k)}) \leftarrow 2 \sum_{i \in S^{(k)}} (f_w^{(k)}(x_i) - y_i) \cdot k(x_i, \cdot)$ ▷ Regression
- 5: $\nabla_{fg}(t^{(k+1)}, f_w^{(k)}) \leftarrow \sum_{i \in S^{(k)}} (\sigma(f_w^{(k)}(x_i)) - y_i) \cdot k(x_i, \cdot)$ ▷ Binary Classification
- 6: $\nabla_{fg}(t^{(k+1)}, f_w^{(k)}) \leftarrow \sum_{i \in S^{(k)}} (\text{softmax}(f_w^{(k)}(x_i)) - y_i) \cdot k(x_i, \cdot)$ ▷ Multi-Class Classification
- 7: $f_w^{(k+1)} \leftarrow f_w^{(k)} - \alpha_{(k)} \nabla_{fg}(t^{(k+1)}, f_w^{(k)})$ ▷ f_w -update in Section 2
- 8: **end**
- 9: Pick t uniformly at random from $[T]$
- 10: **return** $f_w^{(t)}$

Algorithms	Test RMSE							
	Concrete		Wine Quality		Boston Housing		Drug	
	$\epsilon = 0.2(\downarrow)$	$\epsilon = 0.4(\downarrow)$	$\epsilon = 0.2(\downarrow)$	$\epsilon = 0.4(\downarrow)$	$\epsilon = 0.2(\downarrow)$	$\epsilon = 0.4(\downarrow)$	$\epsilon = 0.2(\downarrow)$	$\epsilon = 0.4(\downarrow)$
KRR	1.355 _(0.0934)	2.282 _(0.2063)	1.437 _(0.0979)	2.272 _(0.1088)	1.285 _(0.0896)	2.266 _(0.0686)	1.478 _(0.0533)	2.381 _(0.0203)
TERM	0.829 _(0.0422)	0.928 _(0.0197)	1.854 _(0.7437)	1.069 _(0.1001)	0.879 _(0.0178)	0.875 _(0.0711)	∞	∞
SEVER	<u>0.533</u> _(0.0347)	<u>0.592</u> _(0.0548)	<u>0.915</u> _(0.0343)	<u>0.841</u> _(0.0413)	<u>0.526</u> _(0.0287)	<u>0.720</u> _(0.1147)	<u>1.172</u> _(0.0542)	<u>1.215</u> _(0.0536)
SUBQUANTILE	0.396 _(0.0216)	0.442 _(0.0468)	0.808 _(0.0389)	0.827 _(0.0216)	0.446 _(0.1230)	0.456 _(0.1055)	1.074 _(0.0378)	1.132 _(0.0892)
Oracle ERM	∞	∞	∞	∞	∞	∞	∞	∞

Table 1: Boston Housing, Concrete Data, Wine Quality, and Drug and Polynomial Synthetic Dataset. $R = 10000$ for all datasets. Label Noise: $y_{\text{noise}} \sim \mathcal{N}(5, 5)$. Feature Noise: $y_{\text{noise}} = 10000 y_{\text{original}}$ and $x_{\text{noise}} = 100 x_{\text{original}}$. Polynomial Regression Synthetic Dataset. 1000 samples, $x \sim \mathcal{N}(0, 1)$, $y \sim \mathcal{N}(\sum_{i=0} a_i x^i, 0.01)$ where $a_i \sim \mathcal{N}(0, 1)$. The Radial Basis Function is used in first three experiments and polynomial kernel with degree 3 and $C = 1$ is used in the last experiment.

⁴In Progress

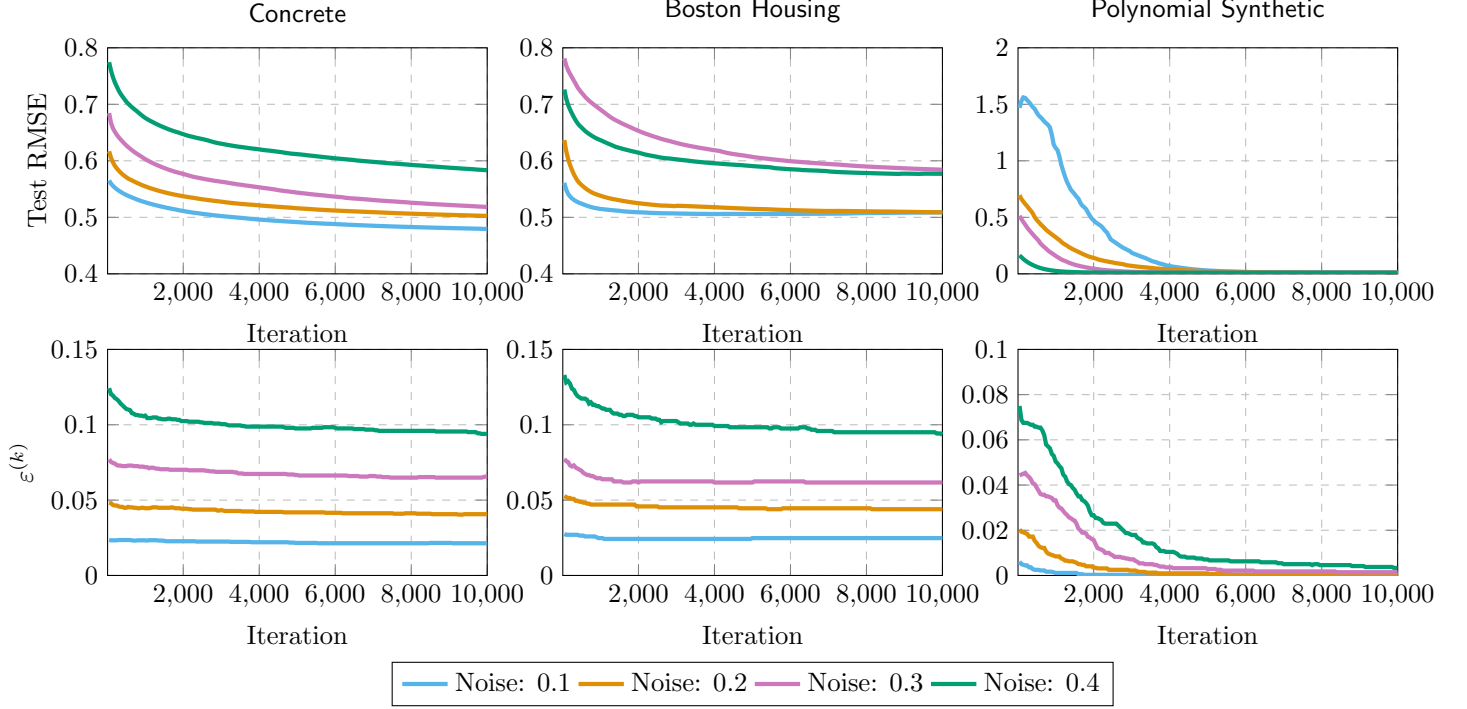


Figure 1: Test RMSE over the iterations in Concrete, Boston Housing, and Polynomial Datasets for SUBQUANTILE at different noise levels

In Figure 1, we see the final subquantile has significantly less outliers than the original corruption in the data set. Furthermore, we see there is a greater decrease in the higher outlier settings.

4.1 Linear Regression

In this section, we give experimental results for datasets using the linear kernel. This section will serve as a comparison to the various Robust Linear Regression Algorithms developed which are not meta-algorithms.

Algorithms	Test RMSE							
	Boston Housing		Wine Quality		Concrete		Drug	
	Label(↓)	Label+Feature	Label	Label+Feature	Label	Label+Feature	Label	Label+Feature
KRR	0.907 _(0.2724)	90.799 _(5.7170)	0.894 _(0.0404)	62.913 _(7.4959)	0.825 _(0.0943)	77.383 _(5.5692)	2.679 _(0.1286)	141.690 _(3.5297)
RANSAC	1.167 _(0.6710)	22.460 _(19.1987)	1.489 _(0.2730)	39.630 _(13.0294)	0.870 _(0.2308)	23.629 _(16.1023)	2.801 _(0.2004)	117.389 _(8.3915)
CRR	0.636 _(0.0905)	88.626 _(5.7380)	0.818 _(0.0224)	58.488 _(3.5612)	0.710 _(0.0919)	73.932 _(4.7867)	1.887 _(0.1463)	152.827 _(6.6038)
STIR	<u>0.562</u> _(0.0626)	78.878 _(8.0164)	0.828 _(0.0293)	58.352 _(4.6700)	<u>0.684</u> _(0.0245)	76.555 _(4.5927)	1.721 _(0.1520)	144.975 _(5.4953)
SEVER	0.601 _(0.0979)	5.980 _(8.2603)	0.814 _(0.0207)	9.065 _(13.7632)	<u>0.684</u> _(0.0438)	4.119 _(8.2436)	1.469 _(0.1162)	156.043 _(4.5543)
TERM	0.608 _(0.1357)	<u>0.569</u> _(0.0620)	0.840 _(0.0563)	<u>0.827</u> _(0.0255)	0.780 _(0.0734)	<u>0.808</u> _(0.0726)	<u>1.185</u> _(0.1077)	1.147 _(0.1258)
SUBQUANTILE	0.503 _(0.0470)	0.548* _(0.0286)	0.813 _(0.0357)	0.821 _(0.0305)	0.632 _(0.0275)	0.703 _(0.0427)	1.074 _(0.1848)	<u>2.413</u> _(0.6737)
Oracle ERM	0.630 _(0.1015)	0.665 _(0.1134)	0.838 _(0.0130)	0.865 _(0.0222)	0.763 _(0.0390)	0.768 _(0.0181)	0.988 _(0.0823)	0.985 _(0.0838)

Table 2: For only Label Noise, $y_{\text{noisy}} \sim \mathcal{N}(5, 5)$. For Label and Feature Noise $x_{\text{noisy}} = 100x_{\text{original}}$ and $y_{\text{noisy}} = 10000y_{\text{original}}$. * As indicated by the theory, when encountering feature noise, we require more gradient descent iterations to achieve the same bound between the returned point and the stationary point. Therefore, we train the label noise perturbed dataset for 10000 iterations, and the feature noise perturbed dataset for 100000 iterations.

4.2 Kernel Binary Classification

In this section we will give the algorithm for subquantile minimization for the kernel classification problem and then give some experimental results on state of the art datasets comparing against other state of the art robust algorithms.

Algorithms	Test Accuracy							
	Heart Disease				Breast Cancer			
	Label		Label+Feature		Label		Label+Feature	
	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$
SVM	0.777 _(0.0396)	0.639 _(0.0762)	0.534 _(0.0766)	0.538 _(0.0626)	0.926 _(0.0331)	0.548 _(0.1194)	0.649 _(0.0254)	0.618 _(0.0507)
SEVER	0.793 _(0.0422)	0.695 _(0.0636)	0.784 _(0.0432)	0.816 _(0.0562)	0.904 _(0.0356)	0.575 _(0.1456)	0.956 _(0.0164)	0.974 _(0.0062)
TERM	0.741 _(0.0393)	0.620 _(0.0699)	0.803 _(0.0613)	0.810 _(0.0286)	0.940 _(0.0378)	0.763 _(0.0364)	0.986 _(0.0143)	0.986 _(0.0119)
SUBQUANTILE	0.803 _(0.0293)	0.790 _(0.0350)	0.833 _(0.0318)	0.807 _(0.0468)	0.928 _(0.0129)	0.916 _(0.0185)	0.972 _(0.0187)	0.963 _(0.0170)
Oracle ERM	∞	∞	∞	∞	∞	∞	∞	∞

Table 3: Heart Disease and Breast Cancer Dataset. Label Noise: $y_{\text{noise}} = \mathbb{I}\{y_{\text{original}} = 0\}$. Feature Noise: $x_{\text{noise}} = 100x_{\text{original}}$. The Linear Kernel is used in all experiments.

4.3 Kernel Multi-Class Classification

In this section we will provide some experimental results on the multi-class classification task.

Algorithms	Test Accuracy							
	Iris		Glass		Wine		Satimage	
	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$
SVC	0.977 _(0.0300)	0.757 _(0.1155)	0.553 _(0.0969)	0.435 _(0.0721)	0.928 _(0.0484)	0.678 _(0.1368)	0.882 _(0.0056)	0.732 _(0.0168)
TERM	∞	∞	∞	∞	∞	∞	∞	∞
SEVER	∞	∞	∞	∞	∞	∞	∞	∞
SUBQUANTILE	0.987 _(0.0163)	0.820 _(0.1720)	0.656 _(0.0804)	0.598 _(0.0889)	0.975 _(0.0262)	0.867 _(0.1971)	0.899 _(0.0076)	0.861 _(0.0297)
Oracle ERM	∞	∞	∞	∞	∞	∞	∞	∞

Table 4: Iris ($R = 1$), Glass ($R = 10$), Wine ($R = 100$), and Satimage ($R = 10000$) Datasets. Label Noise is a randomly chosen incorrect label. Feature Noise: $y_{\text{noise}} = 10000y_{\text{original}}$ and $x_{\text{noise}} = 100x_{\text{original}}$. The Radial Basis Function is used in all experiments.

Results. We will clearly state our main findings.

- **Label Noise vs. Label and Feature Noise.** As suggested by our developed theory, for linear regression or using unbounded kernels, a large multiplicative term increases β and therefore requires more gradient descent iterations to achieve the same distance from a Moreau stationary point. Therefore, from simply increasing the number of gradient descent iterations, we are able to achieve similar RMSE in practice. This happens because the distance from a stationary point and the optimal is not affected by feature noise. This is one of the strengths of our theoretical analysis.
- **Error vs. ϵ .** We find approximately linear increase in the error with increasing ϵ . This can be seen in the γ term, which is upper bounded $\sqrt{\epsilon/(1-2\epsilon)}$. When $\epsilon \rightarrow 0.5$, the denominator approaches 0 and therefore our worst case bound increases.
- **Kernel.** Our error bounds are stronger when the dimension of the kernel is lower, i.e. we need more data to obtain the same error bounds. However, in practice, we find many datasets are better approximated by polynomial or RBF kernels, and therefore the γ term is significantly lower.

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. We also present accelerated gradient methods for the two-step algorithm to solve the nonconvex-concave optimization problem and give novel theoretical bounds.

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, \mathbb{P} , in the Gaussian Design. In expectation and with high probability, given sufficient data dependent on the kernel, we obtain a near minimax optimal error bound for a general positive definite continuous kernel. Our theoretical analysis is novel in that it utilizes the Moreau Envelope from a min-max formulation of the iterative thresholding algorithm.

Experiments. From our experiments, we see Subquantile Minimization is competitive with algorithms developed solely for robust linear regression as well as other meta-algorithms. Our theoretical analysis is through the lens of kernel-learning, but the generalization to linear regression from a non-kernel perspective can be done. In kernelized regression, we see SUBQUANTILE is the strongest of the meta-algorithms. Furthermore, in binary and multi-class classification, SUBQUANTILE is very strong. Thus, we can see empirically SUBQUANTILE is the strongest meta-algorithm across all kernelized regression and classification tasks and also the strongest algorithm in linear regression.

Interpretability. One of the strengths in Subquantile Optimization is the high interpretability. Once training is finished, we can see the $n(1-p)$ points with highest error to find the outliers and the features follow Gaussian Design. Furthermore, there is only hyperparameter p , which should be chosen to be approximately the percentage of inliers in the data and thus is not very difficult to tune for practical purposes. Our theory suggests for a problem where the amount of corruptions is unknown,

General Assumptions. The general assumption is the majority of the data should inliers. This is not a very strong assumption, as by the definition of outlier it should be in the minority. Furthermore, we assume the feature maps have a Gaussian Design. Such a design in many prior works in kernel learning and we therefore find it suitable.

The analysis of Subquantile Minimization can be extended to neural networks. This generalization will be appear in subsequent work.

References

- [1] 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](#), [2](#)
- [2] 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](#)
- [3] Kush Bhatia, Prateek Jain, and Purushottam Kar. Robust regression via hard thresholding. *Advances in neural information processing systems*, 28, 2015. [1](#)
- [4] 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](#)
- [5] Yu Cheng, Ilias Diakonikolas, Rong Ge, and David P Woodruff. Faster algorithms for high-dimensional robust covariance estimation. In *Conference on Learning Theory*, pages 727–757. PMLR, 2019. [1](#)
- [6] 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. [19](#)

- [7] Damek Davis and Dmitriy Drusvyatskiy. Stochastic model-based minimization of weakly convex functions. *SIAM Journal on Optimization*, 29(1):207–239, 2019. 8
- [8] 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
- [9] Ilias Diakonikolas and Daniel M Kane. *Algorithmic high-dimensional robust statistics*. Cambridge University Press, 2023. 1
- [10] 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
- [11] 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
- [12] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010. 8
- [13] Arthur Gretton. Introduction to rkhs, and some simple kernel algorithms. *Adv. Top. Mach. Learn. Lecture Conducted from University College London*, 16(5-3):2, 2013. 15
- [14] Arthur Gretton. Introduction to rkhs, and some simple kernel algorithms. *Adv. Top. Mach. Learn. Lecture Conducted from University College London*, 16(5-3):2, 2013. 18
- [15] Shu Hu, Yiming Ying, Siwei Lyu, et al. Learning by minimizing the sum of ranked range. *Advances in Neural Information Processing Systems*, 33:21013–21023, 2020. 1, 2
- [16] Peter J. Huber and Elvezio Ronchetti. *Robust statistics*. Wiley series in probability and statistics. Wiley, Hoboken, N.J., 2nd ed. edition, 2009. 1
- [17] 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
- [18] Ashish Khetan, Zachary C. Lipton, and Anima Anandkumar. Learning from noisy singly-labeled data. In *International Conference on Learning Representations*, 2018. 1
- [19] Yassine Laguel, Krishna Pillutla, Jérôme Malick, and Zaid Harchaoui. Superquantiles at work: Machine learning applications and efficient subgradient computation. *Set-Valued and Variational Analysis*, 29(4):967–996, Dec 2021.
- [20] 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. 1
- [21] Tian Li, Ahmad Beirami, Maziar Sanjabi, and Virginia Smith. Tilted empirical risk minimization. In *International Conference on Learning Representations*, 2021. 1
- [22] James Mercer. Xvi. functions of positive and negative type, and their connection the theory of integral equations. *Philosophical transactions of the royal society of London. Series A, containing papers of a mathematical or physical character*, 209(441-458):415–446, 1909. 24
- [23] Jean-Jacques Moreau. Proximité et dualité dans un espace hilbertien. *Bulletin de la Société mathématique de France*, 93:273–299, 1965. 4
- [24] 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

- [25] 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
- [26] Adarsh Prasad, Sivaraman Balakrishnan, and Pradeep Ravikumar. A unified approach to robust mean estimation. *arXiv preprint arXiv:1907.00927*, 2019. 1
- [27] Adarsh Prasad, Arun Sai Suggala, Sivaraman Balakrishnan, and Pradeep Ravikumar. Robust estimation via robust gradient estimation. *arXiv preprint arXiv:1802.06485*, 2018. 1
- [28] 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. 2
- [29] Ralph Tyrell Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, 1970. 8
- [30] R.T. Rockafellar, J.O. Royset, and S.I. Miranda. Superquantile regression with applications to buffered reliability, uncertainty quantification, and conditional value-at-risk. *European Journal of Operational Research*, 234(1):140–154, 2014.
- [31] Gabriele Santin and Robert Schaback. Approximation of eigenfunctions in kernel-based spaces. *Advances in Computational Mathematics*, 42:973–993, 2016. 19
- [32] Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry. Adversarially robust generalization requires more data. *Advances in neural information processing systems*, 31, 2018. 6
- [33] 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. 2
- [34] Hermann Weyl. Das asymptotische verteilungsgesetz der eigenwerte linearer partieller differentialgleichungen (mit einer anwendung auf die theorie der hohlraumstrahlung). *Mathematische Annalen*, 71(4):441–479, 1912. 17
- [35] Rahul Yedida, Snehanishu Saha, and Tejas Prashanth. Lipschitzlr: Using theoretically computed adaptive learning rates for fast convergence. *Applied Intelligence*, 51:1460–1478, 2021. 5

A Concentration Inequalities

In this section we will give various concentration inequalities on the inlier data for functions in the Reproducing Kernel Hilbert Space. We will first give our assumptions for robust kernelized regression.

Assumption 25 (Gaussian Design). We assume for $x_i \sim \mathbb{P} \in \mathcal{X}$, then it follows for the feature map, $\phi(\cdot) : \mathcal{X} \rightarrow \mathcal{H}$,

$$\phi(x_i) \sim \mathcal{N}(0, \Sigma)$$

where Σ is a possibly infinite dimensional covariance operator.

Assumption 26 (Normal Residuals). The residual is defined as $\mu_i \triangleq f_w^*(x_i) - y_i$. Then we assume for some $\sigma > 0$, it follows

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

Lemma 27 (Maximum of Gaussians). Let $\mu_1, \dots, \mu_n \sim \mathcal{N}(0, \sigma^2)$ for some $\sigma > 0$. Then it follows

$$\mathbb{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \max_{i \in [n]} |\mu_i| \leq O\left(\sigma \sqrt{\log n}\right)$$

Proof. We will integrate over the CDF to make our claim.

$$\begin{aligned} \mathbb{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \max_{i \in [n]} |\mu_i| &= \int_0^\infty \mathbb{P}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \left\{ \max_{i \in [n]} |\mu_i| > t \right\} dt \stackrel{(i)}{\leq} c_1 + n \int_{c_1}^\infty \mathbb{P}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \{|\mu_i| \geq t\} dt \\ &\stackrel{(ii)}{=} c_1 + 2n \int_{c_1}^\infty \mathbb{P}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \{\mu_i \geq t\} dt = c_1 + 2n \int_{c_1}^\infty \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{t}{\sigma}\right)^2} dx \\ &\leq c_1 + \frac{n}{\sigma} \sqrt{\frac{2}{\pi}} \int_{c_1}^\infty \frac{x}{c_1} e^{-\frac{1}{2}\left(\frac{x}{\sigma}\right)^2} dx = c_1 + n\sigma \sqrt{\frac{2}{\pi}} \frac{e^{-\left(\frac{c_1}{\sigma\sqrt{2}}\right)^2}}{c_1} \end{aligned}$$

(i) follows from a union bound and noting for a i.i.d sequence of random variables $\{X_i\}_{i \in [n]}$ and a constant C , it follows $\mathbb{P}\{\max_{i \in [n]} X_i \geq C\} = n\mathbb{P}\{X \geq C\}$ where X is sampled from the same distribution as each X_i .

(ii) follows from the symmetricity of the Gaussian distribution. From here, we choose $c_1 \triangleq \sigma \sqrt{2 \log n}$. Then we have,

$$\mathbb{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \max_{i \in [n]} |\mu_i| \leq \sigma \sqrt{2 \log n} + \frac{1}{\sqrt{\pi \log n}}$$

This completes the proof. ■

Lemma 28 (⁵Expected Maximum P_k). Let $x_i \sim \mathbb{P}$ such that $\phi(x_i) \sim \mathcal{N}(0, \Sigma)$ from Assumption 25. Then it follows

$$\mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \left[\max_{i \in [n]} k(x_i, x_i) \right] \leq O(\log n \operatorname{Tr}(\Sigma))$$

Proof. We once integrate over the CDF to make our claim.

$$\begin{aligned} \mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \left[\max_{i \in [n]} k(x_i, x_i) \right] &= \int_0^\infty \mathbb{P}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \left\{ \max_{i \in [n]} k(x_i, x_i) \geq t \right\} dt \\ &\stackrel{(i)}{\leq} c_2 + n \int_{c_2}^\infty \mathbb{P}_{\phi(x) \sim \mathcal{N}(0, \Sigma)} \{k(x, x) \geq t\} dt \\ &\stackrel{(ii)}{\leq} c_2 + n \int_{c_2}^\infty \frac{\operatorname{Tr}(\Sigma)}{t} dt \\ &= c_2 + \end{aligned}$$

⁵In Progress

Lemma 29 (Norm of Functions with Gaussian Design in the Reproducing Kernel Hilbert Space). *Let $x_i \sim \mathbb{P}$ such that $\phi(x_i) \sim \mathcal{N}(0, \Sigma)$ from Assumption 25 and Assumption 26. Then, it follows*

$$\mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \left\| \sum_{i=1}^n \mu_i \phi(x_i) \right\|_{\mathcal{H}} \leq O\left(\sigma \sqrt{n \log n \operatorname{Tr}(\Sigma)}\right)$$

Proof. Our proof follows standard ideas from High-Dimensional Probability. Let ξ_i for $i \in [n]$ denote i.i.d Rademacher variables such that for $\xi_i \sim \mathcal{R}$, it follows $\mathbb{P}\{\xi_i = 1\} = \mathbb{P}\{\xi_i = -1\} = \frac{1}{2}$. We then have,

$$\begin{aligned} & \mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \left\| \sum_{i=1}^n \mu_i \phi(x_i) \right\|_{\mathcal{H}} \leq \mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \max_{i \in [n]} |\mu_i| \left\| \sum_{i=1}^n \phi(x_i) \right\|_{\mathcal{H}} \\ & \stackrel{\text{lem. 27}}{\leq} O\left(\sigma \sqrt{\log n}\right) \mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(x_i) \right\|_{\mathcal{H}} \\ & \stackrel{(i)}{\leq} O\left(\sigma \sqrt{\log n}\right) \left(\mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(x_i) \right\|_{\mathcal{H}}^2 \right)^{1/2} \\ & \stackrel{(ii)}{=} O\left(\sigma \sqrt{\log n}\right) \left(\mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\langle \sum_{i=1}^n \xi_i \phi(x_i), \sum_{j=1}^n \xi_j \phi(x_j) \right\rangle_{\mathcal{H}} \right)^{1/2} \\ & = O\left(\sigma \sqrt{\log n}\right) \left(\mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \sum_{i=1}^n \sum_{j=1}^n \xi_i \xi_j k(x_i, x_j) \right)^{1/2} \\ & = O\left(\sigma \sqrt{\log n}\right) \left(\mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \sum_{i=1}^n k(x_i, x_i) \right)^{1/2} = O\left(\sigma \sqrt{n \log n}\right) \left(\mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} [k(x_i, x_i)] \right)^{1/2} \\ & = O\left(\sigma \sqrt{n \log n \operatorname{Tr}(\Sigma)}\right) \end{aligned}$$

(i) follows from Jensen's Inequality. (ii) follows from the reproducing property [13]. ■

Lemma 30 (Infinite Dimensional Covariance Estimation in the Hilbert-Schmidt Norm). *Let $\Sigma \triangleq \mathbb{E}_{\phi(x_i) \sim \mathbb{P}} [\phi(x_i) \otimes \phi(x_i)]$. Then let x_1, \dots, x_n be i.i.d sampled from \mathbb{P} such that $\phi(x_i) \sim \mathcal{N}(0, \Sigma)$ from Assumption 25, we then have*

$$\mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \left\| \frac{1}{n} \sum_{i=1}^n \phi(x_i) \otimes \phi(x_i) - \Sigma \right\|_{\text{HS}} \leq O\left(n^{-1/2} \operatorname{Tr}(\Sigma)\right)$$

Proof. Our proof follows standard ideas from High-Dimensional Probability. Let ξ_i for $i \in [n]$ denote i.i.d Rademacher variables such that for $\xi_i \sim \mathcal{R}$, it follows $\mathbb{P}\{\xi_i = 1\} = \mathbb{P}\{\xi_i = -1\} = \frac{1}{2}$. We then have,

$$\begin{aligned} & \mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \left\| \frac{1}{n} \sum_{i=1}^n \phi(x_i) \otimes \phi(x_i) - \Sigma \right\|_{\text{HS}} \\ & \stackrel{(i)}{\leq} \mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\tilde{\phi}(x_i) \sim \mathcal{N}(0, \Sigma)} \left\| \frac{1}{n} \sum_{i=1}^n \left(\phi(x_i) \otimes \phi(x_i) - \tilde{\phi}(x_i) \otimes \tilde{\phi}(x_i) \right) \right\|_{\text{HS}} \\ & = \mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\tilde{\phi}(x_i) \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \frac{1}{n} \sum_{i=1}^n \xi_i \left(\phi(x_i) \otimes \phi(x_i) - \tilde{\phi}(x_i) \otimes \tilde{\phi}(x_i) \right) \right\|_{\text{HS}} \\ & \stackrel{(ii)}{\leq} \frac{2}{n} \mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(x_i) \otimes \phi(x_i) \right\|_{\text{HS}} \\ & \stackrel{(iii)}{\leq} \frac{2}{n} \left(\mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(x_i) \otimes \phi(x_i) \right\|_{\text{HS}}^2 \right)^{1/2} \end{aligned}$$

(i) follows from noticing $\phi(x_i) \otimes \phi(x_i) - \Sigma$ is a mean 0 operator in $\mathcal{H} \otimes \mathcal{H}$, then for $X, Y \in \mathcal{H} \otimes \mathcal{H}$ s.t. $\mathbb{E}[Y] = 0$ it follows $\|X\|_{\text{HS}} = \|X - \mathbb{E}[Y]\|_{\text{HS}} = \|\mathbb{E}_Y[X - Y]\|_{\text{HS}}$ and finally applying Jensen's Inequality. (ii) follows from the triangle inequality. (iii) follows from Jensen's Inequality. Let e_k for $k \in [p]$ represent an orthonormal basis for the Hilbert Space \mathcal{H} . By expanding out the Hilbert-Schmidt Norm, we then have

$$\begin{aligned}
& \frac{2}{n} \left(\mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(x_i) \otimes \phi(x_i) \right\|_{\text{HS}}^2 \right)^{1/2} \\
&= \frac{2}{n} \left(\mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \sum_{k=1}^p \left\langle \sum_{i=1}^n \xi_i \phi(x_i) \otimes \phi(x_i) e_k, \sum_{j=1}^n \xi_j \phi(x_j) \otimes \phi(x_j) e_k \right\rangle \right)^{1/2} \\
&= \frac{2}{n} \left(\mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \sum_{k=1}^p \sum_{i=1}^n \sum_{j=1}^n \xi_i \xi_j \langle \phi(x_i) \otimes \phi(x_i) e_k, \phi(x_j) \otimes \phi(x_j) e_k \rangle \right)^{1/2} \\
&\stackrel{(iv)}{\leq} \frac{2}{n} \left(\mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \sum_{k=1}^p \sum_{i=1}^n \langle \phi(x_i) \otimes \phi(x_i) e_k, \phi(x_i) \otimes \phi(x_i) e_k \rangle \right)^{1/2} \\
&= \frac{2}{n} \left(\sum_{i=1}^n \mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \|\phi(x_i) \otimes \phi(x_i)\|_{\text{HS}}^2 \right)^{1/2} \stackrel{(v)}{=} \frac{2}{n} \left(\sum_{i=1}^n \mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} \|\phi(x_i)\|_{\mathcal{H}}^4 \right)^{1/2} \\
&= \frac{2}{n} \left(\sum_{i=1}^n \mathbb{E}_{\phi(x_i) \sim \mathcal{N}(0, \Sigma)} [k^2(x_i, x_i)] \right)^{1/2} = \frac{2}{\sqrt{n}} \left(2 \text{Tr}(\Sigma^2) + \text{Tr}(\Sigma)^2 \right)^{1/2} \leq 2\sqrt{3}n^{-1/2} \text{Tr}(\Sigma)
\end{aligned}$$

(iv) follows from noticing $\mathbb{E}_{\xi_i, \xi_j \sim \mathcal{R}} [\xi_i \xi_j] = \delta_{i,j}$. (v) follows from expanding the Hilbert-Schmidt Norm and applying Parseval's Identity. We note $\text{Tr}(\Sigma) < \infty$ and therefore even though the covariance operator is infinite-dimensional we are able to get a finite bound on the covariance approximation. This completes the proof. \blacksquare

Lemma 31 (Finite Dimensional Covariate Estimation in the Spectral Norm). *Let $x_1, \dots, x_n \sim \mathcal{N}(0, \Sigma)$. It then follows,*

$$\mathbb{E}_{x_i \sim \mathcal{N}(0, \Sigma)} \left\| \frac{1}{n} \sum_{i=1}^n x_i x_i^\top - \Sigma \right\|_2 \leq \frac{2\sqrt{3} \text{Tr}(\Sigma)}{\sqrt{n}}$$

Proof. From similar steps in Lemma 30. We have,

$$\begin{aligned}
\mathbb{E}_{x_i \sim \mathcal{N}(0, \Sigma)} \left\| \frac{1}{n} \sum_{i=1}^n x_i x_i^\top - \Sigma \right\|_2 &\leq \frac{2}{n} \left(\mathbb{E}_{x_i \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i x_i x_i^\top \right\|_2^2 \right)^{1/2} \\
&= \frac{2}{n} \left(\mathbb{E}_{x_i \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \sum_{j=1}^n \xi_i \xi_j x_i x_i^\top x_j x_j^\top \right\|_2 \right)^{1/2} \\
&\leq \frac{2}{n} \left(\mathbb{E}_{x_i \sim \mathcal{N}(0, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \text{Tr} \left(\sum_{i=1}^n \sum_{j=1}^n \xi_i \xi_j x_i x_i^\top x_j x_j^\top \right) \right)^{1/2} \\
&= \frac{2}{n} \left(\mathbb{E}_{x_i \sim \mathcal{N}(0, \Sigma)} \sum_{i=1}^n \|x_i\|^4 \right)^{1/2} = \frac{2\sqrt{\text{Tr}(\Sigma)^2 + 2 \text{Tr}(\Sigma^2)}}{\sqrt{n}}
\end{aligned}$$

Noting $\text{Tr}(\Sigma^2) \leq \text{Tr}(\Sigma)^2$ as $\Sigma \succ 0$ and applying the triangle inequality completes the proof. \blacksquare

B Proofs for Section 3

In this section we give the deferred proofs of our main structural results.

B.1 Proof of Lemma 3

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

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

Appendix B.1 is equal to 0 when $t = \hat{\nu}_{np}$. This is equivalent to the p -quantile of the Risk. ■

B.2 Proof of Lemma 5

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

$$\begin{aligned} \nabla_f g(t^{(k+1)}, f_w^{(k)}) &= \nabla_f \left(\hat{\nu}_{np} - \frac{1}{np} \sum_{i=1}^n \left(\hat{\nu}_{np} - \ell(x_i; f_w^{(k)}, y_i) \right)^+ \right) \\ &= -\frac{1}{np} \sum_{i=1}^{np} \nabla_f \left(\hat{\nu}_{np} - \ell(x_i; f_w^{(k)}, y_i) \right)^+ = \frac{1}{np} \sum_{i=1}^{np} \nabla_f \ell(x_i; f_w^{(k)}, y_i) \begin{cases} 1 & \text{if } t > \hat{\nu}_i \\ 0 & \text{if } t < \hat{\nu}_i \\ [0, 1] & \text{if } t = \hat{\nu}_i \end{cases} \end{aligned}$$

Now we note $\nu_{np} \leq t^{(k+1)} \leq \nu_{np+1}$

$$\nabla_f g(t^{(k+1)}, f_w^{(k)}) = \frac{1}{np} \sum_{i=1}^{np} \nabla_f \ell(x_i; f_w^{(k)}, y_i)$$

This concludes the proof. ■

B.3 Proof of Lemma 6

Proof. We will upper bound the operator norm of the Hessian. Let $\vartheta \triangleq \sigma(\lambda(t - \ell(f_w; x_i, y_i)))$, we then have

$$\begin{aligned} \nabla_f^2 \tilde{g}_\lambda(t, f_w) &= \nabla_f \left(\frac{1}{np} \sum_{i=1}^n \vartheta \nabla_f \ell(f_w; x_i, y_i) \right) \\ &= \frac{1}{np} \sum_{i=1}^n \left(\vartheta \nabla_f^2 \ell(f_w; x_i, y_i) - \vartheta(1 - \vartheta) (\nabla_f \ell(f_w; x_i, y_i) \otimes \nabla_f \ell(f_w; x_i, y_i)) \right) \end{aligned}$$

Now we will upper bound the operator norm of the Hessian.

$$\begin{aligned} &\lim_{\lambda \rightarrow \infty} \sup_{f_w \in \mathcal{K}} \left\| \nabla_f^2 \tilde{g}_\lambda(t, f_w) \right\|_{\text{op}} \|f_w - f_{\bar{w}}\|_{\mathcal{H}} \\ &\stackrel{\text{appendix B.3}}{=} \lim_{\lambda \rightarrow \infty} \sup_{f_w} \left\| \frac{1}{np} \sum_{i=1}^n \vartheta \nabla_f^2 \ell(f_w; x_i, y_i) - \vartheta(1 - \vartheta) \nabla_f \ell(f_w; x_i, y_i) \otimes \nabla_f \ell(f_w; x_i, y_i) \right\|_{\text{op}} \|f_w - f_{\bar{w}}\|_{\mathcal{H}} \\ &\stackrel{(i)}{\leq} \lim_{\lambda \rightarrow \infty} \sup_{f_w \in \mathcal{K}} \frac{1}{np} \sum_{i=1}^n \left\| \vartheta \nabla_f^2 \ell(f_w; x_i, y_i) \right\|_{\text{op}} \|f_w - f_{\bar{w}}\|_{\mathcal{H}} \stackrel{(ii)}{\leq} \sup_{f_w \in \mathcal{K}} \frac{1}{np} \sum_{i=1}^n \left\| \nabla_f^2 \ell(f_w; x_i, y_i) \right\|_{\text{op}} \|f_w - f_{\bar{w}}\|_{\mathcal{H}} \end{aligned}$$

(i) follows from applying the Triangle Inequality and then Weyl's Inequality [34]. (ii) follows from noting $\vartheta \in (0, 1)$. We now note that removing ϑ also removes the dependence on λ which allows us to take the limit out of the expression. ■

B.4 Proof of Lemma 11

Proof. By the definition of stationary point, we have

$$\begin{aligned} f_{\hat{w}} &= \lim_{\lambda \rightarrow \infty} \arg \inf_{f_w \in \mathcal{K}} \left\{ \Phi_\lambda(f_w) + \frac{1}{2\rho} \|f_w - f_{\hat{w}}\|_{\mathcal{H}}^2 \right\} \\ &\stackrel{(i)}{=} \arg \inf_{f_w \in \mathcal{K}} \left\{ \lim_{\lambda \rightarrow \infty} \Phi_\lambda(f_w) + \frac{1}{2\rho} \|f_w - f_{\hat{w}}\|_{\mathcal{H}}^2 \right\} \end{aligned}$$

(i) holds as we ρ is independent of λ as shown in the proof of Lemma 6. This implies then for any $f_w \in \mathcal{K}$ and noting $\rho \leq \beta^{-1}$, it follows

$$\lim_{\lambda \rightarrow \infty} \Phi_\lambda(f_{\hat{w}}) \leq \lim_{\lambda \rightarrow \infty} \Phi_\lambda(f_w) + \beta \|f_w - f_{\hat{w}}\|_{\mathcal{H}}^2$$

where we choose $\rho \triangleq 1/(2\beta)$. We can then plug in the optimal, f_w^* for f_w and rearrange and we have the desired result. \blacksquare

B.5 Proof of Lemma 12

Proof. For any $f_{w_1}, f_{w_2} \in \mathcal{K}$, we will first show the gradient is bounded.

$$\begin{aligned} |g(t, f_{w_1}) - g(t, f_{w_2})| &= \left| \int_0^1 \nabla_f g(t, (1-\lambda)f_{w_1} + \lambda f_{w_2})(f_{w_1} - f_{w_2}) d\lambda \right| \\ &\leq \|f_{w_1} - f_{w_2}\|_{\mathcal{H}} \left| \int_0^1 \nabla_f g(t, (1-\lambda)f_{w_1} + \lambda f_{w_2}) d\lambda \right| \\ &\stackrel{(a)}{\leq} \|f_{w_1} - f_{w_2}\|_{\mathcal{H}} \max_{f_w \in \mathcal{K}} \|\nabla_f g(t, f_w)\|_{\mathcal{H}} \end{aligned}$$

In (a), we note that since \mathcal{K} is convex, then by definition as $f_{w_1}, f_{w_2} \in \mathcal{K}$, we have for $\lambda \in [0, 1]$, the convex combination $(1-\lambda)f_{w_1} + \lambda f_{w_2} \in \mathcal{K}$. We use the \mathcal{H} norm of the gradient to bound L from above for an element in the convex closed set \mathcal{K} .

$$\|\nabla g(t, f_w)\|_{\mathcal{H}} = \left\| \frac{2}{np} \sum_{i=1}^n \mathbb{I}\{t \geq (f_w(x_i) - y_i)^2\} (f_w(x_i) - y_i) \cdot k(x_i, \cdot) \right\|_{\mathcal{H}}$$

W.L.O.G, let x_1, x_2, \dots, x_m where $0 \leq m \leq n$, represent the data vectors such that $t \geq (f_w(x_i) - y_i)^2$.

$$\begin{aligned} &= \left\| \frac{2}{np} \sum_{i=1}^m (f_w(x_i) - y_i) \cdot k(x_i, \cdot) \right\|_{\mathcal{H}} \leq \frac{2}{np} \left(\left\| \sum_{i=1}^m f_w(x_i) \cdot k(x_i, \cdot) \right\|_{\mathcal{H}} + \left\| \sum_{i=1}^m y_i k(x_i, \cdot) \right\|_{\mathcal{H}} \right) \\ &\stackrel{(a)}{\leq} \frac{2}{np} \left(\left\| \sum_{i=1}^m \left\langle \sum_{j=1}^n w_j k(x_j, \cdot), k(x_i, \cdot) \right\rangle_{\mathcal{H}} \cdot k(x_i, \cdot) \right\|_{\mathcal{H}} + \left\| \sum_{i=1}^m y_i \left\| \sum_{i=1}^m k(x_i, \cdot) \right\|_{\mathcal{H}} \right\|_{\mathcal{H}} \right) \\ &\leq \frac{2}{np} \left(\left\| \sum_{j=1}^n w_j k(x_j, \cdot), \sum_{i=1}^m k(x_i, \cdot) \right\rangle_{\mathcal{H}} \left\| \sum_{i=1}^m k(x_i, \cdot) \right\|_{\mathcal{H}} + \left\| \sum_{i=1}^m y_i \left\| \sum_{i=1}^m \sqrt{k(x_i, x_i)} \right\|_{\mathcal{H}} \right) \right) \\ &\leq \frac{2}{np} \left(\|f_w\|_{\mathcal{H}} \left(\sum_{i=1}^m \sqrt{k(x_i, x_i)} \right)^2 + \sqrt{n} \|y\|_2 \left(\sum_{i=1}^n \sqrt{k(x_i, x_i)} \right) \right) \\ &\leq \frac{2R}{np} \left(\sum_{i=1}^n \sqrt{k(x_i, x_i)} \right)^2 + \frac{2\|y\|_2}{p\sqrt{n}} \left(\sum_{i=1}^n \sqrt{k(x_i, x_i)} \right) \end{aligned}$$

(a) follows from the reproducing property for RKHS [14]. If we have a normalized kernel such as the Gaussian Kernel, then we have the Lipschitz Constant is finite. Furthermore, if the adversary introduces label corruption that tends to ∞ , then these points will not be in the Subquantile as f_w has bounded norm, so it will have infinite error. This concludes the proof. \blacksquare

B.6 Proof of Lemma 14

Proof. Let S be the set containing the points with the minimum error from X w.r.t to the weights vector w . Define $\eta_i \triangleq (f_w^*(x_i) - y_i)$ where $i \in P$.

$$\begin{aligned}
\lim_{\lambda \rightarrow \infty} (\Phi_\lambda(f_w) - \Phi_\lambda(f_w^*)) &= \sum_{i \in S} (f_w(x_i) - y_i)^2 - \sum_{j \in P} (f_w^*(x_j) - y_j)^2 \\
&= \sum_{i \in S \cap P} (f_w(x_i) - y_i)^2 + \sum_{i \in S \cap Q} (f_w(x_i) - y_i)^2 - \sum_{j \in P} (f_w^*(x_j) - y_j)^2 \\
&\geq \sum_{i \in S \cap P} (f_w(x_i) - y_i)^2 - \sum_{j \in P} (f_w^*(x_j) - y_j)^2 = \sum_{i \in S \cap P} (f_w(x_i) - f_w^*(x_i) - \eta_i)^2 - \sum_{j \in P} \eta_j^2 \\
&\geq \sum_{i \in S \cap P} \underbrace{((f_w - f_w^*)(x_i))^2}_{A_1} - 2 \underbrace{\sum_{i \in S \cap P} \eta_i (f_w - f_w^*)(x_i)}_{A_2} - \underbrace{\sum_{j \in P \setminus S} \eta_j^2}_{A_3}
\end{aligned}$$

Now we will upper bound A_1 . Similar to [6] Let $\mathbb{E}_{x \sim \mathbb{P}}[\varphi(x) \otimes \varphi(x)] = \mathbb{I}_m$ where $\varphi(x) = \{\varphi(x)\}_{k=1}^m$ and m is possibly infinite. We can then rescale the basis features. Then let $\phi(x) = \Sigma^{1/2} \varphi(x)$. We therefore have $\Sigma = \mathbb{E}_{x \sim \mathbb{P}}[\phi(x) \otimes \phi(x)] = \text{diag}(\xi_1, \dots, \xi_n)$. This is the eigenfunction basis described in [31].

$$\begin{aligned}
A_1 &\triangleq \sum_{i \in S \cap P} ((f_w - f_w^*)(x_i))^2 \stackrel{(a)}{=} \sum_{i \in S \cap P} \left\langle \sum_{j \in X} (w_j - w_j^*) k(x_j, \cdot), k(x_i, \cdot) \right\rangle_{\mathcal{H}}^2 \\
&= \sum_{i \in S \cap P} \left\langle \sum_{j \in X} (w_j - w_j^*) \phi(x_j), \phi(x_i) \right\rangle_{\mathcal{H}} \left\langle \phi(x_i), \sum_{j \in X} (w_j - w_j^*) \phi(x_j) \right\rangle_{\mathcal{H}} \\
&= \sum_{i \in S \cap P} \left\langle \sum_{j \in X} (w_j - w_j^*) \phi(x_j), \phi(x_i) \otimes \phi(x_i) \sum_{j \in X} (w_j - w_j^*) \phi(x_j) \right\rangle_{\mathcal{H}} \\
&= \sum_{i \in S \cap P} \left\langle \phi(x) \otimes \phi(x), (f_w - f_w^*) \otimes (f_w - f_w^*) \right\rangle_{\text{HS}} \\
&= \sum_{i \in S \cap P} \left\langle \Sigma + \phi(x) \otimes \phi(x) - \Sigma, (f_w - f_w^*) \otimes (f_w - f_w^*) \right\rangle_{\text{HS}} \\
&\stackrel{\text{lem. 30}}{\geq} \left(n(1 - 2\varepsilon) \lambda_{\min}(\Sigma) - \left\| \sum_{i \in S \cap P} \phi(x) \otimes \phi(x) - \Sigma \right\|_{\text{HS}} \right) \|f_w - f_w^*\|_{\mathcal{H}}^2
\end{aligned}$$

Next we will upper bound A_2 ,

$$\begin{aligned}
A_2 &\triangleq \sum_{i \in S \cap P} \eta_i (f_w - f_w^*)(x_i) = \sum_{i \in S \cap P} \left\langle \sum_{j \in X} (w_j - w_j^*) k(x_j, \cdot), \eta_i k(x_i, \cdot) \right\rangle_{\mathcal{H}} \\
&= \left\langle \sum_{j \in X} (w_j - w_j^*) k(x_j, \cdot), \sum_{i \in S \cap P} \eta_i k(x_i, \cdot) \right\rangle_{\mathcal{H}} \\
&\leq \|f_w - f_w^*\|_{\mathcal{H}} \left\| \sum_{i \in S \cap P} \eta_i k(x_i, \cdot) \right\|_{\mathcal{H}} = \|f_w - f_w^*\|_{\mathcal{H}} \left\| \sum_{i \in S \cap P} \eta_i \phi(x_i) \right\|_{\mathcal{H}}
\end{aligned}$$

Then, combining our bounds, we have

$$\begin{aligned}
\lim_{\lambda \rightarrow \infty} (\Phi_\lambda(f_w) - \Phi_\lambda(f_w^*)) &\stackrel{\text{appendix B.6}}{\geq} \eta^2 \left(n(1 - 2\varepsilon) \lambda_{\min} \left(\mathbb{E}_{x \sim \mathbb{P}} [\phi(x) \otimes \phi(x)] \right) \right. \\
&\quad \left. - \left\| \sum_{i \in S \cap P} \phi(x) \otimes \phi(x) - \Sigma \right\|_{\text{HS}} \right) - 2\eta \left\| \sum_{i \in S \cap P} \eta_i \phi(x_i) \right\|_{\mathcal{H}} - \sum_{j \in P \setminus S} \eta_j^2
\end{aligned}$$

This completes the proof. ■

B.7 Proof of Theorem 15

Proof. First, we give the definition of the Moreau stationary point.

$$\|\nabla \mathbf{M}_{\Phi_\lambda, \rho}(f_w)\|_{\mathcal{H}} = \left\| \frac{1}{\rho} \left(f_w - \arg \min_{f_{\hat{w}} \in \mathcal{K}} \left(\Phi(f_{\hat{w}}) + \frac{1}{2\rho} \|f_w - f_{\hat{w}}\|_{\mathcal{H}}^2 \right) \right) \right\|_{\mathcal{H}} = 0$$

This implies for any $f_{\hat{w}} \in \mathcal{K}$, it follows

$$\lim_{\lambda \rightarrow \infty} (\Phi_\lambda(f_{\hat{w}})) < \lim_{\lambda \rightarrow \infty} (\Phi_\lambda(f_{\hat{w}})) + \frac{1}{2\rho} \|f_{\hat{w}} - f_w^*\|_{\mathcal{H}}^2$$

For any $f_{\hat{w}}$ satisfying above, then the distance from the optimal must be low. Let $\tilde{w} = w^*$, then we have

$$\lim_{\lambda \rightarrow \infty} (\Phi_\lambda(f_{\hat{w}}) - \Phi_\lambda(f_w^*)) \leq \frac{1}{2\rho} \|f_{\hat{w}} - f_w^*\|_{\mathcal{H}}^2$$

We proceed by proof by contradiction. Assume $\|f_{\hat{w}} - f_w^*\| > \eta$, then if $\Phi(f_{\hat{w}}) - \Phi(f_w^*) > \frac{\eta^2}{2\rho}$, then we will have $f_{\hat{w}}$ is not a stationary point, which will imply $\|f_{\hat{w}} - f_w^*\|_{\mathcal{H}} \leq \eta$. Therefore, we attempt to find the minimum value for η . From Lemma 14, we have the expected distance from a stationary point of the Moreau Envelope from the optimal point over the distribution of uncorrupted datasets.

$$\begin{aligned} & \mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} \lim_{\lambda \rightarrow \infty} (\Phi(f_w) - \Phi(f_w^*)) \stackrel{\text{lem. 14}}{\geq} \eta^2 \left(n(1-2\varepsilon) \lambda_{\min} \left(\mathbb{E}_{x \sim \mathbb{P}} [\phi(x) \otimes \phi(x)] \right) \right. \\ & \quad \left. - \mathbb{E}_{x_i \sim \mathbb{P}} \left\| \sum_{i \in S \cap P} \phi(x_i) \otimes \phi(x_i) - \Sigma \right\|_{\text{HS}} \right) - 2\eta \mathbb{E}_{\mu_i, x_i \sim \mathbb{P}} \left\| \sum_{i \in S \cap P} \eta_i \phi(x_i) \right\| - \mathbb{E}_{\mu_j \sim \mathbb{P}} \sum_{j \in P \setminus S} \eta_j^2 \\ & \stackrel{\text{lems. 29 and 30}}{\geq} \eta^2 \left(n(1-2\varepsilon) \lambda_{\min}(\Sigma) - 2 \text{Tr}(\Sigma) \sqrt{n(1-2\varepsilon)} \right) - \eta O \left(\sigma \sqrt{n(1-2\varepsilon) \log(n(1-2\varepsilon)) \text{Tr}(\Sigma)} \right) - \sigma \varepsilon n \end{aligned}$$

From the definition of stationary point, we have

$$\eta^2 \left(n(1-2\varepsilon) \lambda_{\min}(\Sigma) - 2 \text{Tr}(\Sigma) \sqrt{n(1-2\varepsilon)} - \beta \right) - \eta O \left(\sigma \sqrt{n(1-2\varepsilon) \log(n(1-2\varepsilon)) \text{Tr}(\Sigma)} \right) - \sigma \varepsilon n \leq 0$$

Therefore, when Appendix B.7 does not hold, we have a contradiction. It thus follows from upper bounding the positive solution of the quadratic equation,

$$\begin{aligned} \eta & \leq (\sigma \varepsilon n)^{1/2} \left(n(1-2\varepsilon) \left(\lambda_{\min}(\Sigma) - \frac{2 \text{Tr}(\Sigma)}{\sqrt{n(1-2\varepsilon)}} \right) - \beta \right)^{-1/2} \\ & \quad + O \left(\sigma \sqrt{n(1-2\varepsilon) \log(n(1-2\varepsilon)) \text{Tr}(\Sigma)} \right) \left(n(1-2\varepsilon) \left(\lambda_{\min}(\Sigma) - \frac{2 \text{Tr}(\Sigma)}{\sqrt{n(1-2\varepsilon)}} \right) - \beta \right)^{-1} \end{aligned}$$

Then for some constant $c_1 \in (0, 1)$, if $n \geq \frac{8 \text{Tr}(\Sigma)^2}{\lambda_{\min}(\Sigma)(1-c_1)^2(1-2\varepsilon)} + \frac{8\beta}{(1-c_1)^2(1-2\varepsilon)}$, we have

$$\eta \leq \left(\frac{\sigma \varepsilon n}{c_1 n(1-2\varepsilon) \lambda_{\min}(\Sigma)} \right)^{1/2} + \frac{O \left(\sigma \sqrt{\log(n(1-2\varepsilon)) \text{Tr}(\Sigma)} \right)}{c_1 \sqrt{n(1-2\varepsilon) \lambda_{\min}(\Sigma)}}$$

we therefore see as n goes large, $c_1 \rightarrow 1$, and we have in the worst case

$$\mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} \|f_{\hat{w}} - f_w^*\|_{\mathcal{H}} \leq O \left(\sqrt{\frac{\varepsilon}{1-2\varepsilon}} \frac{\sigma}{\lambda_{\min}(\Sigma)} \right)$$

This completes the proof. ■

B.8 Proof of Corollary 16

We follow the same framework as our proof for kernelized linear regression, we will simply give the new constants. Assuming the uncorrupted covariates, $x_i \sim \mathcal{N}(f_W(x_i), \Sigma)$. To simplify notation, let us define $\tilde{n} \triangleq n(1 - 2\varepsilon)$ to represent the absolute minimum number of uncorrupted points in the Subquantile. We then have,

$$\begin{aligned} \mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} \lim_{\lambda \rightarrow \infty} (\Phi_\lambda(w) - \Phi_\lambda(w^*)) &\stackrel{\text{lem. 14}}{\geq} \eta^2 \left(\tilde{n} \lambda_{\min}(\Sigma) - \mathbb{E} \left\| \sum_{i \in S \cap P} x_i x_i^\top - \Sigma \right\|_2 \right) - \mathbb{E}_{\xi, \mu_i \sim \mathbb{P}} \left\| \sum_{i \in S \cap P} \mu_i x_i \right\|_2 - \mathbb{E}_{\mu_i \sim \mathbb{P}} \sum_{i \in P \setminus S} \mu_i^2 \\ &\stackrel{\text{lems. 27, 29 and 31}}{\geq} \eta^2 \left(\tilde{n} \lambda_{\min}(\Sigma) - \sqrt{\tilde{n}} \left(2\sqrt{3} \text{Tr}(\Sigma) \right) \right) - \eta O \left(\sigma \sqrt{\tilde{n} \log(\tilde{n}) \text{Tr}(\Sigma)} \right) - \varepsilon n \sigma^2 \end{aligned}$$

Then from a similar contradiction idea and upper bounding the quadratic, we have in expectation

$$\eta \stackrel{\text{thm. 15}}{\leq} O \left(\sigma \sqrt{\tilde{n} \log(\tilde{n}) \text{Tr}(\Sigma)} \right) \left(\tilde{n} \lambda_{\min}(\Sigma) - \sqrt{\tilde{n}} \left(2\sqrt{3} \text{Tr}(\Sigma) \right) - \beta \right)^{-1} + \sigma \sqrt{\tilde{n} \frac{\varepsilon}{1 - 2\varepsilon}} \left(\tilde{n} \lambda_{\min}(\Sigma) - \sqrt{\tilde{n}} \left(2\sqrt{3} \text{Tr}(\Sigma) \right) - \beta \right)^{-1/2}$$

We then have for a constant $c_2 \in (0, 1)$, if $n \geq \frac{54 \text{Tr}(\Sigma)}{(1 - c_2)^2 (1 - 2\varepsilon) \lambda_{\min}^2(\Sigma)} + 2\beta$, it follows

$$\eta \leq \sqrt{\frac{\sigma^2 \varepsilon}{(1 - 2\varepsilon) c_2 \lambda_{\min}(\Sigma)}} + \frac{O \left(\sigma \sqrt{\log(n(1 - 2\varepsilon)) \text{Tr}(\Sigma)} \right)}{\sqrt{n(1 - 2\varepsilon) c_2 \lambda_{\min}(\Sigma)}}$$

We thus see as n goes large, $c_2 \rightarrow 1$ and we will have in worst case,

$$\mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} \|\hat{w} - w^*\|_2 \leq O \left(\frac{\gamma \sigma}{\sqrt{\lambda_{\min}(\Sigma)}} \right)$$

where $\gamma \triangleq \sqrt{\frac{|P \setminus S|}{|S \cap P|}}$. Obtaining the same asymptotic bound as in the kernelized regression case. This completes the proof. \blacksquare

B.9 Proof of Lemma 17

Proof. We use the \mathcal{H} norm of the gradient to bound L from above. Let S be denoted as the subquantile set. Define the sigmoid function as $\sigma(x) = \frac{1}{1 + e^{-x}}$.

$$\begin{aligned} \|\nabla_f g(t, f_w)\|_{\mathcal{H}} &= \left\| \frac{1}{np} \sum_{i=1}^n \mathbb{I}\{t \geq (1 - y_i) \log(f_w(x_i))\} (y_i - \sigma(f_w(x_i))) \cdot k(x_i, \cdot) \right\|_{\mathcal{H}} \\ &\stackrel{(i)}{\leq} \frac{1}{np} \sum_{i \in S} \|(y_i - \sigma(f_w(x_i))) \cdot k(x_i, \cdot)\|_{\mathcal{H}} \stackrel{(ii)}{\leq} \frac{1}{np} \sum_{i \in S} |y_i - \sigma(f_w(x_i))| \|k(x_i, \cdot)\|_{\mathcal{H}} \stackrel{(iii)}{\leq} \frac{1}{np} \sum_{i=1}^n \sqrt{k(x_i, x_i)} \end{aligned}$$

(i) follows from the triangle inequality. (ii) follows from the Cauchy-Schwarz inequality. (iii) follows from the fact that $y_i \in \{0, 1\}$ and $\text{range}(\sigma) \in [0, 1]$. This completes the proof. \blacksquare

B.10 Proof of Lemma 18

We use the operator norm of second derivative to bound β from above. Let S be the subquantile set.

$$\begin{aligned} \|\nabla_f^2 g(t, f_w)\|_{\text{op}} &= \frac{1}{np} \sum_{i=1}^n \mathbb{I}\{t \geq (1 - y_i) \log(f_w(x_i))\} \sigma(f_w(x_i)) (1 - \sigma(f_w(x_i))) \|\phi(x_i) \otimes \phi(x_i)\|_{\text{op}} \\ &\leq \frac{1}{np} \sum_{i=1}^n |\sigma(f_w(x_i)) (1 - \sigma(f_w(x_i)))| \|\phi(x_i)\|_{\text{op}}^2 \stackrel{(i)}{\leq} \frac{1}{4np} \sum_{i=1}^n k(x_i, x_i) = \frac{1}{4np} \text{Tr}(\mathbf{K}) \end{aligned}$$

(i) follows as for a scalar $\alpha \in [0, 1]$, the maximum value of $\alpha(1 - \alpha)$ is obtained at $\frac{1}{4}$. This completes the proof. \blacksquare

B.11 Proof of Lemma 19

Proof.⁶ By the Lemma statement, we have $f_{\hat{w}}$ is a stationary point, i.e. $0 \in \partial\Phi(f_{\hat{w}})$. This implies for all $f_w \in \mathcal{K}$, we have $\Phi(f_{\hat{w}}) \leq \Phi(f_w)$. As Φ is differentiable, we have the first-order stationary condition, which is $\nabla\Phi(f_{\hat{w}})(f_{\hat{w}} - f_w) \leq 0$ or for all $w \in \mathcal{K}$. We assume $f_w^* \in \mathcal{K}$. Let S be the Subquantile set for $f_{\hat{w}}$. We will proceed by contradiction, assume $\|f_{\hat{w}} - f_w^*\|_{\mathcal{H}} \geq \eta$. Then, we have

$$\begin{aligned} (\nabla_f g(f_{\hat{w}}, t))(f_{\hat{w}} - f_w^*) &= (f_{\hat{w}} - f_w^*) \left(\sum_{i \in S} (\sigma(f_{\hat{w}}(x_i)) - y_i) \phi(x_i) \right) \\ &= (f_{\hat{w}} - f_w^*) \left(\sum_{i \in S} (\sigma(f_{\hat{w}}(x_i)) - \sigma(f_w^*(x_i)) + \sigma(f_w^*(x_i)) - y_i) \phi(x_i) \right) \\ &\stackrel{(i)}{\geq} \underbrace{(f_{\hat{w}} - f_w^*) \left(\sum_{i \in S \cap P} (\sigma(f_{\hat{w}}(x_i)) - \sigma(f_w^*(x_i))) \phi(x_i) \right)}_{B_1} + \underbrace{(f_{\hat{w}} - f_w^*) \left(\sum_{i \in S} (\sigma(f_w^*(x_i)) - y_i) \phi(x_i) \right)}_{B_2} \end{aligned}$$

(i) follows from noting $\sigma(\cdot)$ is a monotonically increasing function. Let us now consider the function $h : \mathcal{H} \rightarrow \mathbb{R}$ defined as $h(f_w) = \sum_{i \in S \cap P} \log(1 + \exp(f_w(x_i)))$. We then have $h'(f_w) = \sum_{i \in S \cap P} \sigma(f_w(x_i)) \phi(x_i)$, from which we have $h''(f_w) = \sum_{i \in S \cap P} \sigma(f_w(x_i))(1 - \sigma(f_w(x_i))) (\phi(x_i) \otimes \phi(x_i))$. We can then note h is strongly convex with $\mu = \Omega(\lambda_{\min}(\sum_{i \in S \cap P} \phi(x_i) \otimes \phi(x_i)))$. Then from the properties of strongly convex functions, we have

$$\sum_{i \in S \cap P} (f_{\hat{w}}(x_i) - f_w^*(x_i)) (\sigma(f_{\hat{w}}(x_i)) - \sigma(f_w^*(x_i))) \gtrsim \lambda_{\min} \left(\sum_{i \in S \cap P} \phi(x_i) \otimes \phi(x_i) \right) \|f_w^* - f_{\hat{w}}\|_{\mathcal{H}}^2$$

Then from the Cauchy-Schwarz Inequality, we have

$$\begin{aligned} \sum_{i \in S} (f_w^*(x_i) - f_{\hat{w}}(x_i)) (y_i - \sigma(f_w^*(x_i))) &\leq \max_{i \in S} |y_i - \sigma(f_w^*(x_i))| \left\langle \sum_{j \in X} (w_j^* - \hat{w}_j) \phi(x_j), \sum_{i \in S} \phi(x_i) \right\rangle \\ &\leq \|f_w^* - f_{\hat{w}}\|_{\mathcal{H}} \left\| \sum_{i \in S} \phi(x_i) \right\|_{\mathcal{H}} \leq \|f_w^* - f_{\hat{w}}\|_{\mathcal{H}} \left(\left\| \sum_{i \in S \cap P} \phi(x_i) \right\|_{\mathcal{H}} + \left\| \sum_{i \in S \cap Q} \phi(x_i) \right\|_{\mathcal{H}} \right) \end{aligned}$$

for a small positive constant we denote c_3 . This completes the proof. \blacksquare

B.12 Proof of Theorem 20

Proof. From Lemma 19, we have in expectation

$$\begin{aligned} \mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} (\nabla_f g(f_{\hat{w}}, t))(f_w^* - f_{\hat{w}}) &\stackrel{\text{lem. 19}}{\geq} c_3 \left(n(1 - 2\varepsilon) \lambda_{\min}(\Sigma) - \mathbb{E}_{x_i \sim \mathbb{P}} \left\| \sum_{i \in S \cap P} \phi(x_i) \otimes \phi(x_i) - \Sigma \right\| \right) \|f_{\hat{w}} - f_w^*\|_{\mathcal{H}}^2 \\ &\quad - \left(\sqrt{n(1 - 2\varepsilon) \text{Tr}(\Sigma)} + \sqrt{n\varepsilon Q_k} \right) \|f_{\hat{w}} - f_w^*\|_{\mathcal{H}} \end{aligned}$$

We will lower bound the constant we introduced in Appendix B.11 and call it c_3 , recall for $f \in \mathcal{K}$, we have $\|f\|_{\mathcal{H}} \leq R$ and $P_k \triangleq \max_{i \in P} k(x_i, x_i)$.

$$\begin{aligned} \mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} c_3 &\stackrel{\text{appendix B.11}}{=} \mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} \min_{i \in S \cap P} (1 - \sigma(f_{\hat{w}}(x_i))) \sigma(f_{\hat{w}}(x_i)) \\ &\geq \mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} \left(1 - \sigma(R \max_{i \in P} k(x_i, x_i)) \right) \sigma(R \max_{i \in P} k(x_i, x_i)) \\ &\geq \mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} \frac{\sigma(-R \max_{i \in P} k(x_i, x_i))}{2} \stackrel{(i)}{\gtrsim} \exp \left(-R \mathbb{E} \left[\max_{i \in P} k(x_i, x_i) \right] \right) \end{aligned}$$

⁶Work in Progress

$$\stackrel{\text{lem. 28}}{\geq}$$

(i) follows from Jensen's Inequality as $\exp(-x)$ is a convex function. Then we have from the definition of a stationary point, $\nabla_{fg}(f_{\hat{w}}, t)(f_{\hat{w}} - f_w^*) \leq 0$ when

$$\mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} \|f_{\hat{w}} - f_w^*\|_{\mathcal{H}} \leq O \left(\frac{\sqrt{n(1-2\varepsilon) \text{Tr}(\Sigma)} + \sqrt{n\varepsilon Q_k}}{n(1-2\varepsilon)\lambda_{\min}(\Sigma) - 2\sqrt{n(1-2\varepsilon) \text{Tr}(\Sigma)}} \right)$$

If $n \geq \frac{4 \text{Tr}(\Sigma)}{\lambda_{\min}(\Sigma)(1-2\varepsilon)(1-c_4)}$ for $c_4 \in (0, 1)$, then we have

$$\mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} \|f_{\hat{w}} - f_w^*\|_{\mathcal{H}} \leq O \left(\frac{\sqrt{n(1-2\varepsilon) \text{Tr}(\Sigma)} + \sqrt{n\varepsilon Q_k}}{c_3 c_4 n(1-2\varepsilon)\lambda_{\min}(\Sigma)} \right) = O \left(\frac{\sqrt{\text{Tr}(\Sigma)} + \sqrt{Q_k}}{\sqrt{n(1-2\varepsilon)\lambda_{\min}(\Sigma)}} \right)$$

This completes the proof as we see we have $O(1/\sqrt{n})$ convergence. \blacksquare

B.13 Proof of Lemma 21

Proof. We use the Hilbert Space norm of the gradient to bound L from above. Let S be denoted as the subquantile set. We first give some derivatives.

$$\frac{\partial}{\partial w_k} (\ell(x_i, y_i; f_W)) = \begin{cases} -\phi(x_i) \text{softmax}(f_W(x_i))_k & \text{if } k = y_i \\ \phi(x_i) (1 - \text{softmax}(f_W(x_i))_k) & \text{if } k \neq y_i \end{cases}$$

Our proof then follows very similarly to the proof for Lemma 17.

$$\begin{aligned} \|\nabla_{fg}(t, f_W)\|_{\mathcal{H}} &= \left\| \frac{1}{np} \sum_{i=1}^n \mathbb{I} \left\{ -\log \left(\text{softmax}(f_W(x_i))_{y_i} \right) \geq t \right\} (\text{softmax}(f_W(x_i)) - y_i) \odot k(x_i, \cdot) \right\| \\ &\leq \frac{1}{np} \sum_{i=1}^n \|(\text{softmax}(f_W(x_i)) - y_i) \odot k(x_i, \cdot)\| \leq \frac{1}{np} \sum_{i=1}^n \|k(x_i, \cdot)\|_{\mathcal{H}} = \frac{1}{np} \sum_{i=1}^n \sqrt{k(x_i, x_i)} \end{aligned}$$

This completes the proof. \blacksquare

B.14 Proof of Lemma 22

Proof. We upper bound the operator norm of the Hessian to make our claim. We will give some derivatives.

$$\frac{\partial}{\partial w_k \partial w_j} (\text{softmax}(f_W(x_i))_k) = \{$$

$$\|\nabla_{fg}(t, f_W)\|_{\text{op}} = \left\| \frac{1}{np} \sum_{i=1}^n \mathbb{I} \left\{ -\log \left(\text{softmax}(f_W(x_i))_{y_i} \right) \geq t \right\} \right\|_{\text{op}}$$

\blacksquare

B.15 Proof of Lemma 23

Proof. We follow the similar set up as in Appendix B.11. We have $f_{\hat{W}}$ is a stationary point, i.e. $0 \in \partial\Phi(f_{\hat{W}})$. This implies for all $f_W \in \mathcal{K}$, we have $\Phi(f_{\hat{W}}) \leq \Phi(f_W)$. As Φ is differentiable, we have the first-order stationary condition, which is $\nabla\Phi(f_{\hat{W}})(f_{\hat{W}} - f_W) \leq 0$ or for all $W \in \mathcal{K}$. We assume $f_W^* \in \mathcal{K}$. Let S be the Subquantile set for $f_{\hat{W}}$. We will proceed by contradiction, assume $\|f_{\hat{W}} - f_W^*\|_{\mathcal{H}} \geq \eta$. Then, we have

$$\begin{aligned} (\nabla_{fg}(f_{\hat{W}}, t))(f_{\hat{W}} - f_W^*) &= (f_{\hat{W}} - f_W^*) \left(\sum_{i \in S} (\text{softmax}(f_{\hat{W}}(x_i)) - y_i) \odot k(x_i, \cdot) \right) \\ &= (f_{\hat{W}} - f_W^*) \left(\sum_{i \in S} (\text{softmax}(f_{\hat{W}}(x_i)) - \text{softmax}(f_W^*(x_i))) \odot k(x_i, \cdot) \right) + (f_{\hat{W}} - f_W^*) \left(\sum_{i \in S} (\text{softmax}(f_W^*(x_i)) - y_i) \odot k(x_i, \cdot) \right) \end{aligned}$$

Let us now consider the function $h : \mathcal{H} \times \dots \times \mathcal{H} \rightarrow \mathbb{R}^n$ defined as $h(f_W) = \sum_{i \in S} \log(\sum_{j=1}^{|\mathcal{Y}|} \exp(f_{w_j}(x_i) - y_{i,j}))$. We then have $h'(f_W) = \sum_{i \in S} \text{softmax}(f_W(x_i) - y_i) \odot \phi(x_i)$. From which it follows $h''(x_i) = \sum_{i \in S}$

C Necessary Results

Theorem 32 (Mercer, [22]). *If k is a positive definite and continuous kernel on a compact set Ω , the operator T has a countable set of eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq 0$ and eigenfunctions $\{\varphi_j\}_{j \in \mathbb{N}}$ with $T\varphi_j = \lambda_j \varphi_j$ s.t. $\|\varphi_j\| = \lambda_j^{-1}$. The kernel then admits the following representation,*

$$k(x, y) = \sum_{j=1}^{\infty} \lambda_j \varphi_j(x) \varphi_j(y)$$