

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 [LPMH21, RRM14]. 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- ϵ Contamination Problem for Kernel Learning where the distribution is formed as $\hat{\mathcal{P}} = (1 - \epsilon)\mathcal{P} + \epsilon\mathcal{Q}$, and we want to find the function $\inf_{f_{\mathbf{w}} \in \mathcal{H}} \mathbb{E}_{\mathcal{D} \sim \mathcal{P}} [\ell(f_{\mathbf{w}}; \mathbf{X}, \mathbf{y})]$, from the noisy distribution, $\hat{\mathcal{P}}$. We assume the adversary has knowledge of the true distribution of \mathcal{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. We empirically test Kernel Regression and Kernel Classification on various state of the art datasets and show Subquantile Minimization gives strong results in comparison to the state of the art robust algorithms.

*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), [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 \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 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 in [HYwL20]. 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.

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. In this work, we adopt a similar framework, and with the power of Mercer's Theorem [Mer09], we are able to say $\text{Tr}(\Sigma) < \infty$. We use this fact extensively in our infinite-dimensional concentration inequalities.

Theorem 2. (Informal). Let the dataset be given as $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ such that the labels and features of ϵn samples are arbitrarily corrupted by an adversary. Assume Subquantile Minimization returns $\hat{f}_{\mathbf{w}}$ 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:

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

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

Kernel Binary Classification:

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

Assume $n \geq \left(\frac{\text{Tr}(\Sigma)}{\lambda_{\min}(\Sigma)(1-c)} \right)^2$ for a constant $0 < c < 1$.

Kernel Multi-Class Classification:

$$\mathbb{E}_{\mathcal{D} \sim \hat{\mathcal{P}}} \| \hat{f}_{\mathbf{w}} - f_{\mathbf{w}}^* \| \leq O \left(\frac{\sqrt{n(1-\epsilon)\text{Tr}(\Sigma)} + \sqrt{n\epsilon Q_k}}{n(1-2\epsilon)\lambda_{\min}(\Sigma)} \right) \quad (3)$$

1.1 Related Work

The idea of iterative thresholding algorithms for robust learning tasks dates back to 1806 by Legendre [Leg06]. 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.

Robust Regression via Hard Thresholding [BJK15]. 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 [BJK15] (up to constants) for linear regression with and without feature corruption, which is one of our key contributions.

Learning with bad training data via iterative trimmed loss minimization [SS19a]. 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.

Trimmed Maximum Likelihood Estimation for Robust Generalized Linear Model [ADKS22]. 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.

Sum of Ranked Range Loss for Supervised Learning [HYwL20]. Hu et al. proposed learning over the bottom k losses, this is an alternative formulation of our algorithm. This is an extension of previous work studying the learning of the top k losses, [FLYH17]. 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. Subsequent work on analyzing the middle k losses is analyzed in [HYW+23].

The iterative trimmed loss framework with batch Stochastic Gradient Descent (SGD) is analyzed in [SS19b]. They experimentally test their design in deep learning applications such as image classification and Generative Adversarial Networks (GANs).

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 corruption with a general Gaussian Design.
3. We perform experiments on state-of-the-art matrices and show the effectiveness of our algorithm compared to other robust learning procedures. Furthermore, we use our experiments to demonstrate the practicality of our theory.

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 .

$$\mathcal{Q}_p(U) \leq t \iff \mathbb{P}\{U \leq t\} \geq p \quad (4)$$

The p -subquantile of the risk is then given by

$$\mathbb{L}_p(U) = \frac{1}{p} \int_0^p \mathcal{Q}_q(U) dq = \mathbb{E}[U | U \leq \mathcal{Q}_p(U)] = \max_{t \in \mathbb{R}} \left\{ t - \frac{1}{p} \mathbb{E}(t - U)^+ \right\} \quad (5)$$

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

$$\min_{f_{\mathbf{w}} \in \mathcal{K}} \max_{t \in \mathbb{R}} \left\{ g(t, f_{\mathbf{w}}) \triangleq t - \sum_{i=1}^n (t - (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2)^+ \right\} \quad (6)$$

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

$$t^{(k+1)} = \arg \max_{t \in \mathbb{R}} g(f_{\mathbf{w}}^{(k)}, t) \quad (7)$$

$$f_{\mathbf{w}}^{(k+1)} = \text{Proj}_{\mathcal{K}} \left(f_{\mathbf{w}}^{(k)} - \alpha \nabla_f g(f_{\mathbf{w}}^{(k)}, t^{(k+1)}) \right) \quad (8)$$

We provide an algorithm for Subquantile Minimization of the ridge regression and classification kernel learning algorithm.

Algorithm 1 (Subquantile Minimization for Kernel Learning).

Input: Iterations: T , Quantile: p ; Data Matrix: $\mathbf{X} \in \mathbb{R}^{n \times d}$, $n \gg d$; Labels: $\mathbf{y} \in \mathbb{R}^{n \times 1}$; Learning Rate schedule: $\eta^{(1)}, \dots, \eta^{(T)}$

Output: Function in RKHS: $\hat{f}_{\mathbf{w}}$

1. Initialize the weights $w_i^{(0)} \sim \text{Unif} \left[-\sqrt{\frac{6}{n}}, \sqrt{\frac{6}{n}} \right]$ for all $i \in [n]$.
2. **for** $k = 1, 2, \dots, T$ **do**
 3. Find the Subquantile denoted as $S^{(k)}$ as the set of $(1 - \epsilon)n$ elements with the lowest error with respect to the loss function.
 4. Update the gradient in accordance with the kernel learning problem.

$$\nabla_f g \left(t^{(k+1)}, f_{\mathbf{w}}^{(k)} \right) \leftarrow \frac{2}{n(1 - \epsilon)} \sum_{i \in S^{(k)}} \left(f_{\mathbf{w}}^{(k)}(\mathbf{x}_i) - y_i \right) \cdot k(\mathbf{x}_i, \cdot) \quad (\text{Regression})$$

$$\nabla_f g \left(t^{(k+1)}, f_{\mathbf{w}}^{(k)} \right) \leftarrow \frac{1}{n(1 - \epsilon)} \sum_{i \in S^{(k)}} \left(\sigma \left(f_{\mathbf{w}}^{(k)}(\mathbf{x}_i) \right) - y_i \right) \cdot k(\mathbf{x}_i, \cdot) \quad (\text{Binary Classification})$$

$$\nabla_f g \left(t^{(k+1)}, f_{\mathbf{w}}^{(k)} \right) \leftarrow \frac{1}{n(1 - \epsilon)} \sum_{i \in S^{(k)}} \left(\text{softmax} \left(f_{\mathbf{w}}^{(k)}(\mathbf{x}_i) \right) - \mathbf{y}_i \right) \odot k(\mathbf{x}_i, \cdot) \quad (\text{Multi-Class Classification})$$

5. Perform Projected Standard Gradient Descent to find the next iterate

$$f_{\mathbf{w}}^{(k+1)} \leftarrow \text{Proj}_{\mathcal{K}} \left[f_{\mathbf{w}}^{(k)} - \eta^{(k)} \nabla_f g \left(t^{(k+1)}, f_{\mathbf{w}}^{(k)} \right) \right]$$

Return: Function in RKHS: $f_{\mathbf{w}}^{(T)}$

3 Theory

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 (7) and (8).

Lemma 3. *Let $f(\mathbf{x}; \mathbf{w})$ be a convex loss function. Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ denote the n data points ordered such that $f(\mathbf{x}_1; \mathbf{w}, y_1) \leq f(\mathbf{x}_2; \mathbf{w}, y_2) \leq \dots \leq f(\mathbf{x}_n; \mathbf{w}, y_n)$. If we denote $\hat{v}_i \triangleq f(\mathbf{x}_i; \mathbf{w}, y_i)$, it then follows $\hat{v}_{np} \in \arg \max_{t \in \mathbb{R}} g(t, \mathbf{w})$.*

Proof is given in Appendix B.1. From Lemma 3, we see that 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 4. Let $\hat{\nu}_i \triangleq f(\mathbf{x}_i; \mathbf{w}, y_i)$ 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 $\nabla_{\mathbf{w}} g(t^{(k)}, f_{\mathbf{w}}^{(k)}) = \frac{1}{np} \sum_{i=1}^{np} \nabla f(\mathbf{x}_i; f_{\mathbf{w}}^{(k)}, y_i)$

Proof is given in Appendix B.2.

3.1 Kernelized Regression

The loss for the Kernel Ridge Regression problem for a single training pair $(\mathbf{x}_i, y_i) \in \mathcal{D}$ is given by the following equation

$$\ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i) = (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 \quad (9)$$

Our goals throughout the proofs will be to obtain approximation bounds for infinite-dimensional kernels. The key challenge is the obvious undetermined problem, i.e. considering an infinite eigenfunction basis, we require infinite samples to obtain an accurate approximation. Instead, we will calculate the approximation bounds for the rank- m approximation of $f_{\mathbf{w}}^*$ and push $m \rightarrow \infty$.

Definition 5 (Strong Projection Property). Let $f_{\mathbf{w}}^*$ be the optimal function for the uncorrupted dataset, \mathbb{P} . Then, we have for a finite m and an absolute constant $c > 0$,

$$\|\text{Proj}_{\Psi_m} f_{\mathbf{w}}^* - f_{\mathbf{w}}^*\|_{\mathcal{H}} = 0, \quad \|\text{Proj}_{\Psi_m} f_{\mathbf{w}} - f_{\mathbf{w}}^*\|_{\mathcal{H}} > c \quad (10)$$

The *Strong Projection Property* is important as $\lambda_{\min}(\Sigma)$ is not well defined for an infinite dimensional feature space, e.g. Gaussian Kernel. The implication of the *Strong Projection Property* is given in the following lemma.

Lemma 6 (Strong Projection Property Implication). Assume the *Strong Projection Property* (Definition 5) holds for $f_{\mathbf{w}}^{(t)}$ for all $t \in [T]$, where $f_{\mathbf{w}}^{(t)}$ are iterates from Algorithm 1. Then, it follows for a $m \in \mathbb{N}$ and a constant $0 < C \leq 1$,

$$\left\langle \Sigma, \left(f_{\mathbf{w}}^{(t)} - f_{\mathbf{w}}^* \right) \otimes \left(f_{\mathbf{w}}^{(t)} - f_{\mathbf{w}}^* \right) \right\rangle_{\text{HS}} \geq C \lambda_m \|f_{\mathbf{w}}^{(t)} - f_{\mathbf{w}}^*\|_{\mathcal{H}}^2, \quad (11)$$

where we define

$$C \triangleq \frac{\|f_{\mathbf{w}}^{(t)} - f_{\mathbf{w}}^*\|_{\mathcal{H}}^2 - \|\text{Proj}_{\Psi_m} f_{\mathbf{w}}\|_{\mathcal{H}}^2}{\|f_{\mathbf{w}}^{(t)} - f_{\mathbf{w}}^*\|_{\mathcal{H}}^2} \quad (12)$$

Proof is given in Appendix C.2.

Theorem 7 (Subquantile Minimization for Kernelized Regression is Good with High Probability). Let Algorithm 1 be run on a dataset $\mathcal{D} \sim \hat{\mathcal{P}}$ with learning rate $\eta \triangleq \beta^{-1}$ where β is the Lipschitz Gradient Constant given in Lemma 25. Suppose

$$n = O\left(\frac{\text{Tr}(\Sigma)^2}{(1 - 2\epsilon)\lambda_m^2(\Sigma)}\right) \quad (13)$$

Then $\|\text{Proj}_{\Psi_m} f_{\mathbf{w}} - \text{Proj}_{\Psi_m} f_{\mathbf{w}}^*\|_{\mathcal{H}} \leq \epsilon$, after

$$T = O\left(\frac{\log\left(\|f_{\mathbf{w}}^*\|_{\mathcal{H}} + \sigma \sqrt{n \log n \text{Tr}(\Sigma)} + \|\xi\| \sqrt{\text{Tr}(\mathbf{K}_Q)}\right)}{\epsilon \log\left(\frac{1}{1 - \lambda_m}\right)}\right) \quad (14)$$

iterations with high probability.

Full proof with explicit constants is given in Appendix C.3. A direct application of Theorem 7 is that learning an infinite dimensional function $f_{\mathbf{w}}^*$ to within ε error in the Hilbert Space Norm requires infinite data. Furthermore, we see that given covariate noise and label noise, our bound requires more iterations dependent on the magnitude of the corruption. Such a result is corroborated in [SST+18]. 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. The effect of Lemma 6 can be seen in the denominator of both terms. Instead of $\lambda_{\min}(\mathbf{\Sigma})$ we have $c_4\lambda_m$ for a finite m . This difference will be clear in the following corollary, where we utilize the theory developed for kernelized regression to imply a result for regularized linear regression.

Corollary 8 (Linear Regression Expected Error Bound). *Consider Subquantile Minimization for Linear Regression on the data X with optimal parameters \mathbf{w}^* . Assume $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$ 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*

$$\mathbb{E}\|\mathbf{w}^{(T)} - \mathbf{w}^*\|_2 \leq O\left(\frac{\gamma\sigma}{\sqrt{\lambda_{\min}(\mathbf{\Sigma})}}\right) \quad (15)$$

Proof given in ???. Let us note for the case where p is finite, i.e. the feature mapping is finite-dimensional, e.g. linear or polynomial kernel. Then we have that $\text{Proj}_{\Psi_m^\perp}$ where $m = p$ is equal to zero as $\{\varphi_i\}_{i=1}^m$ spans the finite-dimensional space, in which we case we have the absolute constant given in Definition 5 is equal to zero. It is important to note in all our bounds, $\gamma \leq \sqrt{\frac{\epsilon}{1-2\epsilon}}$ is a theoretical worst case bound when the Subquantile contains the minimum possible number of uncorrupted points. In other words, we have $\gamma \triangleq \frac{|P \setminus S|}{|S \cap P|} \leq \frac{n\epsilon}{n(1-2\epsilon)} = \frac{\epsilon}{1-2\epsilon}$. 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.

3.2 Kernelized Binary Classification

The Negative Log Likelihood for the the Kernel Classification problem is given by the following equation for a single training pair (\mathbf{x}_i, y_i)

$$\ell(\mathbf{x}_i, y_i; f_{\mathbf{w}}) = -y_i \log(\sigma(f_{\mathbf{w}}(\mathbf{x}_i))) - (1 - y_i) \log(1 - \sigma(f_{\mathbf{w}}(\mathbf{x}_i))) \quad (16)$$

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

$$\mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} \|f_{\hat{\mathbf{w}}} - f_{\mathbf{w}}^*\|_{\mathcal{H}} \leq O\left(\frac{\sqrt{\text{Tr}(\mathbf{\Sigma})} + \sqrt{Q_k}}{\sqrt{n(1-2\epsilon)} \exp(-R(\text{Tr}(\mathbf{\Sigma}) + \log n)) \lambda_{\min}(\mathbf{\Sigma})}\right) \quad (17)$$

Proof is given in ???. This result although shows consistency, i.e. when $n \rightarrow \infty$, then we have in expectation $\|f_{\mathbf{w}} - f_{\mathbf{w}}^*\| \rightarrow 0$, 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. To mitigate the effect of a large Q_k , a filtering algorithm can be used to remove points, \mathbf{x}_i , such that $k(\mathbf{x}_i, \mathbf{x}_i)$ is far from the mean.

3.3 Kernelized 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 (\mathbf{x}_i, y_i) , note $\mathbf{W} \in \mathbb{R}^{n \times |\mathcal{Y}|}$.

$$\ell(\mathbf{x}_i, y_i; f_{\mathbf{w}}) = -\sum_{j=1}^{|\mathcal{Y}|} \mathbb{I}\{j = y_i\} \log\left(\frac{\exp(f_{\mathbf{w}_j}(\mathbf{x}_i))}{\sum_{k=1}^{|\mathcal{Y}|} \exp(f_{\mathbf{w}_k}(\mathbf{x}_i))}\right) \quad (18)$$

Theorem 10 (Stationary Point for Kernelized Multi-Class Classification is Good). *Let $f_{\hat{\mathbf{w}}}$ be a stationary point defined in ?? for the function Φ defined in ?. Then for a constant $c_1 \in (0, 1)$, if*

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

$$\mathbb{E}_{\mathcal{D} \sim \hat{\mathbb{P}}} \|f_{\hat{\mathbf{w}}} - f_{\mathbf{w}}^*\|_{\mathcal{H}} \leq O \left(\frac{\sqrt{n(1-\epsilon) \text{Tr}(\mathbf{\Sigma})} + \sqrt{n\epsilon Q_k}}{n(1-2\epsilon)\lambda_{\min}(\mathbf{\Sigma}) - \sqrt{n(1-2\epsilon) \text{Tr}(\mathbf{\Sigma})}} \right) \quad (19)$$

where β is the Lipschitz Gradient Constant given in ??.

3.4 Optimization

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 11 ([Roc70, DD19]). *Assume the function Φ is β -weakly convex. Let $\lambda < \frac{1}{\beta}$, and let $f_{\hat{\mathbf{w}}} = \arg \min_{f_{\mathbf{w}} \in \mathcal{K}} (\Phi(f_{\mathbf{w}}) + \frac{1}{2\lambda} \|f_{\mathbf{w}} - f_{\hat{\mathbf{w}}}\|_{\mathcal{H}}^2)$, then $\|\nabla \Phi_{\lambda}(f_{\mathbf{w}})\|_{\mathcal{H}} \leq \epsilon$ implies:*

$$\|f_{\hat{\mathbf{w}}} - f_{\mathbf{w}}\|_{\mathcal{H}} = \lambda \epsilon \quad \text{and} \quad \min_{\mathbf{g} \in \partial \Phi(f_{\hat{\mathbf{w}}}) + \partial \mathcal{I}_{\mathcal{K}}(f_{\hat{\mathbf{w}}})} \|\mathbf{g}\|_{\mathcal{H}} \leq \epsilon \quad (20)$$

With Lemma 11 in hand, it suffices to show that $\|\nabla \Phi_{\lambda}(f_{\mathbf{w}})\|_{\mathcal{H}}$ is small, as it then follows that $f_{\mathbf{w}}$ is close to a stationary point of the Moreau Envelope. It has been shown in optimization theory that utilizing standard gradient descent, $\|\nabla \Phi_{\lambda}(f_{\mathbf{w}})\|_{\mathcal{H}}$ decreases at a rate of $O(T^{-1/2})$. The exact theorem and proof can be seen in [DD19] and a proof where the maximum of the inner problem can be calculated to within $(1+\epsilon)$ optimality can be seen in [JNJ20] and [CDGS20].

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 $f_{\mathbf{w}}$ with the Glorot Initialization Scheme [GB10].

Algorithms	Test RMSE							
	Concrete [Yeh07]		Wine Quality [CR09]		Boston Housing [DG17]		Drug [OSB+18]	
	$\epsilon = 0.2(\downarrow)$	$\epsilon = 0.4(\downarrow)$	$\epsilon = 0.2(\downarrow)$	$\epsilon = 0.4(\downarrow)$	$\epsilon = 0.2(\downarrow)$	$\epsilon = 0.4(\downarrow)$	$\epsilon = 0.2(\downarrow)$	$\epsilon = 0.4(\downarrow)$
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	0.533 _(0.0347)	0.592 _(0.0548)	0.915 _(0.0343)	0.841 _(0.0413)	0.526 _(0.0287)	0.720 _(0.1147)	1.172 _(0.0542)	1.215 _(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: $y_{\text{noise}} = 10000y_{\text{original}}$ and $\mathbf{x}_{\text{noise}} = 100\mathbf{x}_{\text{original}}$. 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.

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.

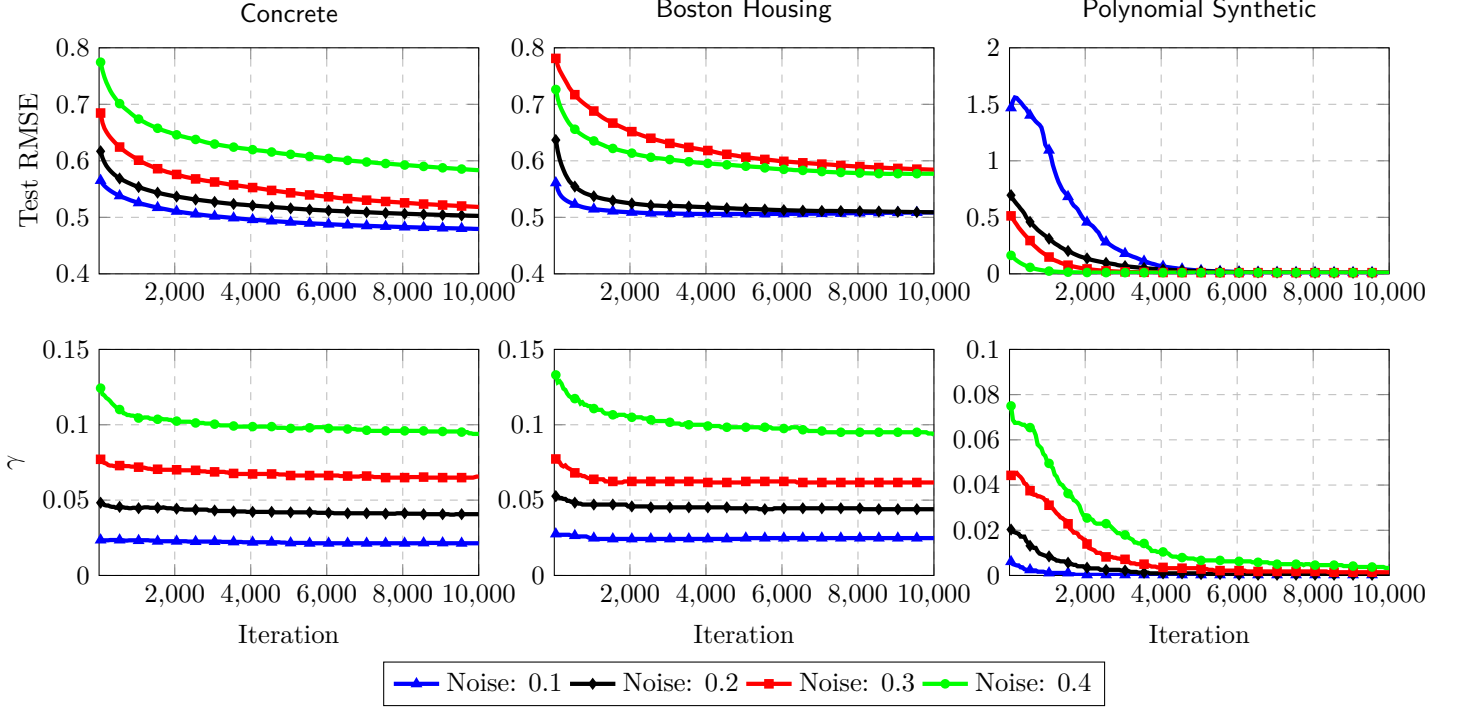


Figure 1: Test RMSE over the iterations in Concrete, Boston Housing, and Polynomial Datasets for SUB-QUANTILE at different noise levels

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 state of the art algorithms developed specifically for the Robust Linear Problem. In particular, we compare against Kernel Ridge Regression (KRR) implemented in the sklearn package [PVG+11], Consistent Robust Regression (CRR) [BJKK17], Globally-convergent iteratively reweighted least squares (STIR) [MGJK19]. We also compare with several robust meta-algorithms, i.e. algorithms which work for multiple robust learning tasks, e.g classification and regression. We compare with SEVER [DKK+19] and Tilted Empirical Risk Minimization [LBSS21].

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.

Algorithms	Test RMSE							
	Boston Housing [DG17]		Wine Quality [CR09]		Concrete [Yeh07]		Drug [OSB ⁺ 18]	
	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 $\mathbf{x}_{\text{noisy}} = 100\mathbf{x}_{\text{original}}$ 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 10^5 iterations, and the feature noise perturbed dataset for 10^6 iterations.

Algorithms	Test Accuracy							
	Heart Disease [JD88]				Breast Cancer [WS95]			
	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	<u>0.793</u> _(0.0422)	<u>0.695</u> _(0.0636)	0.784 _(0.0432)	0.816 _(0.0562)	0.904 _(0.0356)	0.575 _(0.1456)	0.956 _(0.0164)	<u>0.974</u> _(0.0062)
TERM	0.741 _(0.0393)	0.620 _(0.0699)	<u>0.803</u> _(0.0613)	<u>0.810</u> _(0.0286)	<u>0.940</u> _(0.0378)	<u>0.763</u> _(0.0364)	0.986 _(0.0143)	0.986 _(0.0119)
SUBQUANTILE	0.803 _(0.0293)	0.790 _(0.0350)	0.833 _(0.0318)	<u>0.807</u> _(0.0468)	0.928 _(0.0129)	0.916 _(0.0185)	<u>0.972</u> _(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: $\mathbf{x}_{\text{noise}} = 100\mathbf{x}_{\text{original}}$. The Linear Kernel is used in all experiments.

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

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 $\mathbf{x}_{\text{noise}} = 100\mathbf{x}_{\text{original}}$. The Radial Basis Function is used in all experiments.

Algorithms	Test Accuracy							
	Iris [Fis88]		Glass [Ger87]		Wine [AF91]		Satimage [Sri93]	
	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$
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	∞	∞	∞	∞	∞	∞	∞	∞

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.

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 be 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, 2
- [AF91] Stefan Aeberhard and M. Forina. Wine. UCI Machine Learning Repository, 1991. DOI: <https://doi.org/10.24432/C5PC7J>. 9
- [AMMIL12] Yaser S Abu-Mostafa, Malik Magdon-Ismail, and Hsuan-Tien Lin. *Learning from data*, volume 4. AMLBook New York, 2012. 1
- [Ber24] Sergei Bernstein. On a modification of chebyshev’s inequality and of the error formula of laplace. *Ann. Sci. Inst. Sav. Ukraine, Sect. Math*, 1(4):38–49, 1924. 18
- [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. 2
- [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, 7

- [Bog98] Vladimir Igorevich Bogachev. *Gaussian measures*. Number 62. American Mathematical Soc., 1998. 13
- [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, 6
- [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
- [CR09] Cerdeira A. Almeida F. Matos T. Cortez, Paulo and J. Reis. Wine Quality. UCI Machine Learning Repository, 2009. DOI: <https://doi.org/10.24432/C56S3T>. 6, 8
- [DD19] Damek Davis and Dmitriy Drusvyatskiy. Stochastic model-based minimization of weakly convex functions. *SIAM Journal on Optimization*, 29(1):207–239, 2019. 6
- [DG17] Dheeru Dua and Casey" Graff. UCI machine learning repository, 2017. 6, 8
- [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, 7
- [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
- [Fis88] R. A. Fisher. Iris. UCI Machine Learning Repository, 1988. DOI: <https://doi.org/10.24432/C56C76>. 9
- [FLYH17] Yanbo Fan, Siwei Lyu, Yiming Ying, and Baogang Hu. Learning with average top-k loss. 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. 2
- [FWZ18] Jianqing Fan, Weichen Wang, and Yiqiao Zhong. An l_1 eigenvector perturbation bound and its application to robust covariance estimation. *Journal of Machine Learning Research*, 18(207):1–42, 2018. 1
- [GB10] 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. 6
- [Ger87] B. German. Glass Identification. UCI Machine Learning Repository, 1987. DOI: <https://doi.org/10.24432/C5WW2P>. 9

- [Gre13] 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. [16](#)
- [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](#)
- [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](#), [2](#)
- [JD88] Steinbrunn William Pfisterer Matthias Janosi, Andras and Robert Detrano. Heart Disease. UCI Machine Learning Repository, 1988. DOI: <https://doi.org/10.24432/C52P4X>. [8](#)
- [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. [6](#)
- [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](#), [7](#)
- [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. [1](#)
- [LPMH21] 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.
- [Mer09] 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. [1](#)
- [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](#), [7](#)
- [OSB⁺18] Ivan Olier, Nouredin Sadawi, G Richard Bickerton, Joaquin Vanschoren, Crina Grosan, Larisa Soldatova, and Ross D King. Meta-qsar: a large-scale application of meta-learning to drug design and discovery. *Machine Learning*, 107:285–311, 2018. [6](#), [8](#)
- [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
- [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
- [PVG⁺11] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011. 7
- [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
- [Roc70] Ralph Tyrell Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, 1970. 6
- [RRM14] 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.
- [Sri93] Ashwin Srinivasan. Statlog (Landsat Satellite). UCI Machine Learning Repository, 1993. DOI: <https://doi.org/10.24432/C55887>. 9
- [SS19a] 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
- [SS19b] Yanyao Shen and Sujay Sanghavi. Learning with bad training data via iterative trimmed loss minimization. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5739–5748. PMLR, 09–15 Jun 2019. 2
- [SST⁺18] 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. 5
- [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. 21
- [Ver20] Roman Vershynin. High-dimensional probability. *University of California, Irvine*, 2020. 18
- [Wey12] 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. 22
- [WS95] Mangasarian Olvi Street Nick Wolberg, William and W. Street. Breast Cancer Wisconsin (Diagnostic). UCI Machine Learning Repository, 1995. DOI: <https://doi.org/10.24432/C5DW2B>. 8
- [Yeh07] I-Cheng Yeh. Concrete Compressive Strength. UCI Machine Learning Repository, 2007. DOI: <https://doi.org/10.24432/C5PK67>. 6, 8

A Probability Theory

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

$$\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma) \quad (21)$$

where Σ is a possibly infinite dimensional covariance operator.

Assumption 13 (Normal Residuals). The residual is defined as $\mu_i \triangleq f_{\mathbf{w}}^*(\mathbf{x}_i) - y_i$. Then we assume for some $\sigma > 0$, it follows

$$\mu_i \sim \mathcal{N}(0, \sigma^2) \quad (22)$$

Proposition 14 (Concentration for functions of a Gaussian Vector [Bog98]). Suppose h is a Lipschitz function on vectors, i.e.

$$|h(\mathbf{x}) - h(\mathbf{y})| \leq \|h\|_{\text{lip}} \|\mathbf{x} - \mathbf{y}\| \quad (23)$$

for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. Then,

$$\mathbb{P}\{|h(\mathbf{g}) - \mathbb{E}h(\mathbf{g})| \geq T\} \leq 2 \exp \left[-\frac{t^2}{2 \|h\|_{\text{lip}}^2} \right] \quad (24)$$

Lemma 15 (Maximum of Gaussians). Let $\mu_1, \dots, \mu_n \sim \mathcal{N}(0, \sigma^2)$ for some $\sigma > 0$. Then it follows

$$\mathbb{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \max_{i \in [n]} |\mu_i| \leq \sigma \sqrt{2 \log n} + \frac{\sigma^2}{\sqrt{\pi \log n}} \quad (25)$$

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

$$\mathbb{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \max_{i \in [n]} |\mu_i| = \int_0^\infty \mathbb{P}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \left\{ \max_{i \in [n]} |\mu_i| > t \right\} dt \stackrel{(i)}{\leq} c_1 + n \int_{c_1}^\infty \mathbb{P}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \{|\mu_i| \geq t\} dt \quad (26)$$

$$\stackrel{(ii)}{=} c_1 + 2n \int_{c_1}^\infty \mathbb{P}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \{\mu_i \geq t\} dt = c_1 + 2n \int_{c_1}^\infty \int_t^\infty \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x}{\sigma}\right)^2} dx dt \quad (27)$$

$$\leq c_1 + \frac{n}{\sigma} \sqrt{\frac{2}{\pi}} \int_{c_1}^\infty \int_t^\infty \left(\frac{x}{t}\right) e^{-\frac{1}{2} \left(\frac{x}{\sigma}\right)^2} dx dt = c_1 + n\sigma \sqrt{\frac{2}{\pi}} \int_{c_1}^\infty \frac{e^{-\frac{1}{2} \left(\frac{t}{\sigma}\right)^2}}{t} dt \quad (28)$$

$$\leq c_1 + n\sigma \sqrt{\frac{2}{\pi}} \int_{c_1}^\infty \left(\frac{t}{c_1}\right) e^{-\frac{1}{2} \left(\frac{t}{\sigma}\right)^2} dt = c_1 + n\sigma^3 \sqrt{\frac{2}{\pi}} \frac{e^{-\frac{1}{2} \left(\frac{c_1}{\sigma}\right)^2}}{c_1} \quad (29)$$

(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 about zero. From here, we choose $c_1 \triangleq \sigma \sqrt{2 \log n}$. Then we have,

$$\mathbb{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \max_{i \in [n]} |\mu_i| \leq \sigma \sqrt{2 \log n} + \frac{\sigma^2}{\sqrt{\pi \log n}} \quad (30)$$

This completes the proof. ■

Proposition 16. Let $\mu_1, \dots, \mu_n \sim \mathcal{N}(0, \sigma^2)$ for some $\sigma > 0$, then it follows for any $s \geq 1$

$$\mathbb{P} \left\{ \max_{i \in [n]} |\mu_i| \geq \sigma \sqrt{2 \log n} \cdot s \right\} \leq \frac{\sqrt{2}}{\log n} e^{-s^2} \quad (31)$$

Proof. The proof follows simply using similar steps as in the proof of Lemma 15. Let C be a positive constant to be determined.

$$\mathbb{P}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \left\{ \max_{i \in n} |\mu_i| \geq C \cdot s \right\} = 2n \mathbb{P}_{\mu \sim \mathcal{N}(0, \sigma^2)} \{ \mu \geq C \cdot s \} = n \sqrt{\frac{2}{\pi}} \int_{C \cdot s}^{\infty} \left(\frac{1}{\sigma} \right) e^{-\frac{1}{2} \left(\frac{x}{\sigma} \right)^2} dx \quad (32)$$

$$\leq 2\sigma n \left(\frac{1}{C \cdot s} \right) e^{-\frac{1}{2} \left(\frac{C \cdot s}{\sigma} \right)^2} \leq \frac{\sqrt{2} n^{1-s^2}}{s \log n} \leq \frac{\sqrt{2}}{\log n} e^{-s^2} \quad (33)$$

In the second to last inequality, we plug in $C \triangleq \sigma \sqrt{2 \log n}$. Our proof is now complete. \blacksquare

Proposition 17. Let $\mu_1, \dots, \mu_n \sim \mathcal{N}(0, \sigma^2)$ for some $\sigma > 0$, then it follows for any $s \geq 1$,

$$\mathbb{P} \left\{ \sum_{i=1}^n \mu_i^2 \geq n\sigma^2 \cdot s \right\} \leq e^{-s/2} \quad (34)$$

Proof. Concatenate all the samples μ_i into a vector $\boldsymbol{\mu} \in \mathbb{R}^n$. Our proof generalizes for a $\boldsymbol{\mu} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ where $\boldsymbol{\Sigma} \triangleq \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top$ for a unitary \mathbf{U} and positive diagonal $\boldsymbol{\Lambda}$. Let C be a positive to be determined constant, we then have

$$\mathbb{P}_{\boldsymbol{\mu} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} \left\{ \|\boldsymbol{\mu}\|^2 \geq C \cdot s \right\} = \mathbb{P}_{\boldsymbol{\mu} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} \left\{ \|\boldsymbol{\mu}\| \geq \sqrt{C \cdot s} \right\} = \mathbb{P}_{\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left\{ \|\mathbf{U} \boldsymbol{\Lambda}^{1/2} \mathbf{g}\| \geq \sqrt{C \cdot s} \right\} \quad (35)$$

$$= \mathbb{P}_{\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left\{ \sqrt{\sum_{i=1}^n \lambda_i g_i^2} \geq \sqrt{C \cdot s} \right\} \leq \mathbb{P}_{\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left\{ \sum_{i=1}^n \sqrt{\lambda_i} g_i \geq \sqrt{C \cdot s} \right\} \quad (36)$$

$$\leq \inf_{\theta > 0} \mathbb{E}_{\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\prod_{i=1}^n \exp \left(\theta \sqrt{\lambda_i} g_i \right) \right] \exp \left[-\theta \sqrt{C \cdot s} \right] \quad (37)$$

$$= \inf_{\theta > 0} \exp \left[\frac{\text{Tr}(\boldsymbol{\Lambda})}{2} \theta^2 - \theta \sqrt{C \cdot s} \right] = \exp \left[-\frac{(C \cdot s)}{2 \text{Tr}(\boldsymbol{\Lambda})} \right] \quad (38)$$

The second to last equality follows from the MGF of a Gaussian. Then, plugging in $C \triangleq \text{Tr}(\boldsymbol{\Lambda})$ completes the proof. \blacksquare

Lemma 18 (Maximum of Squared Gaussians). Let $\mu_1, \dots, \mu_n \sim \mathcal{N}(0, \sigma^2)$ for $\sigma > 0, n > 1$. Then it follows

$$\mathbb{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \max_{i \in [n]} \mu_i^2 \leq 2\sigma^2 \log(n) + \left(\sigma^3 \sqrt{\frac{8}{\pi}} \right) \left(1 + \frac{1}{\log(n)} \right) \quad (39)$$

Proof. Our proof follows similarly to the proof for Lemma 15.

$$\mathbb{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \max_{i \in \mathbb{N}} \mu_i^2 = \int_0^\infty \mathbb{P}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \left\{ \max_{i \in [n]} \mu_i^2 \geq t \right\} dt \leq c_2 + n \int_{c_2}^\infty \mathbb{P}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \left\{ |\mu_i| \geq \sqrt{t} \right\} dt \quad (40)$$

$$= c_2 + 2n \int_{c_2}^\infty \mathbb{P}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \left\{ \mu_i \geq \sqrt{t} \right\} dt = c_2 + 2n \int_{c_2}^\infty \int_{\sqrt{t}}^\infty \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x}{\sigma} \right)^2} dx dt \quad (41)$$

$$= c_2 + n\sigma \sqrt{\frac{2}{\pi}} \int_{c_2}^\infty \frac{e^{-\frac{1}{2} \left(\frac{t}{\sigma^2} \right)}}{\sqrt{t}} dt \stackrel{(i)}{\leq} c_2 + n\sigma \sqrt{\frac{2}{\pi}} \int_{c_2}^\infty \left(\frac{t}{c_2} \right) e^{-\frac{1}{2} \left(\frac{t}{\sigma^2} \right)} dt \quad (42)$$

$$\leq c_2 + \left(\sqrt{\frac{2}{\pi}} \right) \frac{n\sigma (4\sigma^4 + 2c_2\sigma^2) e^{-\frac{c_2}{2\sigma^2}}}{c_2} \quad (43)$$

(i) holds for $c_2 > 1$. Then, setting $c_2 \triangleq 2\sigma^2 \log(n)$, we have

$$\mathbb{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \max_{i \in [n]} \mu_i^2 \leq 2\sigma^2 \log(n) + \left(2\sigma^3 \sqrt{\frac{2}{\pi}} \right) \left(1 + \frac{1}{\log(n)} \right) \quad (44)$$

This completes the proof. \blacksquare

Lemma 19 (Expected Maximum P_k). *Let $\mathbf{x}_i \sim \mathbb{P}$ such that $\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)$ from Assumption 12. Then it follows for any $s \geq 1$*

$$\mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left[\max_{i \in [n]} k(\mathbf{x}_i, \mathbf{x}_i) \right] \leq 2 \text{Tr}(\Sigma) \log \left(\frac{n \cdot s}{2 \text{Tr}(\Sigma)} \right) + \frac{1}{s} \quad (45)$$

Proof. We will use the integral identity of the expectation of a random variable to make our claim. Throughout the proof, let C be a positive to be determined constant.

$$\mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left[\max_{i \in [n]} k(\mathbf{x}_i, \mathbf{x}_i) \right] \leq C + \int_C^\infty \mathbb{P}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \max_{i \in [n]} k(\mathbf{x}_i, \mathbf{x}_i) \geq t \right\} dt \quad (46)$$

$$\stackrel{(i)}{\leq} C + n \int_C^\infty \mathbb{P}_{\phi(\mathbf{x}) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \{k(\mathbf{x}, \mathbf{x}) \geq t\} dt \quad (47)$$

$$\stackrel{(ii)}{=} C + n \int_C^\infty \mathbb{P}_{\phi(\mathbf{x}) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \|\phi(\mathbf{x})\|_{\mathcal{H}} \geq \sqrt{t} \right\} dt \quad (48)$$

$$\stackrel{(iii)}{=} C + n \int_C^\infty \mathbb{P}_{\psi(\mathbf{x}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left\{ \left\| \Psi \Lambda^{1/2} \psi(\mathbf{x}) \right\|_{\mathcal{H}} \geq \sqrt{t} \right\} dt \quad (49)$$

$$\leq C + n \int_C^\infty \mathbb{P}_{\psi_i(\mathbf{x}) \sim \mathcal{N}(0, 1)} \left\{ \sum_{i=1}^p \sqrt{\lambda_i} \psi_i(\mathbf{x}) \geq \sqrt{t} \right\} dt \quad (50)$$

$$\leq C + n \int_C^\infty \inf_{\theta > 0} \mathbb{E}_{\psi(\mathbf{x}) \sim \mathcal{N}(0, 1)} \left[\prod_{i=1}^p \exp \left(\theta \sqrt{\lambda_i} \psi_i(\mathbf{x}) \right) \right] \exp \left(-\theta \sqrt{t} \right) dt \quad (51)$$

$$= C + n \int_C^\infty \inf_{\theta > 0} \exp \left[\theta^2 \frac{\text{Tr}(\Sigma)}{2} - \theta \sqrt{t} \right] dt = C + n \int_C^\infty \exp \left[-\frac{t}{2 \text{Tr}(\Sigma)} \right] dt \quad (52)$$

$$= C + \frac{n}{2 \text{Tr}(\Sigma)} \exp \left[-\frac{C}{2 \text{Tr}(\Sigma)} \right] \quad (53)$$

See (i) from the proof of Lemma 15. (ii) follows from the reproducing property. In (iii) we define $\psi(\mathbf{x})$ as the whitened RKHS function. Setting $C \triangleq 2 \text{Tr}(\Sigma) \log(s \cdot n / (2 \text{Tr}(\Sigma)))$ completes the proof. ■

Lemma 20 (Norm of Functions with Gaussian Design in the Reproducing Kernel Hilbert Space). *Let $\mathbf{x}_i \sim \mathbb{P}$ such that $\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)$ from Assumption 12 and Assumption 13. Then, it follows*

$$\mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbb{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \left\| \sum_{i=1}^n \mu_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \leq O \left(\sigma \sqrt{n \log n \text{Tr}(\Sigma)} \right) \quad (54)$$

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,

$$\mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbb{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \left\| \sum_{i=1}^n \mu_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \quad (55)$$

$$\leq \mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbb{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \max_{i \in [n]} |\mu_i| \left\| \sum_{i=1}^n \phi(\mathbf{x}_i) \right\|_{\mathcal{H}}$$

$$\stackrel{(i)}{\leq} O \left(\sigma \sqrt{\log n} \right) \mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \quad (56)$$

$$\stackrel{(ii)}{\leq} O \left(\sigma \sqrt{\log n} \right) \left(\mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}}^2 \right)^{1/2} \quad (57)$$

$$= O \left(\sigma \sqrt{\log n} \right) \left(\mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\langle \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i), \sum_{j=1}^n \xi_j \phi(\mathbf{x}_j) \right\rangle_{\mathcal{H}} \right)^{1/2} \quad (58)$$

$$\stackrel{(iii)}{=} O\left(\sigma\sqrt{\log n}\right) \left(\mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \mathbb{E}_{\xi_i \sim \mathcal{R}} \sum_{i=1}^n \sum_{j=1}^n \xi_i \xi_j k(\mathbf{x}_i, \mathbf{x}_j) \right)^{1/2} \quad (59)$$

$$\stackrel{(iv)}{=} O\left(\sigma\sqrt{\log n}\right) \left(\mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \sum_{i=1}^n k(x_i, x_i) \right)^{1/2} \quad (60)$$

$$= O\left(\sigma\sqrt{n \log n \operatorname{Tr}(\mathbf{\Sigma})}\right) \quad (61)$$

(i) follows from applying Lemma 15. (ii) follows from Jensen's Inequality. (iii) follows from the definition of the kernel [Gre13]. (iv) holds as we have $\mathbb{E}[\xi_i \xi_j] = \delta_{i,j}$, where δ is the Kronecker Delta function. ■

Proposition 21 (Probabilistic bound on Norm of Functions with Gaussian Design in the Reproducing Kernel Hilbert Space). *Let $\mathbf{x}_i \sim \mathbb{P}$ such that $\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$ from Assumption 12. Then it follows*

$$\mathbb{P}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\{ \left\| \sum_{i=1}^n \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \geq \sqrt{n \operatorname{Tr}(\mathbf{\Sigma})} \cdot u \right\} \leq e^{-u^2/2} \quad (62)$$

Proof. Our proof will utilize a symmetrization argument similar to our previous expected covariance approximation proof, the proof then follows similarly to Lemma 19. on the Let C be a positive constant to be determined and $u \geq 1$. We then have

$$\mathbb{P}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\{ \left\| \sum_{i=1}^n \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \geq C \cdot u \right\} = \mathbb{P}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\{ \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \geq C \cdot u \right\} \quad (63)$$

$$\stackrel{(i)}{\leq} \mathbb{P}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\{ \sqrt{\mathbb{E}_{\xi_i \sim \mathcal{R}} \left[\left\| \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}}^2 \right]} \geq C \cdot u \right\} \quad (64)$$

$$= \mathbb{P}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\{ \sqrt{\mathbb{E}_{\xi_i \sim \mathcal{R}} \sum_{i=1}^n \sum_{j=1}^n \xi_i \xi_j k(\mathbf{x}_i, \mathbf{x}_j)} \geq C \cdot u \right\} \quad (65)$$

$$= \mathbb{P}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\{ \sqrt{\sum_{i=1}^n k(\mathbf{x}_i, \mathbf{x}_i)} \geq C \cdot u \right\} \quad (66)$$

$$\leq \mathbb{P}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\{ \sum_{i=1}^n \|\phi(\mathbf{x}_i)\|_{\mathcal{H}} \geq C \cdot u \right\} \quad (67)$$

$$\stackrel{(ii)}{\leq} \inf_{\theta > 0} \prod_{i=1}^n \prod_{j=1}^p \exp \left[\frac{1}{2} \theta^2 \lambda_j \right] \exp [-\theta C \cdot u] = \exp \left[-\frac{C^2 \cdot u^2}{2 \operatorname{Tr}(\mathbf{\Sigma})} \right] \quad (68)$$

In (i) we use Jensen's Inequality. For (ii) see the proof for Lemma 19 with an additional product term over n . Finally, setting $C = \sqrt{n \operatorname{Tr}(\mathbf{\Sigma})}$ completes the proof.

Lemma 22 (Infinite Dimensional Covariance Estimation in the Hilbert-Schmidt Norm). *Let $\mathbf{\Sigma} \triangleq \mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathbb{P}}[\phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i)]$. Then let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be i.i.d sampled from \mathbb{P} such that $\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$ from Assumption 12, we then have*

$$\mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\| \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) - \mathbf{\Sigma} \right\|_{\text{HS}} \leq C \left(\frac{\operatorname{Tr}(\mathbf{\Sigma})}{\sqrt{n}} \right) \quad (69)$$

where $C \leq 2\sqrt{3}$.

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,

$$\mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\| \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) - \mathbf{\Sigma} \right\|_{\text{HS}}$$

$$\stackrel{(i)}{\leq} \mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbb{E}_{\tilde{\phi}(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\| \frac{1}{n} \sum_{i=1}^n (\phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) - \tilde{\phi}(\mathbf{x}_i) \otimes \tilde{\phi}(\mathbf{x}_i)) \right\|_{\text{HS}} \quad (70)$$

$$= \mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbb{E}_{\tilde{\phi}(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \frac{1}{n} \sum_{i=1}^n \xi_i (\phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) - \tilde{\phi}(\mathbf{x}_i) \otimes \tilde{\phi}(\mathbf{x}_i)) \right\|_{\text{HS}} \quad (71)$$

$$\stackrel{(ii)}{\leq} \frac{2}{n} \mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right\|_{\text{HS}} \quad (72)$$

$$\stackrel{(iii)}{\leq} \frac{2}{n} \left(\mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right\|_{\text{HS}}^2 \right)^{1/2} \quad (73)$$

(i) follows from noticing $\phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) - \Sigma$ is a mean $\mathbf{0}$ operator in $\mathcal{H} \otimes \mathcal{H}$, then for $X, Y \in \mathcal{H} \otimes \mathcal{H}$ s.t. $\mathbb{E}[Y] = \mathbf{0}$ 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

$$\begin{aligned} & \frac{2}{n} \left(\mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right\|_{\text{HS}}^2 \right)^{1/2} \\ &= \frac{2}{n} \left(\mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \sum_{k=1}^p \left\langle \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) e_k, \sum_{j=1}^n \xi_j \phi(\mathbf{x}_j) \otimes \phi(\mathbf{x}_j) e_k \right\rangle_{\mathcal{H}} \right)^{1/2} \end{aligned} \quad (74)$$

$$= \frac{2}{n} \left(\mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbb{E}_{\xi_i \sim \mathcal{R}} \sum_{k=1}^p \sum_{i=1}^n \sum_{j=1}^n \xi_i \xi_j \langle \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) e_k, \phi(\mathbf{x}_j) \otimes \phi(\mathbf{x}_j) e_k \rangle_{\mathcal{H}} \right)^{1/2} \quad (75)$$

$$\stackrel{(iv)}{\leq} \frac{2}{n} \left(\mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \sum_{k=1}^p \sum_{i=1}^n \langle \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) e_k, \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) e_k \rangle_{\mathcal{H}} \right)^{1/2} \quad (76)$$

$$= \frac{2}{n} \left(\sum_{i=1}^n \mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \|\phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i)\|_{\text{HS}}^2 \right)^{1/2} \stackrel{(v)}{=} \frac{2}{n} \left(\sum_{i=1}^n \mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \|\phi(\mathbf{x}_i)\|_{\mathcal{H}}^4 \right)^{1/2} \quad (77)$$

$$= \frac{2}{n} \left(\sum_{i=1}^n \mathbb{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} [k^2(x_i, x_i)] \right)^{1/2} = \frac{2}{\sqrt{n}} \left(2 \text{Tr}(\Sigma^2) + \text{Tr}(\Sigma)^2 \right)^{1/2} \leq 2\sqrt{3}n^{-1/2} \text{Tr}(\Sigma) \quad (78)$$

(iv) follows from noticing $\mathbb{E}_{\xi_i, \xi_j \sim \mathcal{R}} [\xi_i \xi_j] = \delta_{ij}$. (v) follows from expanding the Hilbert-Schmidt Norm and applying Parseval's Identity. We 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

Proposition 23 (Probabilistic Bound on Infinite Dimensional Covariance Estimation in the Hilbert-Schmidt Norm). *Let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be i.i.d sampled from \mathbb{P} such that $\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)$ from Assumption 12, we then have for any $u \geq \sqrt{\text{Tr}(\Sigma)}$*

$$\mathbb{P}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \left\| \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) - \Sigma \right\|_{\text{HS}} \geq \frac{\text{Tr}(\Sigma)}{\sqrt{n}} \cdot u \right\} \leq e^{-u^2 \text{Tr}(\Sigma)/2} \quad (79)$$

Proof. Let C be a to be determined positive constant and u be a positive constant.

$$\mathbb{P}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\{ \left\| \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) - \mathbf{\Sigma} \right\|_{\text{HS}} \geq C \cdot u \right\} \quad (80)$$

$$\begin{aligned} &\leq \mathbb{P}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\{ \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \frac{1}{n} \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right\|_{\text{HS}} \right. \\ &\quad \left. + \mathbb{E}_{\tilde{\phi}(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \mathbb{E}_{\xi_i \sim \mathcal{R}} \left\| \frac{1}{n} \sum_{i=1}^n \xi_i \tilde{\phi}(\mathbf{x}_i) \otimes \tilde{\phi}(\mathbf{x}_i) \right\|_{\text{HS}} \geq C \cdot u \right\} \\ &\stackrel{(i)}{\leq} \mathbb{P}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\{ \frac{\sqrt{\text{Tr}(\mathbf{K})}}{n} + \frac{\text{Tr}(\mathbf{\Sigma})}{\sqrt{n}} \geq C \cdot u \right\} \end{aligned} \quad (81)$$

$$\leq \mathbb{P}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\{ \sum_{i=1}^n \|\phi(\mathbf{x}_i)\|_{\mathcal{H}} \geq nC \cdot u - \sqrt{n} \text{Tr}(\mathbf{\Sigma}) \right\} \quad (82)$$

$$\leq \exp \left[-\frac{(nC \cdot u + \sqrt{n} \text{Tr}(\mathbf{\Sigma}))^2}{n \text{Tr}(\mathbf{\Sigma})} \right] \stackrel{(ii)}{\leq} e^{-u^2 \text{Tr}(\mathbf{\Sigma})/2} \quad (83)$$

In (i) we apply Jensen's Inequality to both expectation terms and denote $\mathbf{K} \triangleq \mathbf{\Phi}^\top \mathbf{\Phi}$. In (ii) we chose $C \triangleq \text{Tr}(\mathbf{\Sigma})/\sqrt{n}$ and then simplify the resultant probability bound. ■

Lemma 24 (Finite Dimensional Covariate Estimation in the Spectral Norm). *Let $\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$. It then follows,*

$$\mathbb{E}_{\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top - \mathbf{\Sigma} \right\|_2 \leq C \|\mathbf{\Sigma}\| \left(\sqrt{\frac{d}{n}} + \frac{1}{\sqrt{dn}} \right) \quad (84)$$

where $C \leq 62.82$.

Proof. Our proof combines multiple results in High-Dimensional Probability for Sub-Gaussian vectors and adapting it for Gaussian-Design. We have,

$$\mathbb{E}_{\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top - \mathbf{\Sigma} \right\|_2 \leq \|\mathbf{\Sigma}\| \mathbb{E}_{\tilde{\mathbf{x}}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left\| \frac{1}{n} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} - \mathbf{I} \right\|_2 \quad (85)$$

$$= \|\mathbf{\Sigma}\| \int_0^\infty \mathbb{P}_{\tilde{\mathbf{x}}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left\{ \left\| \frac{1}{n} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} - \mathbf{I} \right\|_2 \geq t \right\} dt \quad (86)$$

Let \mathcal{M} be an ε -net of \mathbb{S}^{d-1} for $\varepsilon = \frac{1}{4}$, then $|\mathcal{M}| \leq 9^d$. It then follows from [Ver20] Corollary 4.2.13,

$$\mathbb{P}_{\tilde{\mathbf{x}}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left\{ \left\| \frac{1}{n} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} - \mathbf{I} \right\|_2 \geq t \right\} \leq \mathbb{P}_{\tilde{\mathbf{x}}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left\{ \max_{\mathbf{y} \in \mathcal{M}} \left| \frac{1}{n} \|\tilde{\mathbf{X}} \mathbf{y}\|_2^2 - 1 \right| \geq \frac{t}{2} \right\} \quad (87)$$

Denote $K \triangleq 16\sqrt{\frac{8}{3}}$, then we have from a union bound and Bernstein's Inequality [Ber24].

$$\mathbb{P}_{\tilde{\mathbf{x}}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left\{ \max_{\mathbf{y} \in \mathcal{M}} \left| \frac{1}{n} \|\tilde{\mathbf{X}} \mathbf{y}\|_2^2 - 1 \right| \geq \frac{t}{2} \right\} \leq 9^d \exp \left[-\frac{n}{2} \left(\frac{t^2}{K^2} \wedge \frac{t}{K} \right) \right] \quad (88)$$

Let $\delta \in (0, 1)$, then we find RHS Equation (88) is less than δ when

$$t \geq K \left(\frac{2d \log(9) + 2 \log(2/\delta)}{n} \vee \left(\frac{2d \log(9) + 2 \log(2/\delta)}{n} \right)^{1/2} \right) \quad (89)$$

Furthermore, we note we have equality with one, when $t = K\sqrt{(2d\log(9) + \log(4))/n} \triangleq C$ all $t \leq C$ occur with probability also equal to one. Therefore, plugging this back into the RHS of Equation (86).

$$\int_0^\infty \mathbb{P}_{\tilde{\mathbf{x}}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left\{ \left\| \frac{1}{n} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} - \mathbf{I} \right\|_2 \geq t \right\} dt \quad (90)$$

$$\leq K\sqrt{\frac{2d\log(9)}{n} + \frac{2\log(2)}{n}} + \int_C^\infty 9^d \exp\left[-\frac{n}{2} \left(\frac{t^2}{K^2}\right)\right] dt$$

$$\leq K\sqrt{\frac{2d\log(9)}{n} + \frac{2\log(2)}{n}} + \frac{K9^d \exp\left[-\left(K\sqrt{\frac{2d\log(9)}{n} + \frac{2\log(2)}{n}}\right)^2 \left(\frac{n}{2K^2}\right)\right]}{2n\sqrt{\frac{2d\log(9)}{n} + \frac{2\log(2)}{n}}} \quad (91)$$

$$\leq K\sqrt{\frac{2d\log(9) + \log(4)}{n}} + \frac{K}{4\sqrt{n(2d\log(9) + \log(4))}} \quad (92)$$

$$\leq \sqrt{\log(324)}K \left(\sqrt{\frac{d}{n}} + \frac{1}{\sqrt{nd}} \right) \quad (93)$$

In the second inequality, we use the integral inequality $\int_c^\infty e^{-x^2} dx \leq \int_c^\infty (\frac{x}{c})e^{-x^2} dx = e^{-c^2}/(2c)$. \blacksquare

B Proofs for Structural Results

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

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 .

$$\partial(-g(t, f_{\mathbf{w}})) = \partial\left(-t + \frac{1}{np} \sum_{i=1}^n (t - \hat{\nu}_i)^+\right) = -1 + \frac{1}{np} \sum_{i=1}^{np} \begin{cases} 1 & \text{if } t > \hat{\nu}_i \\ 0 & \text{if } t < \hat{\nu}_i \\ [0, 1] & \text{if } t = \hat{\nu}_i \end{cases} \quad (94)$$

We observe when setting $t = \hat{\nu}_{np}$, it follows that $0 \in \partial(-g(t, f_{\mathbf{w}}))$. This is equivalent to the p -quantile of the Risk. \blacksquare

B.2 Proof of Lemma 4

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

$$\nabla_f g(t^{(k+1)}, f_{\mathbf{w}}^{(k)}) = \nabla_f \left(t^{(k+1)} - \frac{1}{np} \sum_{i=1}^n \left(t^{(k+1)} - \ell(\mathbf{x}_i; f_{\mathbf{w}}^{(k)}, y_i) \right)^+ \right) \quad (95)$$

$$= -\frac{1}{np} \sum_{i=1}^{np} \nabla_f \left(t^{(k+1)} - \ell(\mathbf{x}_i; f_{\mathbf{w}}^{(k)}, y_i) \right)^+ = \frac{1}{np} \sum_{i=1}^n \nabla_f \ell(\mathbf{x}_i; f_{\mathbf{w}}^{(k)}, y_i) \begin{cases} 1 & \text{if } t > \hat{\nu}_i \\ 0 & \text{if } t < \hat{\nu}_i \\ [0, 1] & \text{if } t = \hat{\nu}_i \end{cases} \quad (96)$$

Now we note $\nu_{np} \leq t^{(k+1)} \leq \nu_{np+1}$. Then, plugging this into Equation (96), we have

$$\nabla_f g(t^{(k+1)}, f_{\mathbf{w}}^{(k)}) = \frac{1}{np} \sum_{i=1}^{np} \nabla_f \ell(\mathbf{x}_i; f_{\mathbf{w}}^{(k)}, y_i) \quad (97)$$

This concludes the proof. \blacksquare

C Proofs for Kernelized Regression

We will first give a simple calculation of the β -smoothness parameter of the subquantile objective. We then will give proofs for our approximation error bounds.

C.1 Subquantile Smoothness

Lemma 25. (β -Smoothness of $g(t, f_{\mathbf{w}})$ w.r.t $f_{\mathbf{w}}$). Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ represent the rows of the data matrix \mathbf{X} . It then follows:

$$\|\nabla_{f_{\mathbf{w}}} g(t, f_{\mathbf{w}}) - \nabla_{f_{\mathbf{w}}} g(t, f_{\tilde{\mathbf{w}}})\|_{\mathcal{H}} \leq \beta \|f_{\mathbf{w}} - f_{\tilde{\mathbf{w}}}\|_{\mathcal{H}} \quad (98)$$

where $\beta = \frac{2}{n(1-\epsilon)} \|\mathbf{K}\|$

Proof. We will upper bound the operator norm of the Hessian Operator. We have from Equation (6),

$$\begin{aligned} \|\nabla_{f_{\mathbf{w}}}^2 g(t, f_{\mathbf{w}})\|_{\text{op}} &= \frac{2}{n(1-\epsilon)} \left\| \sum_{i=1}^n \mathbb{I}\{t \geq \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i)\} \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right\|_{\text{op}} \\ &\leq \frac{2}{n(1-\epsilon)} \left\| \sum_{i=1}^n \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right\|_{\text{op}} = \frac{2}{n(1-\epsilon)} \|\Phi \Phi^\top\|_{\text{op}} = \frac{2}{n(1-\epsilon)} \|\mathbf{K}\| \end{aligned} \quad (99)$$

This completes the proof. ■

C.2 Proof of Lemma 6

Proof. We will first expand the expression in the Lemma statement. Let λ_i and φ_i for $i \in \mathbb{N}$ represent the eigenvalues and eigenfunctions for $\mathbb{E}_{\mathbf{x} \sim \mathcal{P}}[\phi(\mathbf{x}) \otimes \phi(\mathbf{x})] \triangleq \Sigma$.

$$\langle f_{\mathbf{w}} - f_{\mathbf{w}}^*, \Sigma(f_{\mathbf{w}} - f_{\mathbf{w}}^*) \rangle = \lim_{p \rightarrow \infty} \sum_{i=1}^p \lambda_i \langle f_{\mathbf{w}} - f_{\mathbf{w}}^*, \varphi_i \rangle_{\mathcal{H}}^2 \quad (100)$$

Therefore, for some $m \in \mathbb{N}$, we want the projection of $f_{\mathbf{w}} - f_{\mathbf{w}}^*$ to be non-zero for m . We will show, we only need to make an assumption on $f_{\mathbf{w}}^*$. Let Ψ represent the concatenation of the eigenfunctions φ_i for $i \in 1, \dots, p$ as $p \rightarrow \infty$. Furthermore, let Ψ_m represent the m eigenfunctions corresponding to $\lambda_1, \dots, \lambda_m$ and Ψ_m^\perp represents the eigenfunctions corresponding to λ_{m+1}, \dots . The projection in the Reproducing Kernel Hilbert Space is given as the following,

$$\|\text{Proj}_{\Psi_m} f_{\mathbf{w}}^*\|_{\mathcal{H}} \triangleq \left\| \sum_{i=1}^m \langle \varphi_i, f_{\mathbf{w}}^* \rangle \varphi_i \right\|_{\mathcal{H}} = \sum_{i=1}^m |\langle \varphi_i, f_{\mathbf{w}}^* \rangle_{\mathcal{H}}| \quad (101)$$

Let $f_{\mathbf{w}}^{(T)}$ be the T iterate from Subquantile Kernel Algorithm. From our assumption that $f_{\mathbf{w}}^* \in \text{Span}(\{\varphi_i\}_{i=1}^m) \triangleq \text{Span}(\Psi_m)$, it suffices to prove $\text{Proj}_{\Psi_m} f_{\mathbf{w}}^{(t)} \neq f_{\mathbf{w}}^*$ for all $t \in [T]$. In this case, we can note that $f_{\mathbf{w}} - f_{\mathbf{w}}^*$ will in some part be in the Span of Ψ_m .

$$\langle f_{\mathbf{w}}^{(t)} - f_{\mathbf{w}}^*, \Sigma(f_{\mathbf{w}}^{(t)} - f_{\mathbf{w}}^*) \rangle_{\mathcal{H}} \stackrel{(100)}{=} \lim_{p \rightarrow \infty} \sum_{i=1}^p \lambda_i \langle f_{\mathbf{w}}^{(t)} - f_{\mathbf{w}}^*, \varphi_i \rangle_{\mathcal{H}}^2 \quad (102)$$

$$\stackrel{(i)}{=} \sum_{i=1}^m \lambda_i \langle f_{\mathbf{w}}^{(t)} - f_{\mathbf{w}}^*, \varphi_i \rangle_{\mathcal{H}}^2 + \lim_{p \rightarrow \infty} \sum_{i=m+1}^p \lambda_i \langle f_{\mathbf{w}}^{(t)}, \varphi_i \rangle_{\mathcal{H}}^2 \quad (103)$$

$$= \sum_{i=1}^m \lambda_i \langle f_{\mathbf{w}}^{(t)} - f_{\mathbf{w}}^*, \varphi_i \rangle_{\mathcal{H}}^2 + \left\| \text{Proj}_{\Psi_m^\perp} f_{\mathbf{w}}^{(t)} \right\|_{\mathcal{H}}^2 \quad (104)$$

$$\stackrel{(ii)}{\geq} \lambda_m C \|f_{\mathbf{w}}^{(t)} - f_{\mathbf{w}}^*\|_{\mathcal{H}}^2 \quad (105)$$

Note in (i), since $\|\text{Proj}_{m, \perp} f_{\mathbf{w}}^*\|_{\mathcal{H}} = 0$, then it holds for all $i \in \mathbb{N}$ s.t. $i > m$ that $\langle f_{\mathbf{w}}^*, \varphi_i \rangle_{\mathcal{H}} = 0$. In (ii), we define $C \triangleq (\|f_{\mathbf{w}}^{(t)} - f_{\mathbf{w}}^*\|_{\mathcal{H}}^2 - \|\text{Proj}_{\Psi_m^\perp} f_{\mathbf{w}}^{(t)}\|_{\mathcal{H}}^2) / \|f_{\mathbf{w}}^{(t)} - f_{\mathbf{w}}^*\|_{\mathcal{H}}^2 > 0$ from our assumption that $f_{\mathbf{w}}^* \neq \text{Proj}_{\Psi_m} f_{\mathbf{w}}^{(t)}$ for all $t \in [T]$. This concludes the proof. ■

C.3 Proof of Theorem 7

Proof. From Algorithm 1, we have for kernelized linear regression

$$\text{Proj}_{\Psi_m} f_{\mathbf{w}}^{(t+1)} = \text{Proj}_{\Psi_m} \left[f_{\mathbf{w}}^{(t)} - 2\eta \sum_{i \in S \cap P} (f_{\mathbf{w}}(\mathbf{x}_i) - y_i) \cdot k(\mathbf{x}_i) \right] \quad (106)$$

Next, we note that we can partition $S = (S \cap P) \cup (S \cap Q) \triangleq \text{TP} \cup \text{FP}$. Then we have

$$\begin{aligned} \|\text{Proj}_{\Psi_m} f_{\mathbf{w}}^{(t+1)} - \text{Proj}_{\Psi_m} f_{\mathbf{w}}^*\|_{\mathcal{H}} &= \left\| \text{Proj}_{\Psi_m} f_{\mathbf{w}}^{(t)} - \text{Proj}_{\Psi_m} f_{\mathbf{w}}^* \right. \\ &\quad - \frac{2\eta}{n(1-\epsilon)} \text{Proj}_{\Psi_m} \sum_{i \in S \cap P} \left(f_{\mathbf{w}}^{(t)}(\mathbf{x}_i) - f_{\mathbf{w}}^*(\mathbf{x}_i) - \mu_i \right) \cdot k(\mathbf{x}_i, \cdot) \\ &\quad \left. - \frac{2\eta}{n(1-\epsilon)} \text{Proj}_{\Psi_m} \sum_{i \in S \cap Q} \left(f_{\mathbf{w}}^{(t)}(\mathbf{x}_i) - y_i \right) \cdot k(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}} \quad (107) \end{aligned}$$

We now note the following.

$$\sum_{i \in S \cap P} \left(f_{\mathbf{w}}^{(t)}(\mathbf{x}_i) - f_{\mathbf{w}}^*(\mathbf{x}_i) - \mu_i \right) \cdot k(\mathbf{x}_i, \cdot) = \left(f_{\mathbf{w}}^{(t)} - f_{\mathbf{w}}^* \right) \left[\sum_{i \in S \cap P} \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right] - \sum_{i \in S \cap P} \mu_i \phi(\mathbf{x}_i) \quad (108)$$

We will now bound the final term in the norm. Define $\xi_i \triangleq f_{\mathbf{w}}^*(\mathbf{x}_i) - y_i$ for all $i \in Q$. Then we have

$$\begin{aligned} \sum_{i \in S \cap Q} \left(f_{\mathbf{w}}^{(t)}(\mathbf{x}_i) - y_i \right) \cdot k(\mathbf{x}_i, \cdot) &= \sum_{i \in S \cap Q} \left(f_{\mathbf{w}}^{(t)}(\mathbf{x}_i) - f_{\mathbf{w}}^*(\mathbf{x}_i) - \xi_i \right) \cdot \phi(\mathbf{x}_i) \\ &= \left(f_{\mathbf{w}}^{(t)} - f_{\mathbf{w}}^* \right) \left[\sum_{i \in S \cap Q} \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right] - \sum_{i \in S \cap Q} \xi_i \phi(\mathbf{x}_i) \quad (109) \end{aligned}$$

For simplified notation, let $\Phi_{\text{TP}} \in \mathbb{R}^{\infty \times |\text{TP}|}$ represent the quasimatrix¹ containing the functions of $S \cap P$ and let $\Phi_{\text{FN}} \in \mathbb{R}^{\infty \times |\text{FN}|}$ represent the quasimatrix containing the functions in $P \setminus S$. Furthermore, let \mathbf{I} be the identity operator in \mathcal{H} . We then have from rearranging Equation (108) and Equation (109),

$$\begin{aligned} \|\text{Proj}_{\Psi_m} f_{\mathbf{w}}^{(t+1)} - \text{Proj}_{\Psi_m} f_{\mathbf{w}}^*\|_{\mathcal{H}} &\leq \|\text{Proj}_{\Psi_m} f_{\mathbf{w}}^{(t)} - \text{Proj}_{\Psi_m} f_{\mathbf{w}}^*\|_{\mathcal{H}} \left\| \mathbf{I} - \left(\frac{2\eta}{n(1-\epsilon)} \right) (\Phi_{\text{TP}} \Phi_{\text{TP}}^{\top} + \Phi_{\text{FP}} \Phi_{\text{FP}}^{\top}) \right\|_{\text{op}} \\ &\quad + \frac{2\eta}{n(1-\epsilon)} \left(\left\| \sum_{i \in S \cap P} \mu_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} + \left\| \sum_{i \in S \cap Q} \xi_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \right) \quad (110) \end{aligned}$$

Since we are making no distributional assumptions on the corrupted covariates, we can use the fact $\Phi_{\text{TP}} \Phi_{\text{TP}}^{\top} + \Phi_{\text{FP}} \Phi_{\text{FP}}^{\top} \succeq \Phi_{\text{TP}} \Phi_{\text{TP}}^{\top}$. Then, setting $\eta = \beta^{-1}$ (see Lemma 25), we obtain the following

$$\begin{aligned} \|\text{Proj}_{\Psi_m} f_{\mathbf{w}}^{(t+1)} - \text{Proj}_{\Psi_m} f_{\mathbf{w}}^*\|_{\mathcal{H}} &\leq \|\text{Proj}_{\Psi_m} f_{\mathbf{w}}^{(t)} - \text{Proj}_{\Psi_m} f_{\mathbf{w}}^*\|_{\mathcal{H}} \left(1 - \frac{\lambda_m(\Phi_{\text{TP}} \Phi_{\text{TP}}^{\top})}{\|\mathbf{K}\|} \right) \\ &\quad + \|\mathbf{K}\|^{-1} \left(\left\| \sum_{i \in S \cap P} \mu_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} + \left\| \sum_{i \in S \cap Q} \xi_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \right) \quad (111) \end{aligned}$$

Next, consider $\Lambda^{(t)} \triangleq \|\text{Proj}_{\Psi_m} f_{\mathbf{w}}^{(t)} - \text{Proj}_{\Psi_m} f_{\mathbf{w}}^*\|_{\mathcal{H}} + 2\|\sum_{i \in S \cap P} \mu_i \phi(\mathbf{x}_i)\|_{\mathcal{H}} + 2\|\sum_{i \in S \cap Q} \xi_i \phi(\mathbf{x}_i)\|_{\mathcal{H}}$. It then follows

$$\Lambda^{(t+1)} \leq \max \left\{ 1 - \frac{\lambda_m(\Phi_{\text{TP}} \Phi_{\text{TP}}^{\top})}{\|\mathbf{K}\|}, \frac{1}{\|\mathbf{K}\|} \right\} \cdot \Lambda^{(t)} \stackrel{(i)}{\leq} \left(1 - \frac{\lambda_m(\Phi_{\text{TP}} \Phi_{\text{TP}}^{\top})}{\|\mathbf{K}\|} \right) \cdot \Lambda^{(t)} \quad (112)$$

¹The quasimatrix is the infinite-dimensional generalization of a matrix (see e.g. [TT15]) and represents an ordered set of functions.

where in (i) we use the assumption $m > 1$. We have then showed for any $t \in [T]$, that both (i) $\|\text{Proj}_{\Psi_m} f_{\mathbf{w}}^{(t+1)} - \text{Proj}_{\Psi_m} f_{\mathbf{w}}^*\|_{\mathcal{H}} \leq \Lambda^{(t)}$ and (ii) $\Lambda^{(t+1)} \leq (1 - O(\|\mathbf{K}\|^{-1}) \cdot \Lambda^{(t)})$. Then, to obtain the probabilistic sample complexity we create two parameterized probabilistic events. Since $\mathbf{x}_i \in P$ are sampled i.i.d from \mathcal{P} , it implies that P is mutually independent. From which it follows that all subsets of P are independent.

$$E_u = \left\{ \left\| \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \leq \sqrt{\tilde{n} \text{Tr}(\mathbf{\Sigma})} \cdot u \text{ and } \left\| \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) - \mathbf{\Sigma} \right\|_{\text{HS}} \leq \sqrt{\tilde{n} \text{Tr}(\mathbf{\Sigma})} \cdot u \right\} \quad (113)$$

Invoking Proposition 21 and Proposition 23 together we find $\mathbb{P}\{E_u^c\} \leq e^{-\text{Tr}(\mathbf{\Sigma})u^2/2} + e^{-u^2/2}$ for all $u \geq 1$. Concatenate all the samples of μ_i for $i \in P$ in to a vector $\boldsymbol{\mu} \in \mathbb{R}^{(1-\epsilon)n}$ where $\boldsymbol{\mu}^+ \in \mathbb{R}^{\tilde{n}}$ denotes the samples of μ_i for $i \in S \cap P$ and $\boldsymbol{\mu}^- \in \mathbb{R}^{\gamma\tilde{n}}$ denotes the samples of μ_i for $i \in P \cap S$. We will create a parameterized event over $s \geq 1$.

$$E_s = \left\{ \boldsymbol{\mu} : \|\boldsymbol{\mu}^+\|_2^2 \leq \sigma^2 \gamma \tilde{n} \cdot s \text{ and } \|\boldsymbol{\mu}^-\|_{\infty} \leq \sigma \sqrt{2 \log(\gamma \tilde{n})} \cdot s \right\} \quad (114)$$

First note that $\boldsymbol{\mu} \sim \mathcal{N}(\mathbf{0}_{n(1-\epsilon)}, \sigma^2 \mathbf{I}_{n(1-\epsilon)})$, therefore for any subset of the indices $A \subseteq [n(1-\epsilon)]$, we have $\boldsymbol{\mu}_A \sim \mathcal{N}(\mathbf{0}_{|A|}, \sigma^2 \mathbf{I}_{|A|})$. We then invoke Proposition 16 and Proposition 17 together and obtain $\mathbb{P}\{E_s^c\} \leq e^{-s/2} + (\sqrt{2}/\log(\gamma \tilde{n}))e^{-s^2} \leq 2.05e^{-s/2}$ for all $s \geq 1$ and assuming $\gamma \tilde{n} \geq 2$. We then note from Weyl's Inequality [Wey12],

$$\lambda_m(\boldsymbol{\Phi}_{\text{TP}} \boldsymbol{\Phi}_{\text{TP}}^{\top}) = \lambda_m(\tilde{n} \mathbf{\Sigma} + \boldsymbol{\Phi}_{\text{TP}} \boldsymbol{\Phi}_{\text{TP}}^{\top} - \tilde{n} \mathbf{\Sigma}) \geq \tilde{n} \lambda_m(\mathbf{\Sigma}) - \left\| \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) - \mathbf{\Sigma} \right\|_{\text{HS}} \quad (115)$$

We then have for a constant $C \in (0, 1)$ with probability exceeding $1 - e^{-\text{Tr}(\mathbf{\Sigma})u^2/2}$, that $\lambda_m(\boldsymbol{\Phi}_{\text{TP}} \boldsymbol{\Phi}_{\text{TP}}^{\top}) \geq C \lambda_m(\mathbf{\Sigma})$ when

$$n \geq \frac{\text{Tr}(\mathbf{\Sigma})^2 \cdot u^2}{(1-C)^2(1-2\epsilon)\lambda_m^2} \quad (116)$$

Assume the dataset satisfies Equation (116). Suppose $f_{\mathbf{w}} = \mathbf{0}$, it then follows,

$$\begin{aligned} \mathbb{P}_{\mathcal{D} \sim \mathcal{P}} \left\{ \Lambda^{(0)} \leq \|f_{\mathbf{w}}^*\|_{\mathcal{H}} + \sigma \sqrt{2n(1-\epsilon) \log(n(1-\epsilon)) \text{Tr}(\mathbf{\Sigma})} \cdot su + \|\boldsymbol{\xi}\| \sqrt{\text{Tr}(\mathbf{K}_Q)} \right\} \\ \leq 1 - e^{u^2/2} - e^{-\text{Tr}(\mathbf{\Sigma})^2/2} - e^{-s^2} \end{aligned} \quad (117)$$

Concatenate the corruptions, ξ_i for $i \in Q$, into a vector $\boldsymbol{\xi}$. To obtain a ε accurate rank- m approximation, we require

$$T \geq \frac{\log \left(\|f_{\mathbf{w}}^*\| + \sigma \sqrt{n(1-\epsilon) \log(n(1-\epsilon)) \text{Tr}(\mathbf{\Sigma})} + \|\boldsymbol{\xi}\| \sqrt{\text{Tr}(\mathbf{K}_Q)} \right)}{\varepsilon \log \left(\frac{1}{1-C\lambda_m} \right)} \quad (118)$$

gradient iterations. Our proof is complete. ■