

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

Arvind Rathnashyam*

Alex Gittens[†]

September 12, 2023

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 [12, 21]. 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 $\inf_f \mathbb{E}_{\mathbf{x} \in \mathbb{P}} [\ell_f(\mathbf{x})]$, 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. We theoretically analyze Kernel Ridge Regression and Kernel Classification and empirically show the strength of Subquantile Minimization. Furthermore, we run experiments on various datasets and compare with the state-of-the-art algorithms to show the superior performance of Subquantile Minimization.

*CS, Rensselaer Polytechnic Institute, rathna@rpi.edu

[†]CS, Rensselaer Polytechnic Institute, gitttea@rpi.edu

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), [4, 1, 13, 17, 7] as well as for linear regression [2, 16]. Robust Statistics has been studied extensively [5], problems include high-dimensional mean estimation. Subquantile minimization aims to address the shortcomings of standard ERM in applications of noisy/corrupted data [11, 9]. In many real-world applications, linear models are insufficient to model the data. Therefore, we introduce the problem of Robust Learning for Kernel Learning.

Definition 1. (Huber ϵ -Contamination Model [8]). 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 $n\epsilon$ samples arbitrarily. The algorithm is then given the ϵ -corrupted data matrix \mathbf{X} and \mathbf{y} as training data.

Contributions

1. We propose a gradient-descent based algorithm for robust kernel learning in the Huber ϵ -Contamination Model which is fast.
2. We provide a theoretical analysis and give error bounds for kernel ridge regression and kernel classification.

1.1 Related Work

In this section we will describe previous works in robust algorithms for the Huber ϵ -Contamination Model and works in minimax optimization that will be relevant to our theoretical analysis.

Robust Algorithms

[4] proposed a robust meta-algorithm which filters points based their outlier likelihood score, which they define as the projection of the gradient of the point on to the top right singular vector of the Singular Value Decomposition of the Gradient of Losses. Empirically SEVER is strong in adversarially robust linear regression and Singular Vector Machines. SEVER however requires a base learner execution and SVD calculation for each iteration, thus it does not scale well for large scale applications.

[13] proposed optimization over the Tilted Empirical Loss. This is done by minimization of an exponentially weighted functional of the traditional Empirical Risk. Their involves a hyperparameter t , negative values of t trains more robustly, whereas positive values of t trains more fairly. This empirically works well in machine learning applications such as Noisy Annotation. The issue with introducing the exponential smoothing into the ERM function is the lack of interpretability.

[1] theoretically analyzed the Trimmed Maximum Likelihood Estimator algorithm in General Linear Models, including Gaussian Regression. They were able to show the Trimmed Maximum Likelihood Estimator achieves near optimal error for Gaussian Regression.

[3] studied empirical covariance estimation by gradient descent. They use gradient descent on a minimax formulation of the estimation problem. Their theoretical analysis is based upon the Moreau envelope. They prove their algorithm results in the norm of the gradient of the Moreau Envelope, and the ensuing \mathbf{w} is a good point in the search space. We tend to follow their general framework but we adapt it the Reproducing Kernel Hilbert Space Norm and for our minimax objective.

Minimax Optimization

[10] studied minimax optimization in the non-convex non-concave setting. Furthermore, they study convergence of alternating minimizing-maximizing algorithm with a maximizing oracle. Their research utilizes the Moreau Envelope.

[24] studied minimax optimization in the case of non-strong concavity.

1.2 Notation

The data matrix \mathbf{X} is a fixed $n \times d$ matrix, the matrix \mathbf{K} is the Gram Matrix, where $\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ and $k(\cdot, \cdot)$ represents a kernel function, e.g. Linear kernel: $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y}$, RBF kernel: $k(\mathbf{x}, \mathbf{y}) = \exp(-\gamma \|\mathbf{x} - \mathbf{y}\|_2^2)$. We denote $\mathbf{X}^\top = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$ where $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ represent the data vectors of the data matrix. We denote X as the set of all data vectors, $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$. We represent the data matrix $\mathbf{X} = (\mathbf{P}^\top \quad \mathbf{Q}^\top)^\top$, the labels vector as $\mathbf{y} = (\mathbf{y}_P^\top \quad \mathbf{y}_Q^\top)^\top$, and the dataset $X = P \cup Q = \{(\mathbf{x}_i, y_i)\}_{i=1}^n = \{\mathbf{x}_i, y_i\}_{i \in P} \cup \{\mathbf{x}_i, y_i\}_{i \in Q}$. We denote $\mathbf{I}_{k \times k}$ as the $k \times k$ identity matrix. The spectral norm of \mathbf{A} is $\|\mathbf{A}\|_2 = \max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\| = \sigma_{\max}(\mathbf{A})$. The reproducing Hilbert Space Norm of f is given as $\|f\|_{\mathcal{H}} \triangleq \mathbf{w}^\top \mathbf{K} \mathbf{w}$ where $f(\cdot) = \sum_{i=1}^n w_i k(\mathbf{x}_i, \cdot)$.

We also denote \triangleq as ‘defined as’, to be used when we are defining a variable. We will use $\stackrel{\text{def}}{=}$ to say a variable is defined as a quantity from previous literature.

Uppercase bold ($\mathbf{X}, \mathbf{Y}, \mathbf{Z}, \dots$) are matrices. Uppercase Roman are sets (X, S, P, Q). Lowercase bold are vectors ($\mathbf{x}, \mathbf{y}, \mathbf{z}, \dots$).

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 .

$$\mathcal{Q}_p(U) \leq t \iff \mathbb{P}\{U \leq t\} \geq p \quad (1)$$

The p -subquantile of the risk is then given by

$$\mathbb{L}_p(U) = \frac{1}{p} \int_0^p \mathcal{Q}_p(U) dq = \mathbb{E}[U | U \leq \mathcal{Q}_p(U)] = \max_{t \in \mathbb{R}} \left\{ t - \frac{1}{p} \mathbb{E}(t - U)^+ \right\} \quad (2)$$

Given a convex objective function, f , the learning problem becomes:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d: \|\mathbf{w}\| \leq R} \max_{t \in \mathbb{R}} \left\{ g(t, \mathbf{w}) \triangleq \sum_{i=1}^n (t - (f(\mathbf{x}_i; \mathbf{w}) - y_i)^2)^+ \right\} \quad (3)$$

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 [18]. Let $\Pi_{\mathcal{K}}$ be the projection of a vector on to the convex set $\mathcal{K} \triangleq \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_{\mathcal{H}} \leq R\}$, then our update steps are

$$t^{(k+1)} = \arg \max_{t \in \mathbb{R}} g(\mathbf{w}^{(k)}, t) \quad (4)$$

$$\mathbf{w}^{(k+1)} = \Pi_{\mathcal{K}} \left(\mathbf{w}^{(k)} - \alpha \nabla g(\mathbf{w}^{(k)}, t^{(k+1)}) \right) \quad (5)$$

We note this is a non-convex concave minimax optimization problem. We provide an algorithm for Subquantile Minimization of the ridge regression and classification kernel learning algorithm. ?? is applicable to both kernel ridge regression and kernel classification.

Algorithm 1: SUBQ-GRADIENT

Input: Iterations: T ; Quantile: p ; Data Matrix:

$\mathbf{X}, (n \times d), n \gg d$; Learning schedule:

$\alpha_1, \dots, \alpha_T$; Ridge parameter: λ

Output: Trained Parameters, $\mathbf{w}_{(T)}$

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

Algorithm 2: SUBQUANTILE

Input: Parameters \mathbf{w} , Data Matrix:

$\mathbf{X}, (n \times d)$, Convex Loss Function f

Output: Subquantile Matrix \mathbf{S}

- 1: $\hat{\nu}_i \leftarrow f(\mathbf{x}_i; \mathbf{w}, y_i)$ s.t. $\hat{\nu}_{i-1} \leq \hat{\nu}_i \leq \hat{\nu}_{i+1}$
 - 2: $t \leftarrow \hat{\nu}_{np}$
 - 3: Let $\mathbf{x}_1, \dots, \mathbf{x}_{np}$ be np points such that $f(\mathbf{x}_i; \mathbf{w}, y_i) \leq t$
 - 4: $\mathbf{S} \leftarrow (\mathbf{x}_1^\top \quad \dots \quad \mathbf{x}_{np}^\top)^\top$
 - 5: **return** \mathbf{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 Equations (4) and (5).

Lemma 2. *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 $\hat{\nu}_i \triangleq f(\mathbf{x}_i; \mathbf{w}, y_i)$, it then follows $\arg \max_{t \in \mathbb{R}} g(t, \mathbf{w}) = \hat{\nu}_{np}$.*

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 .

$$\partial(-g(t, \mathbf{w})) = \partial \left(-t + \frac{1}{np} \sum_{i=1}^n (t - \hat{\nu}_i) \right) \quad (6)$$

$$= -1 + \frac{1}{np} \sum_{i=1}^{np} \begin{cases} 1 & \text{if } t > \hat{\nu}_i \\ 0 & \text{if } t < \hat{\nu}_i \\ [0, 1] & \text{if } t = \hat{\nu}_i \end{cases} \quad (7)$$

$$= 0 \text{ when } t = \hat{\nu}_{np} \quad (8)$$

This is equivalent to the p -quantile of the Risk. ■

Interpretation 3. From lemma 2, 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 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 2, it then follows $\nabla_{\mathbf{w}} g(t^{(k)}, \mathbf{w}^{(k)}) = \frac{1}{np} \sum_{i=1}^{np} \nabla f(\mathbf{x}_i; \mathbf{w}^{(k)}, y_i)$*

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

$$\nabla_{\mathbf{w}} g(t^{(k+1)}, \mathbf{w}^{(k)}) = \nabla_{\mathbf{w}} \left(\hat{\nu}_{np} - \frac{1}{np} \sum_{i=1}^n (\hat{\nu}_{np} - f(\mathbf{x}_i; \mathbf{w}, y_i))^+ \right) \quad (9)$$

$$= -\frac{1}{np} \sum_{i=1}^{np} \nabla_{\mathbf{w}} (\hat{\nu}_{np} - f(\mathbf{x}_i; \mathbf{w}, y_i))^+ \quad (10)$$

$$= \frac{1}{np} \sum_{i=1}^n \nabla_{\mathbf{w}} f(\mathbf{x}_i; \mathbf{w}^{(k)}, y_i) \begin{cases} 1 & \text{if } t > \hat{\nu}_i \\ 0 & \text{if } t < \hat{\nu}_i \\ [0, 1] & \text{if } t = \hat{\nu}_i \end{cases} \quad (11)$$

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

$$\nabla_{\mathbf{w}} g(t^{(k+1)}, \mathbf{w}^{(k)}) = \frac{1}{np} \sum_{i=1}^{np} \nabla_{\mathbf{w}} f(\mathbf{x}_i; \mathbf{w}, y_i) \quad (12)$$

This concludes the proof. ■

3.1 Kernel Regression

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

$$f(\mathbf{x}, y_i; \mathbf{w}) = (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 \quad (13)$$

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

Lemma 5. (L -Lipschitz of $g(t, \mathbf{w})$ w.r.t \mathbf{w}). Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, represent the data vectors. It then follows:

$$|g(t, \mathbf{w}) - g(t, \hat{\mathbf{w}})| \leq L \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}} \quad (14)$$

where $L = \frac{2R}{np} \|\sum_{i=1}^n \mathbf{k}_i\|_{\mathcal{H}}^2 + \frac{2}{np} \|\sum_{i=1}^n \mathbf{k}_i\|_{\mathcal{H}} \|\mathbf{y}\|_2$

Lemma 6. (β -Smoothness of $g(t, \mathbf{w})$ w.r.t \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:

$$\|\nabla_{\mathbf{w}} g(t, \mathbf{w}) - \nabla_{\mathbf{w}} g(t, \hat{\mathbf{w}})\| \leq \beta \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}} \quad (15)$$

where $\beta = \frac{2}{np} \|\sum_{i \in X} \mathbf{k}_i\|_{\mathcal{H}}^2$

Proof. W.L.O.G, let S be the set of points such that if $\mathbf{x} \in S$, then $t \geq (\mathbf{k}_{\mathbf{x}}^{\top} \mathbf{w} - y)^2$. Since g is twice differentiable, we will analyze the Hessian.

$$\|\nabla_{\mathbf{w}}^2 g(t, \mathbf{w})\|_{\mathcal{H}} = \left\| \frac{2}{np} \sum_{i \in S} \mathbf{k}_i \mathbf{k}_i^{\top} \right\|_{\mathcal{H}} \stackrel{??}{\leq} \left\| \frac{2}{np} \sum_{i \in X} \mathbf{k}_i \mathbf{k}_i^{\top} \right\|_{\mathcal{H}} \leq \frac{2}{np} \left\| \sum_{i \in X} \mathbf{k}_i \right\|_{\mathcal{H}}^2 \stackrel{\text{lem. 11}}{\leq} \quad (16)$$

This concludes the proof. ■

3.2 Kernel Classification

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

$$f(\mathbf{x}_i, y_i; \mathbf{w}) = (1 - y_i (\mathbf{w}^{\top} \mathbf{k}_i))^+ \quad (17)$$

Similar to § section 3.1, we require the L -Lipschitz constant and β -smoothness constant.

Lemma 7. (L -Lipschitz of $g(t, \mathbf{w})$ w.r.t \mathbf{w}). Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, represent the data vectors. It then follows:

$$|g(t, \mathbf{w}) - g(t, \hat{\mathbf{w}})| \leq L \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}} \quad (18)$$

where $L = \frac{2R}{np} \|\sum_{i=1}^n \mathbf{k}_i\|_{\mathcal{H}}^2 + \frac{2}{np} \|\sum_{i=1}^n \mathbf{k}_i\|_{\mathcal{H}} \|\mathbf{y}\|_2$

Lemma 8. (β -Smoothness of $g(t, \mathbf{w})$ w.r.t \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:

$$\|\nabla_{\mathbf{w}} g(t, \mathbf{w}) - \nabla_{\mathbf{w}} g(t, \hat{\mathbf{w}})\| \leq \beta \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}} \quad (19)$$

where $\beta = \frac{2}{np} \|\sum_{i \in X} \mathbf{k}_i\|_{\mathcal{H}}^2$

Remark 9. Define the function $\Phi(\mathbf{w}) \triangleq \max_{t \in \mathbb{R}} g(t, \mathbf{w})$. This function is a L -weakly convex function, i.e., $\Phi(\mathbf{w}) + \frac{L}{2} \|\mathbf{w}\|^2$ is a convex function over \mathbf{w} .

3.3 Necessary Kernel Inequalities

We will first extend the idea of Resilience [23] to kernel learning.

Definition 10. (**Resilience**) from [23]. Let \mathcal{H} represent a RKHS, then given the feature mapping $\phi : \mathcal{X} \rightarrow \mathcal{H}$, and the set $X = \{\mathbf{x}_i\}_{i=1}^n = P \cup Q$, such that $|P| = n(1 - \epsilon)$ and $|Q| = n\epsilon$, it holds that for all $S \subseteq X$ s.t. $|S| \geq (1 - \epsilon)n$, then $\left\| \frac{1}{|S|} \sum_{i \in S} \phi(\mathbf{x}_i) - \mu \right\| \leq \tau$ then we say the set X has (ϵ, τ) -resilience in the Reproducing Kernel Hilbert Space.

Without the idea of resilience defined in definition 10, we will be unable to put error bounds on our algorithm.

Lemma 11. Let S be the set of elements in the subquantile, then

$$\left\| \frac{1}{np} \sum_{i \in X} \mathbf{k}_i \right\|_{\mathcal{H}} \leq \mathcal{O}() \quad (20)$$

Lemma 12. *Under the same setting as lemma 11,*

$$\left\| \frac{1}{np} \sum_{i \in S} \mathbf{k}_i \right\|_{\mathcal{H}} \leq \mathcal{O}() \quad (21)$$

We denote the matrix \mathbf{K} as the Gram Matrix where $K_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j) \triangleq \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|_2^2)$. Given a parameter set \mathbf{w} , the prediction for a new point will be: $f(\mathbf{x}^*; \mathbf{w}) = \sum_{i=1}^n \mathbf{w}_i \kappa(\mathbf{x}_i, \mathbf{x}^*)$

Theorem 13. *(Monotonically Decreasing). Let f be a convex loss function and \mathbf{X} follows ?? with learning schedule $\alpha_{(1)}, \alpha_{(2)}, \dots, \alpha_{(T)}$. Then it follows at any iteration $k \in \mathbb{N}$:*

$$g(t^{(k+1)}, \mathbf{w}^{(k+1)}) \leq g(t^{(k)}, \mathbf{w}^{(k)}) \quad (22)$$

From our definition of $S^{(k)}$ in theorem 13, we are interested in as $k \rightarrow \infty$ the quantities: $|\mathbf{x} \in S^{(k)} \cap P|$ and $|\mathbf{x} \in S^{(k)} \cap Q|$, where the latter cardinality represents the number of corrupted points in the subquantile set.

Definition 14. (Moreau Envelope, [15]). Let f be proper lower semi-continuous convex function $f : \mathcal{X} \rightarrow \mathbb{R}$, then the Moreau Envelope is defined as:

$$f_\lambda(\mathbf{x}) \triangleq \inf_{\hat{\mathbf{x}} \in \mathcal{X}} \left(f(\hat{\mathbf{x}}) + \frac{1}{2\lambda} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \right) \quad (23)$$

The Moreau Envelope can be interpreted as an infimal convolution of the function f with a quadratic.

Assumption 15. Define $\Phi(\cdot)$ as the function in remark 9. Then it follows $\arg \min_{\mathbf{w} \in \mathbb{R}^d} \Phi(\mathbf{w}) = \mathbf{w}^*$

Definition 16. (First-Order Stationary Point). Let $\Phi(\mathbf{w}) = \max_t g(t, \mathbf{w})$. Then \mathbf{w} is a first-order stationary point if

$$\nabla_{\mathbf{w}} \Phi(\mathbf{w})^\top (\tilde{\mathbf{w}} - \mathbf{w}) \geq 0 \quad \forall \tilde{\mathbf{w}} \in \mathcal{K} \quad (24)$$

Definition 17. (First-Order Stationary Point). Let $\Phi(\mathbf{w}) = \max_t g(t, \mathbf{w})$. Then \mathbf{w} is a first-order stationary point if

$$\|\nabla \Phi_\lambda(\mathbf{w})\|_{\mathcal{H}} = 0 \quad (25)$$

i.e.

$$\mathbf{w} = \arg \min_{\hat{\mathbf{w}} \in \mathcal{K}} \left(\Phi(\hat{\mathbf{w}}) + \frac{1}{2} \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}} \right) \quad (26)$$

Theorem 18. Let $\hat{\mathbf{w}}$ be a stationary point defined in definition 16 for the function Φ defined in ??. Then,

$$\|\hat{\mathbf{w}} - \mathbf{w}^*\|_{\mathcal{H}} \leq \mathcal{O}(\Xi) \quad (27)$$

Proof. First,

$$\|\nabla \Phi_\lambda(\mathbf{w})\|_{\mathcal{H}} = \left\| \frac{1}{\lambda} (\mathbf{w} - \text{prox}_{\lambda, \Phi}(\mathbf{w})) \right\|_{\mathcal{H}} = \left\| \frac{1}{\lambda} \left(\mathbf{w} - \arg \min_{\hat{\mathbf{w}} \in \mathcal{K}} \left(\Phi(\hat{\mathbf{w}}) + \frac{1}{2} \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}} \right) \right) \right\|_{\mathcal{H}} = 0 \quad (28)$$

This implies for any $\tilde{\mathbf{w}} \in \mathcal{K}$, it follows

$$\Phi(\hat{\mathbf{w}}) < \Phi(\tilde{\mathbf{w}}) + \frac{1}{2} \|\tilde{\mathbf{w}} - \hat{\mathbf{w}}\|_{\mathcal{H}} \quad (29)$$

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

$$\Phi(\hat{\mathbf{w}}) - \Phi(\mathbf{w}^*) \leq \frac{1}{2} \|\hat{\mathbf{w}} - \mathbf{w}^*\|_{\mathcal{H}} \quad (30)$$

■

Lemma 19. If $\|\mathbf{w} - \mathbf{w}^*\| \leq \eta$, then it follows

$$\Phi(\mathbf{w}) - \Phi(\mathbf{w}^*) = \mathcal{O}(\eta^2 \Xi^2 + 2\eta\rho\Xi) \quad (31)$$

The proof is deferred to § Appendix B.1. In practice, however, it is important to note that solving for $\|\nabla\Phi_\lambda\| = 0$ is NP-Hard. Thus, we will analyze the approximate stationary point.

Lemma 20. ([20, 19]). Assume the function Φ is ℓ -weakly convex. Let $\lambda < \frac{1}{\ell}$, and denote $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}'} \left(\Phi(\mathbf{w}') + \frac{1}{2\lambda} \|\mathbf{w} - \mathbf{w}'\|_{\mathcal{H}}^2 \right)$, $\|\nabla\Phi_\lambda(\mathbf{w})\|_{\mathcal{H}} \leq \epsilon$ implies:

$$\|\hat{\mathbf{w}} - \mathbf{w}\| = \lambda\epsilon \text{ and } \min_{\mathbf{g} \in \partial\Phi(\hat{\mathbf{w}}) + \partial\mathcal{I}_{\mathcal{K}}(\hat{\mathbf{w}})} \|\mathbf{g}\| \leq \epsilon \quad (32)$$

How to extend this to Hilbert Space Norm?

Note the subdifferential of the support function is the normal cone, i.e.

$$\partial\mathcal{I}_{\mathcal{K}} = \mathcal{N}(\hat{\mathbf{w}}) = \{\tilde{\mathbf{w}} \in \mathbb{R}^n | \langle \tilde{\mathbf{w}}, \bar{\mathbf{w}} - \hat{\mathbf{w}} \rangle \leq \forall \bar{\mathbf{w}} \in \mathcal{K}\} \quad (33)$$

We will define a convex cone.

Definition 21. A set Ω is a cone if $\lambda x \in \Omega$ whenever $x \in \Omega$ and $\lambda \geq 0$, if Ω is convex then it is a convex cone.

Thus there exists $\mathbf{g} = \mathbf{u} + \mathbf{v}$ where $\mathbf{u} \in \partial\Phi(\hat{\mathbf{w}})$ and $\mathbf{v} \in \mathcal{N}(\hat{\mathbf{w}})$.

Theorem 22. Let $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}'} \left(\Phi(\mathbf{w}') + \frac{1}{2\lambda} \|\mathbf{w} - \mathbf{w}'\|^2 \right)$ s.t. $\|\nabla\Phi_\lambda(\mathbf{w})\| \leq \epsilon$, then it follows

$$\|\hat{\mathbf{w}} - \mathbf{w}^*\|_{\mathcal{H}} \leq \Xi \quad (34)$$

Definition 23. (Approximate First-Order Stationary Point) from [3]. For any function f and closed convex set \mathcal{K} consider its associated Moreau envelope $f_\beta(\mathbf{w})$ in definition 14. Then we say that a point \mathbf{w} is a ρ -approximate stationary point if $\|f_\beta(\mathbf{w})\|_2 \leq \rho$.

The approximate stationary point in definition 23 is used in the analysis of the minimax algorithm in [3]. First, if you can prove a stationary point is good, theorem 18, then using lemma 20, you can show an approximate stationary point is good.

We adopt the proof strategy of [1] and [3], and have a two-part proof strategy. First we show an approximate stationary point is close to the true distribution of \mathbb{P} . Then, we analyze the optimization to show ?? 1 converges to an approximate stationary point in a polynomial number of iterations.

4 Optimization Results

Since we are solving a minimax objective, we want a relation between the norm of the gradient of the Moreau Envelope of Φ and $\left(\sum_{\mathbf{x} \in S^{(T)}} \mathbf{k}_{\mathbf{x}} (\mathbf{k}_{\mathbf{x}}^\top \mathbf{w}^{(T)} - y) \right)^\top \left(\frac{\mathbf{w}^{(T)} - \mathbf{w}^*}{\|\mathbf{w}^{(T)} - \mathbf{w}^*\|} \right)$. First, we will show using stepsize of $1/\beta$ returns a μ -approximate stationary point.

Theorem 24. (Algorithm ?? 1 reaches a η -approximate stationary point). Algorithm ?? 1 reaches a η -approximate stationary point in a polynomial number of iterations.

Proof. From [14] Theorem 31 and [3] Lemma 4.2, it follows:

$$\mathbb{E} \left[\|\nabla\Phi_{1/2\ell}(\bar{\mathbf{w}})\|^2 \right] \leq 2 \cdot \frac{(\Phi_{1/2\ell}(\mathbf{w}_0) - \min \Phi(\mathbf{w})) + \ell\beta^2\gamma^2}{\gamma\sqrt{T+1}} \quad (35)$$

where $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}'} \Phi(\mathbf{w}') + \ell \|\mathbf{w} - \mathbf{w}'\|^2$.

Let $\|\nabla\Phi_{1/2\ell}(\mathbf{w}^{(T)})\| \leq \mu$, it then follows from lemma 20, $\|\hat{\mathbf{w}} - \mathbf{w}^{(T)}\| = \mu/2\ell$. ■

4.1 Accelerated Gradient Methods

When working with big data it is often the case we need faster gradient methods as the gradient can be expensive to obtain. In this section, we give results on the convergence rate of accelerated gradient methods on the update of \mathbf{w} . We will analyze the convergence of three popular accelerated gradient methods.

4.1.1 Momentum Accelerated Gradient Descent

4.1.2 Conjugate Gradient Descent

4.1.3 Nesterov Accelerated Gradient Descent

5 Experiments

We perform numerical experiments on state of the art datasets comparing with other state of the art methods.

5.1 Kernel Regression

Algorithm 3: SUBQ-KERNEL-RIDGE-REGRESSION

Input: Iterations: T ; Quantile: p ; Data Matrix: \mathbf{X} , $(n \times d)$, $n \gg d$; Labels: \mathbf{y} , $(n \times 1)$; Learning schedule: $\alpha_1, \dots, \alpha_T$; Ridge parameter: λ

Output: Trained Parameters: $\mathbf{w}_{(T)}$; Base Learner: \mathcal{L}

```

1:  $\mathbf{w}_{(0)} \leftarrow (\mathbf{K}^\top \mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{K}^\top \mathbf{y}$  ▷ Base Learner
2: for  $k \in 1, 2, \dots, T$  do
3:    $\mathbf{S}_{(k)} \leftarrow \text{SUBQUANTILE}(\mathbf{w}^{(k)}, \mathbf{X})$  ▷ ?? 2
4:    $\nabla_{\mathbf{w}} g(t^{(k+1)}, \mathbf{w}^{(k)}) \leftarrow 2 \sum_{i \in S^{(k)}} \mathbf{k}_i (\mathbf{k}_i^\top \mathbf{w}^{(k)} - y_i) + \lambda \mathbf{K} \mathbf{w}^{(k)}$  ▷ Gradient Calculation
5:    $\mathbf{w}^{(k+1)} \leftarrow \mathbf{w}^{(k)} - \alpha_{(k)} \nabla_{\mathbf{w}} g(t^{(k+1)}, \mathbf{w}^{(k)})$  ▷  $\mathbf{w}$ -update in eqn. (5)
6: end
7: return  $\mathbf{w}_{(T)}$ 

```

Objectives	Test RMSE (Polynomial Regression (Degree = 3))			
	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon = 0.4$
KRR	0.460 _(0.2143)	1.171 _(0.7809)	0.950 _(0.3053)	1.230 _(0.4678)
TERM [13]	∞	∞	∞	∞
SEVER [4]	0.071 _(0.0106)	0.015 _(0.0041)	0.056 _(0.0513)	0.101 _(0.0643)
SUBQUANTILE($p = 1 - \epsilon$)	0.010_(0.0004)	0.010_(0.0002)	0.010_(0.0007)	0.012_(0.0030)
Genie ERM	∞	∞	∞	∞

Table 1: **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 $a_i \sim \mathcal{N}(0, 1)$. Oblivious Noise is sampled from $\mathcal{N}(0, 5)$. Subquantile is capped at 10,000 iterations. Polynomial Kernel: $K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top \mathbf{y} + 1)^3$. Regularization parameters is chosen as $\lambda = 1$. SEVER is trained with 16 iterations and $p = 0.02$.

In fig. 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. Looking at table 1 and figures fig. 1, subquantile minimization has near optimal performance in the **Polynomial Regression** Synthetic Dataset.

5.2 Kernel 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

Objectives	Test RMSE (Boston Housing Regression)			
	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon = 0.4$
KRR	0.544 _(0.0712)	0.862 _(0.1199)	0.865 _(0.0811)	1.049 _(0.2249)
TERM [13]	0.888 _(0.1360)	0.891 _(0.1699)	1.023 _(0.1329)	0.931 _(0.0433)
SEVER [4]	0.593 _(0.0478)	0.573 _(0.0559)	0.567 _(0.1191)	∞
SUBQUANTILE($p = 1 - \epsilon$)	0.427_(0.0691)	0.534_(0.1105)	0.510_(0.0695)	0.549_(0.1030)
Genie ERM	∞	∞	∞	∞

Table 2: **Boston Housing Regression** Dataset. Oblivious Noise is sampled from $\mathcal{N}(0, 5)$. Subquantile is capped at 10,000 iterations. Regularization Parameter is chosen as $\lambda = 2$

Objectives	Test RMSE (Concrete Data Regression)			
	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon = 0.4$
KRR	0.802 _(0.0324)	0.929 _(0.0209)	0.993 _(0.0441)	0.775 _(0.0514)
TERM [13]	0.874 _(0.0205)	0.916 _(0.0421)	0.840 _(0.0249)	0.878 _(0.0749)
SEVER [4]	0.532 _(0.0134)	0.516 _(0.0340)	0.526 _(0.0217)	0.552 _(0.0444)
SUBQUANTILE($p = 1 - \epsilon$)	0.468_(0.0220)	0.491_(0.0271)	0.555_(0.0391)	0.566_(0.0405)
Genie ERM	∞	∞	∞	∞

Table 3: **Concrete Data Regression** Dataset. Oblivious Noise is sampled from $\mathcal{N}(0, 5)$. Subquantile is capped at 10,000 iterations. Regularization Parameter is chosen as $\lambda = 2$.

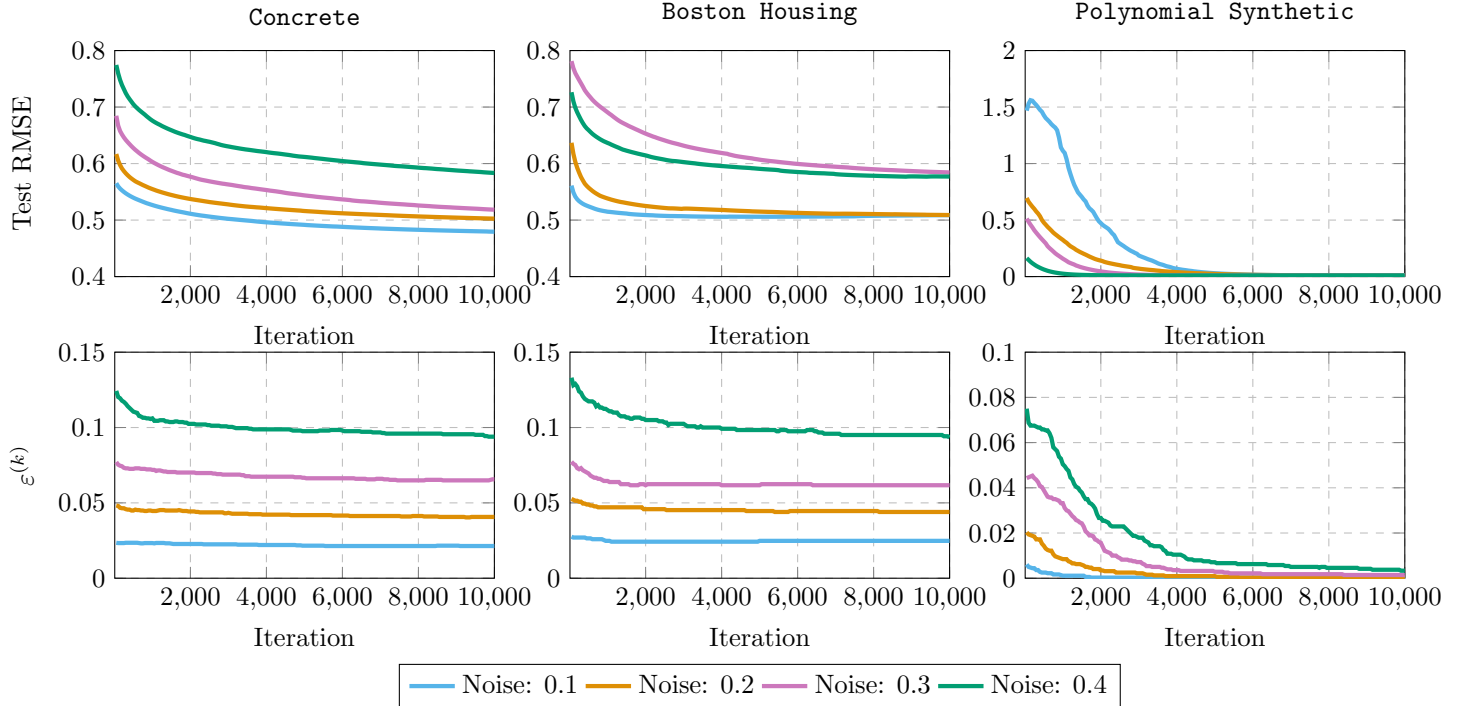


Figure 1: Test RMSE over the iterations in **Concrete**, **Boston Housing**, and **Polynomial** Datasets for SUBQUANTILE at different noise levels

art robust algorithms.

Algorithm 4: SUBQ-KERNEL-CLASSIFICATION

Input: Iterations: T ; Quantile: p ; Data Matrix: \mathbf{X} , $(n \times d)$, $n \gg d$; Labels: \mathbf{y} , $(n \times 1)$; Learning schedule: $\alpha_1, \dots, \alpha_T$; Ridge parameter: λ

Output: Trained Parameters: $\mathbf{w}_{(T)}$; Base Learner: \mathcal{L}

```
1:  $\mathbf{w}_{(0)} \leftarrow (\mathbf{K}^\top \mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{K}^\top \mathbf{y}$  ▷ Base Learner
2: for  $k \in 1, 2, \dots, T$  do
3:    $\mathbf{S}_{(k)} \leftarrow \text{SUBQUANTILE}(\mathbf{w}^{(k)}, \mathbf{X})$  ▷ Algorithm ?? 2
4:    $\nabla_{\mathbf{w}} g(t^{(k+1)}, \mathbf{w}^{(k)}) \leftarrow -\sum_{i \in S^{(k)}} y_i \mathbf{k}_i + \lambda \mathbf{K} \mathbf{w}^{(k)}$  ▷ Gradient Calculation
5:    $\mathbf{w}^{(k+1)} \leftarrow \mathbf{w}^{(k)} - \alpha_{(k)} \nabla_{\mathbf{w}} g(t^{(k+1)}, \mathbf{w}^{(k)})$  ▷  $\mathbf{w}$ -update in eqn. (5)
6: end
7: return  $\mathbf{w}_{(T)}$ 
```

6 Discussion

The main contribution of this paper is the study of a nonconvex-concave optimization algorithm for the robust learning problem for kernel ridge regression and kernel classification.

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. 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.

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.

In future work, the analysis of Subquantile Minimization can be extended to neural networks and other learning algorithms.

References

- [1] Pranjal 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.
- [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.
- [3] 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.
- [4] 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.
- [5] Ilias Diakonikolas and Daniel M Kane. *Algorithmic high-dimensional robust statistics*. Cambridge University Press, 2023.
- [6] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
- [7] 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.
- [8] Peter J. Huber and Elvezio. Ronchetti. *Robust statistics*. Wiley series in probability and statistics. Wiley, Hoboken, N.J., 2nd ed. edition, 2009.
- [9] 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.
- [10] 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.
- [11] Ashish Khetan, Zachary C. Lipton, and Anima Anandkumar. Learning from noisy singly-labeled data. In *International Conference on Learning Representations*, 2018.
- [12] 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.
- [13] Tian Li, Ahmad Beirami, Maziar Sanjabi, and Virginia Smith. Tilted empirical risk minimization. In *International Conference on Learning Representations*, 2021.
- [14] Tianyi Lin, Chi Jin, and Michael I. Jordan. Near-optimal algorithms for minimax optimization. In Jacob Abernethy and Shivani Agarwal, editors, *Proceedings of Thirty Third Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pages 2738–2779. PMLR, 09–12 Jul 2020.
- [15] Jean-Jacques Moreau. Proximité et dualité dans un espace hilbertien. *Bulletin de la Société mathématique de France*, 93:273–299, 1965.

- [16] 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.
- [17] 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.
- [18] 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.
- [19] R Tyrrell Rockafellar. *Convex analysis*. 2015.
- [20] Ralph Tyrell Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, 1970.
- [21] 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.
- [22] Markus Schneider. Probability inequalities for kernel embeddings in sampling without replacement. In *International Conference on Artificial Intelligence and Statistics*, 2016.
- [23] Jacob Steinhardt, Moses Charikar, and Gregory Valiant. Resilience: A criterion for learning in the presence of arbitrary outliers. In Anna R. Karlin, editor, *9th Innovations in Theoretical Computer Science Conference, ITCS 2018, January 11-14, 2018, Cambridge, MA, USA*, volume 94 of *LIPIcs*, pages 45:1–45:21. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2018.
- [24] Junchi Yang, Antonio Orvieto, Aurelien Lucchi, and Niao He. Faster single-loop algorithms for minimax optimization without strong concavity. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera, editors, *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pages 5485–5517. PMLR, 28–30 Mar 2022.

Contents

1	Introduction	2
1.1	Related Work	2
1.2	Notation	3
2	Subquantile Minimization	3
3	Structural Results	4
3.1	Kernel Regression	4
3.2	Kernel Classification	5
3.3	Necessary Kernel Inequalities	5
4	Optimization Results	7
4.1	Accelerated Gradient Methods	8
4.1.1	Momentum Accelerated Gradient Descent	8
4.1.2	Conjugate Gradient Descent	8
4.1.3	Nesterov Accelerated Gradient Descent	8
5	Experiments	8
5.1	Kernel Regression	8
5.2	Kernel Classification	8
6	Discussion	10
A	Kernel Embedding Inequalities	14
A.1	Proof of Lemma 11	14
A.2	Proof of Lemma 12	14
B	Deferred Proofs	15
B.1	Proof of Lemma 19	15
B.2	Proof of Theorem 18	15
B.3	Proof of Theorem 13	16
B.4	Proof of Lemma 5	17
B.5	Proof of Lemma 7	17
C	Base Learner Algorithm	18
D	Experimental Details	19
E	Detailed Related Works	20
E.1	High-dimensional Robust Mean Estimation via Gradient Descent [3]	20
E.2	Trimmed Maximum Likelihood Estimation for Robust Generalized Linear Model [1]	20

A Kernel Embedding Inequalities

Lemma 25. (*Resilience on Inlier Samples*). Let $X = \{\mathbf{x}_i\}_{i=1}^n$ and $P = \{\mathbf{x}_i\}_{i=1}^{np}$, and $[\mathbf{k}_i]_j = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$. If the conditions in assumption definition 10 it then follows:

$$\mathbb{P} \left\{ \left\| \frac{1}{np} \sum_{i \in P} \phi(\mathbf{x}) - \mu_{\mathbb{P}} \right\| > \epsilon \right\} \leq 2 \exp \left(-\frac{2np\epsilon^2}{d^2} \right)$$

where $\|\phi(\mathbf{x})\| \leq d$ for any $\mathbf{x} \in \mathcal{X}$ a.s. A similar inequality can be found in [22], Theorem 2

Lemma 26. Let $\{\xi_i\}_{i=1}^n$ represent realizations of a random variable, $\xi_i \sim \mathcal{N}(\mu, \sigma^2)$. It then follows with high probability:

$$\left\| \sum_{i=1}^n \xi_i \right\| \leq C \text{ with high probability} \quad (36)$$

First, we can note,

$$\sum_{i=1}^n \xi_i \sim \mathcal{N}(n\mu, n\sigma^2) \quad (37)$$

With Hoeffding's Inequality, we have:

$$\mathbb{P} \left\{ \left\| \frac{1}{n} \sum_{i=1}^n \xi_i - \mu \right\| \geq \epsilon \right\} \leq 2 \exp \left(-\frac{2n^2\epsilon^2}{\sum_{i=1}^n \|\xi_i\|_{\psi_2}^2} \right) \quad (38)$$

A.1 Proof of Lemma 11

Proof.

$$\left\| \frac{1}{np} \sum_{i \in S \cap P} \mathbf{k}_i \right\| = \frac{1}{np} \left\| \sum_{i \in S \cap P} \sum_{j \in X} \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) \right\| \quad (39)$$

■

A.2 Proof of Lemma 12

Proof.

■

B Deferred Proofs

In this section we give some deferred proofs.

B.1 Proof of Lemma 19

Proof. Let S be the set containing the points with minimum error from X w.r.t to the weights vector \mathbf{w} .

$$\|\Phi(\mathbf{w}) - \Phi(\mathbf{w}^*)\| = \left\| \sum_{i \in S} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 - \sum_{j \in P} (\mathbf{w}^{*\top} \mathbf{k}_j - y_j)^2 \right\| \quad (40)$$

$$= \left\| \sum_{i \in S \cap P} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 + \sum_{i \in S \cap Q} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 - \sum_{j \in P} (\mathbf{w}^{*\top} \mathbf{k}_j - y_j)^2 \right\| \quad (41)$$

$$\stackrel{(a)}{\leq} \left\| \sum_{i \in S \cap P} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 + \sum_{i \in P \setminus S} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 - \sum_{j \in P} (\mathbf{w}^{*\top} \mathbf{k}_j - y_j)^2 \right\| \quad (42)$$

$$= \left\| \sum_{i \in P} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 - \sum_{j \in P} (\mathbf{w}^{*\top} \mathbf{k}_j - y_j)^2 \right\| \quad (43)$$

$$= \left\| \sum_{i \in P} (\mathbf{w}^\top \mathbf{k}_i)^2 - 2y_i \mathbf{w}^\top \mathbf{k}_i - (\mathbf{w}^{*\top} \mathbf{k}_i)^2 + 2y_i (\mathbf{w}^{*\top} \mathbf{k}_i) \right\| \quad (44)$$

$$= \left\| \sum_{i \in P} \left((\mathbf{w} - \mathbf{w}^*)^\top \mathbf{k}_i \right)^2 + (1 - 2y_i) (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{k}_i - 2 \|\mathbf{w}^{*\top} \mathbf{k}_i\|^2 \right\| \quad (45)$$

$$\stackrel{\text{Cauchy-Schwarz}}{\leq} \sum_{i \in P} \|\mathbf{w} - \mathbf{w}^*\|^2 \|\mathbf{k}_i\|^2 + |1 - 2y_i| \|\mathbf{w} - \mathbf{w}^*\| \|\mathbf{k}_i\| - \|\mathbf{w}^{*\top} \mathbf{k}_i\|^2 \quad (46)$$

$$= \eta^2 \left(\sum_{i \in P} \|\mathbf{k}_i\|^2 \right) + 2\eta \|\mathbf{y}_P\| \left(\sum_{i \in P} \|\mathbf{k}_i\| \right) - \left(\sum_{i \in P} \|\mathbf{w}^{*\top} \mathbf{k}_i\|^2 \right) \quad (47)$$

$$= \mathcal{O}(\eta^2 \Xi^2 + 2\eta \rho \Xi) \quad (48)$$

(a) follows since the points outside the subquantile set in P will have greater error than the points inside the subquantile in Q by the formation of the subquantile set in Equation (4).

This requires distributional assumptions on the elements in P , from which we can derive probabilistic and expected inequalities on the sum of the norms. \blacksquare

B.2 Proof of Theorem 18

Proof. Let $i \in [np]$ be the np indices with the lowest error, then we recall the derivative of Φ is given by,

$$\nabla_{\mathbf{w}} \Phi(\hat{\mathbf{w}}) = \sum_{i=1}^{np} \mathbf{k}_i (\mathbf{k}_i^\top \hat{\mathbf{w}} - y_i) \quad (49)$$

We will do a proof by contrapositive. Assume $\|\hat{\mathbf{w}} - \mathbf{w}^*\|_{\mathcal{H}} > \mathcal{O}(\Xi)$. We will prove $\hat{\mathbf{w}}$ is not a stationary point. Note, the stationary point $\hat{\mathbf{w}}$ satisfies the following property,

$$\nabla_{\mathbf{w}} \Phi(\hat{\mathbf{w}})^\top \hat{\mathbf{w}} \geq \nabla_{\mathbf{w}} \Phi(\hat{\mathbf{w}})^\top \tilde{\mathbf{w}} \quad \forall \hat{\mathbf{w}} \in \mathcal{K} \quad (50)$$

Thus it suffices to prove there exists $\mathbf{w} \in \mathcal{K}$ s.t. Equation (50) is satisfied.

$$\nabla_{\mathbf{w}} \Phi(\hat{\mathbf{w}})^\top \tilde{\mathbf{w}} = \left(\sum_{i=1}^{np} \mathbf{k}_i (\mathbf{k}_i^\top \hat{\mathbf{w}} - y_i) \right)^\top \hat{\mathbf{w}} \quad (51)$$

For what Ξ does there exist $\tilde{\mathbf{w}} \in \mathcal{K}$ s.t.

$$\left(\sum_{i=1}^{np} \mathbf{k}_i (\mathbf{k}_i^\top \hat{\mathbf{w}} - y_i) \right)^\top \hat{\mathbf{w}} \geq \left(\sum_{i=1}^{np} \mathbf{k}_i (\mathbf{k}_i^\top \hat{\mathbf{w}} - y_i) \right)^\top \tilde{\mathbf{w}} \quad (52)$$

This is equivalent to

$$\sum_{i=1}^{np} (\mathbf{k}_i^\top \hat{\mathbf{w}} - y_i) \mathbf{k}_i^\top \hat{\mathbf{w}} \geq \sum_{i=1}^{np} (\mathbf{k}_i^\top \hat{\mathbf{w}} - y_i) \mathbf{k}_i^\top \tilde{\mathbf{w}} \quad (53)$$

■

B.3 Proof of Theorem 13

Proof. We will introduce new notation, let $S^{(k)}$ denote the set of np data points from \mathbf{X} with the lowest objective value $f(\mathbf{x}; \mathbf{w}^{(k)}, y) \triangleq (f(\mathbf{x}; \mathbf{w}^{(k)}) - y)^2$. We will also define $F(\mathbf{w}, S) \triangleq \sum_{\mathbf{x} \in S} (f(\mathbf{x}; \mathbf{w}^{(k)}) - y)^2$. Note this is an equivalent characterization of g from lemma 2.

$$F(\mathbf{w}^{(k+1)}, S^{(k+1)}) \leq F(\mathbf{w}^{(k)}, S^{(k)}) \quad (54)$$

$$F(\mathbf{w}^{(k+1)}, S^{(k+1)}) - F(\mathbf{w}^{(k)}, S^{(k+1)}) \leq F(\mathbf{w}^{(k)}, S^{(k)}) - F(\mathbf{w}^{(k)}, S^{(k+1)}) \quad (55)$$

Upper Bound of LHS

Note that Kernel Ridge Regression is Lipschitz Continuous, let L denote the Lipschitz-Constant.

$$F(\mathbf{w}^{(k+1)}, S^{(k+1)}) - F(\mathbf{w}^{(k)}, S^{(k+1)}) \leq \langle \nabla_{\mathbf{w}} F(\mathbf{w}^{(k)}, S^{(k+1)}), \mathbf{w}^{(k+1)} - \mathbf{w}^{(k)} \rangle + \frac{L}{2} \left\| \mathbf{w}^{(k+1)} - \mathbf{w}^{(k)} \right\|_2^2 \quad (56)$$

$$= \langle \nabla_{\mathbf{w}} F(\mathbf{w}^{(k)}, S^{(k+1)}), -\alpha^{(k)} \nabla_{\mathbf{w}} F(\mathbf{w}^{(k)}, S^{(k+1)}) \rangle + \frac{L}{2} \left\| \mathbf{w}^{(k+1)} - \mathbf{w}^{(k)} \right\|_2^2 \quad (57)$$

$$= -\alpha^{(k)} \left\| \nabla_{\mathbf{w}} F(\mathbf{w}^{(k)}, S^{(k+1)}) \right\|_2^2 + \frac{L}{2} \left\| \mathbf{w}^{(k+1)} - \mathbf{w}^{(k)} \right\|_2^2 \quad (58)$$

$$= -\alpha^{(k)} \left\| \nabla_{\mathbf{w}} F(\mathbf{w}^{(k)}, S^{(k+1)}) \right\|_2^2 + \frac{L\alpha_{(k)}^2}{2} \left\| \nabla_{\mathbf{w}} F(\mathbf{w}^{(k)}, S^{(k+1)}) \right\|_2^2 \quad (59)$$

$$= \left(\frac{L\alpha_{(k)}^2}{2} - \alpha_{(k)} \right) \left\| \nabla_{\mathbf{w}} F(\mathbf{w}^{(k)}, S^{(k+1)}) \right\|_2^2 \quad (60)$$

Lower Bound of RHS

We first note that the points in $S^{(k+1)}$ and not in $S^{(k)}$ have lower residuals.

$$F(\mathbf{w}^{(k)}, S^{(k)}) - F(\mathbf{w}^{(k)}, S^{(k+1)}) = \sum_{\mathbf{x} \in S^{(k)} \setminus S^{(k+1)}} f(\mathbf{x}; \mathbf{w}^{(k)}, y) - \sum_{\mathbf{x} \in S^{(k+1)} \setminus S^{(k)}} f(\mathbf{x}; \mathbf{w}^{(k)}, y) \quad (61)$$

$$\geq |S^{(k)} \setminus S^{(k+1)}| \inf \left\{ f(\mathbf{x}; \mathbf{w}^{(k)}, y) : \mathbf{x} \in S^{(k)} \setminus S^{(k+1)} \right\} - |S^{(k+1)} \setminus S^{(k)}| \sup \left\{ f(\mathbf{x}; \mathbf{w}^{(k)}, y) : \mathbf{x} \in S^{(k+1)} \setminus S^{(k)} \right\} \quad (62)$$

$$\text{Let } \eta \triangleq |S^{(k)} \setminus S^{(k+1)}| = |S^{(k+1)} \setminus S^{(k)}|$$

$$= \eta \left(\inf \left\{ f(\mathbf{x}; \mathbf{w}^{(k)}, y) : \mathbf{x} \in S^{(k)} \setminus S^{(k+1)} \right\} - \sup \left\{ f(\mathbf{x}; \mathbf{w}^{(k)}, y) : \mathbf{x} \in S^{(k+1)} \setminus S^{(k)} \right\} \right) \quad (63)$$

$$= \eta \left(\hat{\nu}_{np+1}^{(k)} - \hat{\nu}_{np}^{(k)} \right) \quad (64)$$

$$\geq 0 \quad (65)$$

Therefore, if $\alpha_{(k)} \leq \frac{2}{L}$, then it follows $F(\mathbf{w}^{(k+1)}, S^{(k+1)}) \leq F(\mathbf{w}^{(k)}, S^{(k)})$. This concludes the proof as we shown descent lemma. ■

B.4 Proof of Lemma 5

Proof. We use the \mathcal{H} norm of the gradient to bound L from above.

$$\|\nabla_{\mathbf{w}} g(t, \mathbf{w})\|_{\mathcal{H}} = \left\| \frac{2}{np} \sum_{i=1}^n \mathbb{1}_{t \geq (\mathbf{k}_i^\top \mathbf{w} - y_i)^2} (\mathbf{k}_i (\mathbf{k}_i^\top \mathbf{w} - y_i)) \right\|_{\mathcal{H}} \quad (66)$$

W.L.O.G, let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$ where $0 \leq m \leq n$, represent the data vectors such that $t \geq (\mathbf{k}_i^\top \mathbf{w} - y_i)^2$.

$$= \left\| \frac{2}{np} \sum_{i=1}^m \mathbf{k}_i (\mathbf{k}_i^\top \mathbf{w} - y_i) \right\|_{\mathcal{H}} \quad (67)$$

$$\stackrel{(a)}{\leq} \frac{2}{np} \left(\left\| \sum_{i=1}^m \mathbf{k}_i (\mathbf{k}_i^\top \mathbf{w}) \right\|_{\mathcal{H}} + \left\| \sum_{i=1}^m \mathbf{k}_i y_i \right\|_{\mathcal{H}} \right) \quad (68)$$

$$\stackrel{(b)}{\leq} \frac{2}{np} \left(\left\| \sum_{i=1}^m \mathbf{k}_i \right\|_{\mathcal{H}} \left\| \sum_{i=1}^m \mathbf{k}_i^\top \mathbf{w} \right\|_{\mathcal{H}} + \left\| \sum_{i=1}^m \mathbf{k}_i \right\|_{\mathcal{H}} \left(\sum_{i=1}^m y_i^2 \right)^{1/2} \right) \quad (69)$$

$$\stackrel{(c)}{\leq} \frac{2}{np} \left(\left\| \sum_{i=1}^m \mathbf{k}_i \right\|_{\mathcal{H}}^2 \|\mathbf{w}\|_{\mathcal{H}} + \left\| \sum_{i=1}^m \mathbf{k}_i \right\|_{\mathcal{H}} \|\mathbf{y}\|_2 \right) \quad (70)$$

$$\stackrel{(d)}{\leq} \frac{2R}{np} \left\| \sum_{i=1}^{np} \mathbf{k}_i \right\|_{\mathcal{H}}^2 + \frac{2}{np} \left\| \sum_{i=1}^{np} \mathbf{k}_i \right\|_{\mathcal{H}} \|\mathbf{y}\|_2 \quad (71)$$

where (a) follows from Triangle Inequality, (b) and (c) follow from Cauchy-Schwarz Inequality, (d) follows from assuming $\|\mathbf{w}\|_{\mathcal{H}} \leq R$, where $R \in \mathbb{R}_+ < \infty$ is some positive constant. This concludes the proof. \blacksquare

B.5 Proof of Lemma 7

Proof. We use the \mathcal{H} norm of the gradient to bound L from above. Let S be denoted as the subquantile set.

$$\|\nabla_{\mathbf{w}} g(t, \mathbf{w})\|_{\mathcal{H}} = \left\| \frac{1}{np} \sum_{i=1}^n \mathbb{1}_{t \geq (1 - y_i \mathbf{k}_i^\top \mathbf{w})^+ - y_i \mathbf{k}_i} \right\|_{\mathcal{H}} \quad (72)$$

$$\leq \frac{1}{np} \left(\sum_{i \in S} y_i^2 \right)^{1/2} \left\| \sum_{i \in S} \mathbf{k}_i \right\|_{\mathcal{H}} \quad (73)$$

$$= \frac{1}{np} \|\mathbf{y}\| \left\| \sum_{i \in S} \mathbf{k}_i \right\|_{\mathcal{H}} \leq \frac{1}{np} \|\mathbf{y}\| \left\| \sum_{i \in S} \mathbf{k}_i \right\|_{\mathcal{H}} = \frac{\sqrt{n}}{np} \left\| \sum_{i \in S} \mathbf{k}_i \right\|_{\mathcal{H}} \quad (74)$$

\blacksquare

C Base Learner Algorithm

Algorithm 5: SUBQ-BASE-LEARNER

Input: Iterations: T ; Quantile: p ; Data Matrix: \mathbf{X} , $(n \times d)$, $n \gg d$; Labels: \mathbf{y} , $(n \times 1)$; Learning schedule: $\alpha_1, \dots, \alpha_T$; Ridge parameter: λ

Output: Trained Parameters: $\mathbf{w}_{(T)}$; Base Learner: \mathcal{L}

```

1:  $\mathbf{w}_{(0)} \leftarrow \mathcal{L}(\mathbf{X}, \mathbf{y})$  ▷ Base Learner
2: for  $k \in 1, 2, \dots, T$  do
3:    $\mathbf{S}_{(k)} \leftarrow \text{SUBQUANTILE}(\mathbf{w}^{(k)}, \mathbf{X})$  ▷ Algorithm ?? 2
4:    $\mathbf{w}^{(k+1)} \leftarrow \mathcal{L}(\mathbf{S}_{(k)}, \mathbf{y}_S)$  ▷  $\mathbf{w}$ -update by base learner
5: end
6: return  $\mathbf{w}_{(T)}$ 

```

Here we can note the similarity of Algorithm ?? 5 to the algorithm described in [1]. This is because the Trimmed Maximum Likelihood Estimator is equivalent to minimizing over the subquantile of the likelihood.

Remark 27. Define the function $\Psi(t) \triangleq \min_{\mathbf{w}} g(t, \mathbf{w})$

D Experimental Details

Our datasets are synthetic and are sourced from [6]

Dataset	Dimension d	Sample Size n	Source
Polynomial	3	1000	Ours
Boston Housing	13	506	[6]
Concrete Data	8	1030	[6]
Wine Quality	11	1599	[6]

Table 4: 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 $a_i \sim \mathcal{N}(0, 1)$. Oblivious Noise is sampled from $\mathcal{N}(0, 5)$. Subquantile is capped at 10,000 iterations.

E Detailed Related Works

In this section we will give a detailed analysis of the relevant works.

E.1 High-dimensional Robust Mean Estimation via Gradient Descent [3]

In this work, Cheng et al. study high dimensional mean estimation when there exists an ϵ -fraction of adversarially corrupted data. They form a non-convex optimization problem based on a lemma from a previous paper of theirs minimize the objective with gradient descent. Let F be the objective function. First they define stationary points. Let $u \in \arg \max f(w)$, then a stationary point is defined as

$$(\nabla_w F(w, u))^\top (\tilde{w} - w) \geq 0 \quad \forall \tilde{w} \in K \quad (75)$$

where K is a closed convex set. They show that any stationary point is a good point, i.e. $\|\mu_w - \mu^*\| = \mathcal{O}(\epsilon \sqrt{\log(1/\epsilon)})$. Next, they show any approximate stationary point is a good point, i.e. if $\|\nabla f_\beta(w)\| = \mathcal{O}(\log(1/\epsilon))$, then $\|\mu_w - \mu^*\| = \mathcal{O}(\epsilon \sqrt{\log(1/\epsilon)})$. Next, they show gradient descent converges to an approximate stationary point in a polynomial number of iterations.

Technical Results:

1. F is L -lipschitz, and β -smooth
2. To prove all stationary points are good, they prove by contradiction by showing if $\|\mu_w - \mu^*\| > \mathcal{O}(\epsilon \sqrt{\log(1/\epsilon)})$, then there exists a corrupted point with a high gradient and a good point with a low gradient.
3. Let $f(w) \triangleq \max_u F(u, w)$ and $f_\beta(w) \triangleq \min_{\tilde{w}} f(\tilde{w}) + \beta \|w - \tilde{w}\|_2^2$ be the Moreau envelope. They then prove $\|\nabla f_\beta(w)\| = \mathcal{O}(\log(1/\epsilon))$.
4. Then prove $\|\nabla f_\beta(w)\| = \mathcal{O}(\log(1/\epsilon))$ in a polynomial number of iterations w.r.t to n the sample size, and d the sample dimension.

E.2 Trimmed Maximum Likelihood Estimation for Robust Generalized Linear Model [1]

First we will give the algorithm

$$S^{(t)} = \arg \min_{T \subset S^{(0)}: |T|=(1-2\epsilon)n} \sum_{i \in T} -\log f(y_i | \langle \beta^{(t)}, \mathbf{x}_i \rangle) \quad (76)$$

$$\beta^{(t+1)} = \arg \min_{\beta, \|\beta\| \leq R} \sum_{i \in S^{(t)}} -\log f(y_i | \langle \beta^{(t)}, \mathbf{x}_i \rangle) \quad (77)$$

In Equation (76), the algorithm chooses the $(1 - 2\epsilon)n$ points giving the least error and put this in the set $S^{(t)}$. Next, in Equation (77), the algorithm then finds β that minimizes the negative log likelihood error for all the points in $S^{(t)}$ s.t. $\|\beta\| \leq R$. For the theoretical analysis, Awasthi et al. consider a different approximation stationary point from [3].

$$\frac{1}{n} \sum_{i \in S} \nabla_\beta \log f(y_i | \langle \beta, \mathbf{x}_i \rangle)^\top \frac{(\beta^* - \beta)}{\|\beta^* - \beta\|} \leq \gamma \quad (78)$$

We see Equation (78) is an upper bound, instead of a lower bound, of Equation (75). Next, they prove their algorithm reaches a η stationary point. Their proof does not use Moreau Envelopes or ideas in concave-non-convex optimization, rather they use the fact their algorithm terminates after it reaches a point when it can no longer make η improvement.