

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

Arvind Rathnashyam* Alex Gittens†

September 16, 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 [13, 25]. 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. In our theoretical analysis, we analyze our non-convex concave objective function with the Moreau Envelope. We show (i) a stationary point with respect to the Moreau Envelope is a good point and (ii) we can reach a stationary point with gradient descent methods. Further, we analyze accelerated gradient methods for the non-convex concave minimax optimization problem. We empirically test Kernel Ridge Regression and Kernel Classification on various state of the art datasets and show Subquantile Minimization gives strong results. Furthermore, we run experiments on various datasets and compare with the state-of-the-art algorithms to show the superior performance of Subquantile Minimization.

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

†CS, Rensselaer Polytechnic Institute, gittes@rpi.edu

Contents

1	Introduction	3
1.1	Related Work	3
1.2	Notation	4
2	Subquantile Minimization	4
3	Structural Results	5
3.1	Kernel Regression	6
3.2	Kernel Binary Classification	7
3.3	Kernel Multi-Class Classification	7
3.4	Necessary Kernel Inequalities	8
4	Optimization Results	9
4.1	Accelerated Gradient Methods	9
4.1.1	Momentum Accelerated Gradient Descent	10
4.1.2	Nesterov Accelerated Gradient Descent	10
5	Experiments	10
5.1	Kernel Regression	10
5.2	Kernel Binary Classification	10
5.3	Kernel Multi-Class Classification	12
5.4	Accelerated Gradient Methods	13
6	Discussion	14
A	Kernel Embedding Inequalities	17
A.1	Proof of Lemma 17	17
A.2	Proof of Lemma 18	17
B	Proofs for Section 3	18
B.1	Proof of Lemma 10	18
B.2	Proof of Lemma 12	18
B.3	Proof of Lemma 15	19
B.4	Proof of Lemma 20	19
B.5	Proof of Lemma 21	20
B.6	Proof of Theorem 22	21
C	Proofs for Section 4	22
C.1	Proof of Theorem 28	22
D	Auxiliary Lemmas	23
E	Base Learner Algorithm	24
F	Experimental Details	25
F.1	Kernel Regression	25
F.2	Kernel Binary Classification	25
G	Detailed Related Works	26
G.1	High-dimensional Robust Mean Estimation via Gradient Descent [3]	26
G.2	Trimmed Maximum Likelihood Estimation for Robust Generalized Linear Model [1]	26

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, 14, 19, 7] as well as for linear regression [2, 17]. 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 [12, 10]. 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.

First, we will give our main error bounds.

Theorem 2. (Informal). *Let the dataset be given as $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ such that ϵn are arbitrarily corrupted by an adversary. Let \mathbf{K} be the kernel matrix and S be the points with the lowest error w.r.t $\hat{\mathbf{w}}$, then Subquantile Minimization returns $\hat{\mathbf{w}}$ such that for Kernel Ridge Regression:*

$$\|\hat{\mathbf{w}} - \mathbf{w}^*\|_{\mathcal{H}} \leq \left(\frac{2L\sigma_{\min}(\sum_{i \in S \cap P} \mathbf{k}_i \mathbf{k}_i^{\top}) - \sqrt{\sigma_{\max}(\mathbf{K})}}{\sqrt{\sigma_{\max}(\mathbf{K})}} \right)^{1/2} + L\epsilon \quad (1)$$

Linear Regression:

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

Kernel Binary Classification:

$$\|\hat{\mathbf{w}} - \mathbf{w}^*\|_{\mathcal{H}} \leq \quad (3)$$

Kernel Multi-Class Classification:

$$\|\hat{\mathbf{w}} - \mathbf{w}^*\|_{\mathcal{H}} \leq \quad (4)$$

We will now state our main contributions clearly. **Contributions**

1. We propose a gradient-descent based algorithm for robust kernel learning in the Huber ϵ -Contamination Model which is fast.
2. We rigorously analyze error bounds for subquantile minimization in the linear regression, kernel regression, kernel binary classification, and kernel multi-class classification tasks.
3. We give new bounds for accelerated gradient methods for accelerated gradient methods in nonconvex-concave minimax optimization.

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.

[14] 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

[11] 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.

[29] 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 .

$$\mathcal{Q}_p(U) \leq t \iff \mathbb{P}\{U \leq t\} \geq p \quad (5)$$

The p -subquantile of the risk is then given by

$$\mathbb{L}_p(U) = \frac{1}{p} \int_0^p \mathcal{Q}_q(U) dq = \mathbb{E}[U | U \leq \mathcal{Q}_p(U)] = \max_{t \in \mathbb{R}} \left\{ t - \frac{1}{p} \mathbb{E}(t - U)^+ \right\} \quad (6)$$

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

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d: \|\mathbf{w}\| \leq R} \max_{t \in \mathbb{R}} \left\{ g(t, \mathbf{w}) \triangleq \sum_{i=1}^n (t - (f(\mathbf{x}_i; \mathbf{w}) - y_i)^2)^+ \right\} \quad (7)$$

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 [22]. 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

$$t^{(k+1)} = \arg \max_{t \in \mathbb{R}} g(\mathbf{w}^{(k)}, t) \quad (8)$$

$$\mathbf{w}^{(k+1)} = \Pi_{\mathcal{K}} \left(\mathbf{w}^{(k)} - \alpha \nabla g(\mathbf{w}^{(k)}, t^{(k+1)}) \right) \quad (9)$$

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 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:  $\mathbf{S} \leftarrow (\mathbf{x}_1^\top \dots \mathbf{x}_{np}^\top)^\top$ 
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 (8) and (9).

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 $\hat{\nu}_i \triangleq f(\mathbf{x}_i; \mathbf{w}, y_i)$, 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 .

$$\partial(-g(t, \mathbf{w})) = \partial \left(-t + \frac{1}{np} \sum_{i=1}^n (t - \hat{\nu}_i) \right) \quad (10)$$

$$= -1 + \frac{1}{np} \sum_{i=1}^{np} \begin{cases} 1 & \text{if } t > \hat{\nu}_i \\ 0 & \text{if } t < \hat{\nu}_i \\ [0, 1] & \text{if } t = \hat{\nu}_i \end{cases} \quad (11)$$

$$= 0 \text{ when } t = \hat{\nu}_{np} \quad (12)$$

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

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

Lemma 5. *Let $\hat{\nu}_i \triangleq f(\mathbf{x}_i; \mathbf{w}, y_i)$ 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 $\nabla_{\mathbf{w}} g(t^{(k)}, \mathbf{w}^{(k)}) = \frac{1}{np} \sum_{i=1}^{np} \nabla f(\mathbf{x}_i; \mathbf{w}^{(k)}, y_i)$*

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

$$\nabla_{\mathbf{w}} g(t^{(k+1)}, \mathbf{w}^{(k)}) = \nabla_{\mathbf{w}} \left(\hat{\nu}_{np} - \frac{1}{np} \sum_{i=1}^n (\hat{\nu}_{np} - f(\mathbf{x}_i; \mathbf{w}, y_i))^+ \right) \quad (13)$$

$$= -\frac{1}{np} \sum_{i=1}^{np} \nabla_{\mathbf{w}} (\hat{\nu}_{np} - f(\mathbf{x}_i; \mathbf{w}, y_i))^+ \quad (14)$$

$$= \frac{1}{np} \sum_{i=1}^n \nabla_{\mathbf{w}} f(\mathbf{x}_i; \mathbf{w}^{(k)}, y_i) \begin{cases} 1 & \text{if } t > \hat{\nu}_i \\ 0 & \text{if } t < \hat{\nu}_i \\ [0, 1] & \text{if } t = \hat{\nu}_i \end{cases} \quad (15)$$

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

$$\nabla_{\mathbf{w}} g(t^{(k+1)}, \mathbf{w}^{(k)}) = \frac{1}{np} \sum_{i=1}^{np} \nabla_{\mathbf{w}} f(\mathbf{x}_i; \mathbf{w}, y_i) \quad (16)$$

This concludes the proof.

■ 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: $f(\mathbf{x}^*; \mathbf{w}) = \sum_{i=1}^n \mathbf{w}_i \kappa(\mathbf{x}_i, \mathbf{x}^*)$. From our definition of $S^{(k)}$ in ??, 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 6. (Moreau Envelope, [16]). Let f be proper lower semi-continuous convex function $f : \mathcal{X} \rightarrow \mathbb{R}$, then the Moreau Envelope is defined as:

$$f_\lambda(\mathbf{x}) \triangleq \inf_{\hat{\mathbf{x}} \in \mathcal{X}} \left(f(\hat{\mathbf{x}}) + \frac{1}{2\lambda} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \right) \quad (17)$$

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

Assumption 7. Define $\Phi(\cdot)$ as the function in remark 14. Then it follows $\arg \min_{\mathbf{w} \in \mathbb{R}^d} \Phi(\mathbf{w}) = \mathbf{w}^*$

Definition 8. (First-Order Stationary Point). Let $\Phi(\mathbf{w}) = \max_t g(t, \mathbf{w})$. Then \mathbf{w} is a first-order stationary point if

$$\nabla_{\mathbf{w}} \Phi(\mathbf{w})^\top (\tilde{\mathbf{w}} - \mathbf{w}) \geq 0 \quad \forall \tilde{\mathbf{w}} \in \mathcal{K} \quad (18)$$

Lemma 9. If $\|\mathbf{w} - \mathbf{w}^*\|_{\mathcal{H}} \geq \eta$, then

$$\nabla_{\mathbf{w}} \Phi(\mathbf{w})^\top (\mathbf{w}^* - \mathbf{w}) \leq f(\eta) \quad (19)$$

Proof.

$$\nabla_{\mathbf{w}} \Phi(\mathbf{w})^\top (\mathbf{w}^* - \mathbf{w}) = \left(\sum_{i \in S} \mathbf{k}_i (\mathbf{w}^\top \mathbf{k}_i - y_i) \right)^\top (\mathbf{w}^* - \mathbf{w}) \quad (20)$$

$$= \left(\sum_{i \in S \cap P} \mathbf{k}_i (\mathbf{w}^\top \mathbf{k}_i - y_i) + \sum_{i \in S \cap Q} \mathbf{k}_i (\mathbf{w}^\top \mathbf{k}_i - y_i) \right)^\top (\mathbf{w}^* - \mathbf{w}) \quad (21)$$

$$= \left(\sum_{i \in S \cap P} \mathbf{k}_i ((\mathbf{w} - \mathbf{w}^*)^\top \mathbf{k}_i - \eta_i) + \sum_{i \in S \cap Q} \mathbf{k}_i (\mathbf{w}^\top \mathbf{k}_i - y_i) \right)^\top (\mathbf{w}^* - \mathbf{w}) \quad (22)$$

■

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

$$f(\mathbf{x}, y_i; \mathbf{w}) = (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 \quad (23)$$

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

Lemma 10. (*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:

$$|g(t, \mathbf{w}) - g(t, \hat{\mathbf{w}})| \leq L \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}} \quad (24)$$

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 11. (*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:

$$\|\nabla_{\mathbf{w}} g(t, \mathbf{w}) - \nabla_{\mathbf{w}} g(t, \hat{\mathbf{w}})\| \leq \beta \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}} \quad (25)$$

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.

$$\|\nabla_{\mathbf{w}}^2 g(t, \mathbf{w})\|_{\mathcal{H}} = \left\| \frac{2}{np} \sum_{i \in S} \mathbf{k}_i \mathbf{k}_i^{\top} \right\|_{\mathcal{H}} \leq \frac{2}{np} \left\| \sum_{i \in S} \mathbf{k}_i \right\|_{\mathcal{H}}^2 \stackrel{\text{lem. 17}}{\leq} \quad (26)$$

This concludes the proof. ■

3.2 Kernel Binary 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)

$$f(\mathbf{x}_i, y_i; \mathbf{w}) = -(y_i \log(\sigma(\mathbf{w}, \mathbf{k}_i)) + (1 - y_i) \log(1 - \sigma(\mathbf{w}, \mathbf{k}_i))) \quad (27)$$

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

Lemma 12. (*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:

$$|g(t, \mathbf{w}) - g(t, \hat{\mathbf{w}})| \leq L \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}} \quad (28)$$

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 13. (*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:

$$\|\nabla_{\mathbf{w}} g(t, \mathbf{w}) - \nabla_{\mathbf{w}} g(t, \hat{\mathbf{w}})\| \leq \beta \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}} \quad (29)$$

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

Remark 14. 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 Kernel Multi-Class Classification

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

$$f(\mathbf{x}_i, y_i; \mathbf{W}) = - \sum_{j=1}^{|\mathcal{C}|} \mathbb{1}_{y_i=j} \log \left(\frac{\exp(\mathbf{W}_k^{\top} \mathbf{k}_i)}{\sum_{h=1}^{|\mathcal{C}|} \exp(\mathbf{W}_h^{\top} \mathbf{k}_i)} \right) \quad (30)$$

Lemma 15. (*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:

$$|g(t, \mathbf{w}) - g(t, \hat{\mathbf{w}})| \leq L \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}} \quad (31)$$

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

3.4 Necessary Kernel Inequalities

We will first extend the idea of Resilience [28] to kernel learning.

Definition 16. (Resilience) from [28]. 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 16, we will be unable to put error bounds on our algorithm.

Lemma 17. *Let S be the set of elements in the subquantile, then*

$$\left\| \frac{1}{np} \sum_{i \in X} \mathbf{k}_i \right\|_{\mathcal{H}} \leq \mathcal{O}() \quad (32)$$

Lemma 18. *Under the same setting as lemma 17,*

$$\left\| \frac{1}{np} \sum_{i \in S} \mathbf{k}_i \right\|_{\mathcal{H}} \leq \mathcal{O}() \quad (33)$$

Definition 19. (First-Order Stationary Point). Let $\Phi(\mathbf{w}) = \max_t g(t, \mathbf{w})$. Then \mathbf{w} is a first-order stationary point if

$$\|\nabla \Phi_{\lambda}(\mathbf{w})\|_{\mathcal{H}} = 0 \quad (34)$$

i.e.

$$\mathbf{w} = \arg \min_{\hat{\mathbf{w}} \in \mathcal{K}} \left(\Phi(\hat{\mathbf{w}}) + \frac{1}{2\lambda} \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}} \right) \quad (35)$$

Lemma 20. *If $\|\mathbf{w} - \mathbf{w}^*\| \leq \eta$, then it follows*

$$\Phi(\mathbf{w}) - \Phi(\mathbf{w}^*) \leq \eta^2 \frac{\sigma_{\max}(\sum_{i \in P} \mathbf{k}_i \mathbf{k}_i^{\top})}{\sqrt{\sigma_{\min}(\mathbf{K})}} - 2\eta \frac{\|\sum_{i \in P} \eta_i \mathbf{k}_i\|}{\sqrt{\sigma_{\max}(\mathbf{K})}} \quad (36)$$

Lemma 21. *If $\|\mathbf{w} - \mathbf{w}^*\| \geq \eta$, then it follows*

$$\Phi(\mathbf{w}) - \Phi(\mathbf{w}^*) \geq \frac{\eta^2 \sigma_{\min} \left(\left(\sum_{i \in S \cap P} \mathbf{k}_i \right) \left(\sum_{i \in S \cap P} \mathbf{k}_i \right)^{\top} \right)}{\sqrt{\sigma_{\max}(\mathbf{K})}} - 2\eta \frac{\|\sum_{i \in S \cap P} \eta_i \mathbf{k}_i\|}{\sqrt{\sigma_{\min}(\mathbf{K})}} - \sum_{j \in P \setminus S} \eta_j^2 \quad (37)$$

The proof is deferred to § Appendix B.4.

Theorem 22. *Let $\hat{\mathbf{w}}$ be a stationary point defined in definition 8 for the function Φ defined in ???. Then,*

$$\eta \leq \left(\sum_{j \in P \setminus S} \eta_j^2 \right)^{1/2} \left(\frac{2\sigma_{\min}(\sum_{i \in S \cap P} \mathbf{k}_i \mathbf{k}_i^{\top}) L - \sqrt{\sigma_{\max}(\mathbf{K})}}{2L\sqrt{\sigma_{\max}(\mathbf{K})}} \right)^{1/2} \quad (38)$$

where L is the Lipschitz Constant given in Lemma 10.

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 23. ([24, 23]). 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:

$$\|\hat{\mathbf{w}} - \mathbf{w}\| = \lambda\epsilon \text{ and } \min_{\mathbf{g} \in \partial\Phi(\hat{\mathbf{w}}) + \partial\mathcal{I}_{\mathcal{K}}(\hat{\mathbf{w}})} \|\mathbf{g}\| \leq \epsilon \quad (39)$$

How to extend this to Hilbert Space Norm?

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

$$\partial\mathcal{I}_{\mathcal{K}} = \mathcal{N}(\hat{\mathbf{w}}) = \{\tilde{\mathbf{w}} \in \mathbb{R}^n | \langle \tilde{\mathbf{w}}, \bar{\mathbf{w}} - \hat{\mathbf{w}} \rangle \leq \forall \bar{\mathbf{w}} \in \mathcal{K}\} \quad (40)$$

We will define a convex cone.

Definition 24. 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 25. Let $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}'} \left(\Phi(\mathbf{w}') + \frac{1}{2\lambda} \|\mathbf{w} - \mathbf{w}'\|^2 \right)$ s.t. $\|\nabla \Phi_{\lambda}(\mathbf{w})\|_{\mathcal{H}} \leq \epsilon$, then it follows

$$\|\hat{\mathbf{w}} - \mathbf{w}^*\|_{\mathcal{H}} \leq \Xi \quad (41)$$

Definition 26. (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 6. 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 26 is used in the analysis of the minimax algorithm in [3]. First, if you can prove a stationary point is good, theorem 22, then using lemma 23, 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 algorithm 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 $(\sum_{\mathbf{x} \in S^{(T)}} \mathbf{k}_{\mathbf{x}} (\mathbf{k}_{\mathbf{x}}^{\top} \mathbf{w}^{(T)} - y))^{\top} \left(\frac{\mathbf{w}^{(T)} - \mathbf{w}^*}{\|\mathbf{w}^{(T)} - \mathbf{w}^*\|} \right)$. First, we will show using stepsize of $1/\beta$ returns a μ -approximate stationary point.

Theorem 27. (Algorithm 1 reaches a η -approximate stationary point). Algorithm 1 reaches a η -approximate stationary point in a polynomial number of iterations.

Proof. From [15] Theorem 31 and [3] Lemma 4.2, it follows:

$$\mathbb{E} \left[\|\nabla \Phi_{1/2\ell}(\bar{\mathbf{w}})\|^2 \right] \leq 2 \cdot \frac{(\Phi_{1/2\ell}(\mathbf{w}_0) - \min \Phi(\mathbf{w})) + \ell\beta^2\gamma^2}{\gamma\sqrt{T+1}} \quad (42)$$

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 23, $\|\hat{\mathbf{w}} - \mathbf{w}^{(T)}\| = \mu/2\ell$. ■

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

In this section we study Momentum Accelerated Gradient Descent [21, 20] with our non-convex-concave optimization algorithm.

$$\mathbf{b}^{(t)} = \mu \mathbf{b}^{(t-1)} + \nabla_{\mathbf{w}} \Phi(\mathbf{w}^{(t-1)}) \quad (43)$$

$$\mathbf{w}^{(t)} = \mathbf{w}^{(t-1)} - \alpha \mathbf{b}^{(t)} \quad (44)$$

Theorem 28. *Momentum Accelerated Gradient Descent given in Equations (43) and (44) reaches a η -approximate stationary point. Algorithm 1 reaches a η -approximate stationary point in a polynomial number of iterations.*

$$\mathbb{E} \left[\|\nabla \Phi_{1/2\ell}(\mathbf{w})\|_{\mathcal{H}}^2 \right] \leq \quad (45)$$

4.1.2 Nesterov Accelerated Gradient Descent

In this section we study Nesterov Accelerated Gradient Descent [18] with our non-convex-concave optimization algorithm.

$$\mathbf{b}^{(t+1)} = (1 + \mu) \mathbf{w}^{(t)} - \mu \mathbf{w}^{(t-1)} \quad (46)$$

$$\mathbf{w}^{(t+1)} = \mathbf{b}^{(t+1)} - \alpha \nabla_{\mathbf{w}} \Phi(\mathbf{w}^{(t)}) \quad (47)$$

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:    $\mathbf{S}_{(k)} \leftarrow \text{SUBQUANTILE}(\mathbf{w}^{(k)}, \mathbf{X})$  ▷ Algorithm 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. (9)
6: end
7: return  $\mathbf{w}_{(T)}$ 

```

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 Binary Classification

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

Now we will give some experimental results.

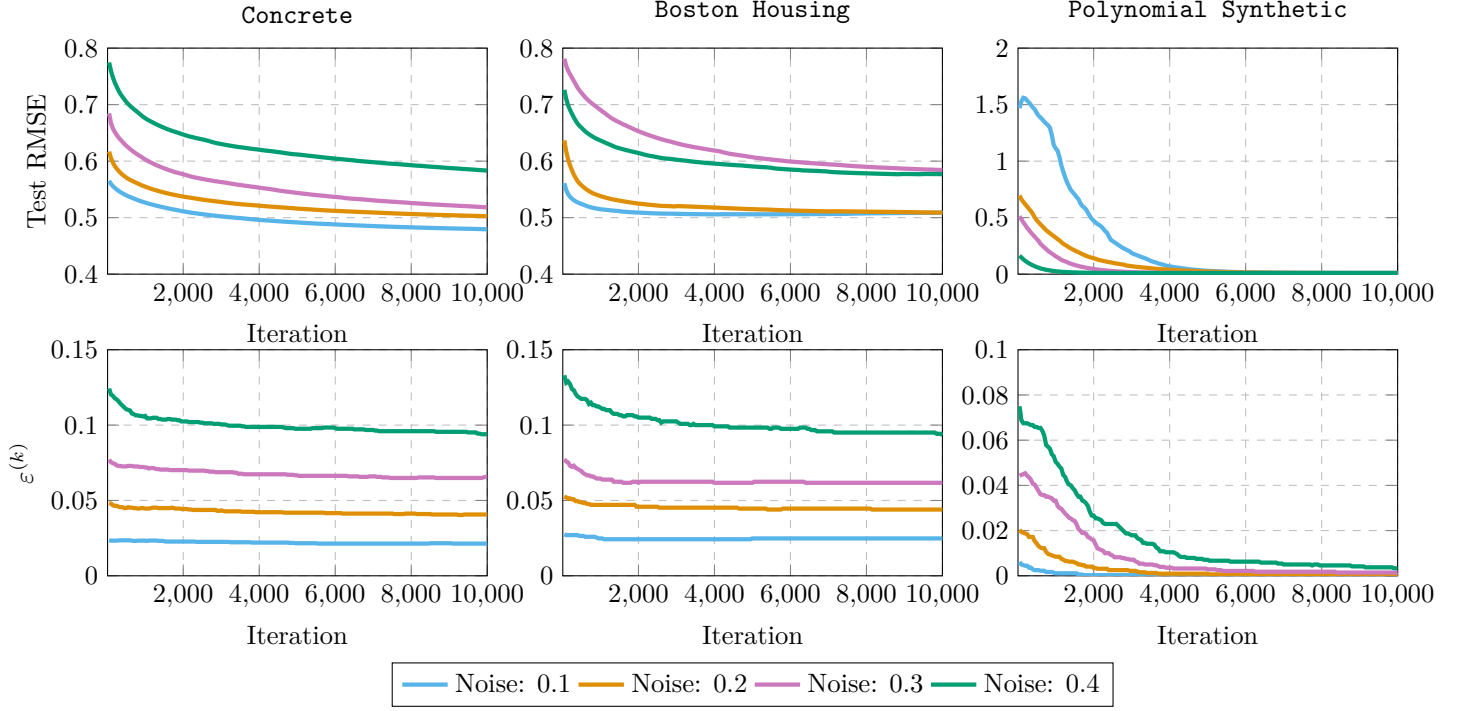


Figure 1: Test RMSE over the iterations in Concrete, Boston Housing, and Polynomial Datasets for SUBQUANTILE at different noise levels

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: $\mathbf{w}_{(0)} \leftarrow (\mathbf{K}^\top \mathbf{K} + \lambda \mathbf{I}) \mathbf{K}^\top \mathbf{y}$ \triangleright Base Learner
 - 2: **for** $k \in 1, 2, \dots, T$ **do**
 - 3: $\mathbf{S}_{(k)} \leftarrow \text{SUBQUANTILE}(\mathbf{w}^{(k)}, \mathbf{X})$ \triangleright Algorithm 2
 - 4: $\nabla_{\mathbf{w}} g(t^{(k+1)}, \mathbf{w}^{(k)}) \leftarrow -\sum_{i \in S^{(k)}} y_i \mathbf{k}_i + \lambda \mathbf{K} \mathbf{w}^{(k)}$ \triangleright Gradient Calculation
 - 5: $\mathbf{w}^{(k+1)} \leftarrow \mathbf{w}^{(k)} - \alpha_{(k)} \nabla_{\mathbf{w}} g(t^{(k+1)}, \mathbf{w}^{(k)})$ \triangleright \mathbf{w} -update in Equation (9)
 - 6: **end**
 - 7: **return** $\mathbf{w}_{(T)}$
-

Objectives	Test RMSE (Polynomial Regression (Degree = 3))			
	$\epsilon = 0.1(\downarrow)$	$\epsilon = 0.2(\downarrow)$	$\epsilon = 0.3(\downarrow)$	$\epsilon = 0.4(\downarrow)$
KRR	0.460 _(0.2143)	1.171 _(0.7809)	0.950 _(0.3053)	1.230 _(0.4678)
TERM [14]	∞	∞	∞	∞
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 $p = 0.02$.

Objectives	Test RMSE (Boston Housing Regression)			
	$\epsilon = 0.1(\downarrow)$	$\epsilon = 0.2(\downarrow)$	$\epsilon = 0.3(\downarrow)$	$\epsilon = 0.4(\downarrow)$
KRR	0.544 _(0.0712)	0.862 _(0.1199)	0.865 _(0.0811)	1.049 _(0.2249)
TERM [14]	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)			
	$\epsilon = 0.1(\downarrow)$	$\epsilon = 0.2(\downarrow)$	$\epsilon = 0.3(\downarrow)$	$\epsilon = 0.4(\downarrow)$
KRR	0.802 _(0.0324)	0.929 _(0.0209)	0.993 _(0.0441)	0.775 _(0.0514)
TERM [14]	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 $\lambda = 2$.

5.3 Kernel Multi-Class Classification

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

Objectives	Test Accuracy (Heart Disease)			
	$\epsilon = 0.1(\uparrow)$	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.3(\uparrow)$	$\epsilon = 0.4(\uparrow)$
SVM	0.866 _(0.0380)	0.826 _(0.0482)	0.705 _(0.0647)	0.625 _(0.0475)
TERM [14]	0.833 _(0.0217)	0.790 _(0.0420)	0.751 _(0.0457)	0.610 _(0.0667)
SEVER [4]	0.803 _(0.0874)	0.728 _(0.1134)	0.702 _(0.1049)	0.531 _(0.0321)
SUBQUANTILE($p = 1 - \epsilon$)	0.869 _(0.0549)	0.846 _(0.0382)	0.734 _(0.0792)	0.597 _(0.0794)
Genie ERM	∞	∞	∞	∞

Table 4: **Heart Disease** Dataset. Oblivious Noise flips the label. The Radial Basis Function is used in all experiments.

Objectives	Test Accuracy (Breast Cancer)			
	$\epsilon = 0.1(\uparrow)$	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.3(\uparrow)$	$\epsilon = 0.4(\uparrow)$
SVM	0.963 _(0.0035)	0.940 _(0.0140)	0.912 _(0.0590)	0.781 _(0.0496)
TERM [14]	0.970 _(0.0119)	0.937 _(0.0151)	0.914 _(0.0238)	0.747 _(0.0650)
SEVER [4]	0.947 _(0.0166)	0.918 _(0.0439)	0.825 _(0.0769)	0.725 _(0.0446)
SUBQUANTILE($p = 1 - \epsilon$)	0.972 _(0.0151)	0.956 _(0.0055)	0.958 _(0.0129)	0.681 _(0.0556)
Genie ERM	∞	∞	∞	∞

Table 5: **Breast Cancer** Dataset. Oblivious Noise flips the label. The Radial Basis Function is used in all experiments.

5.4 Accelerated Gradient Methods

In this section we give some empirical results on using different accelerated gradient methods described in § Section 4. In Figure 2, we see using momentum can give significantly faster convergence.

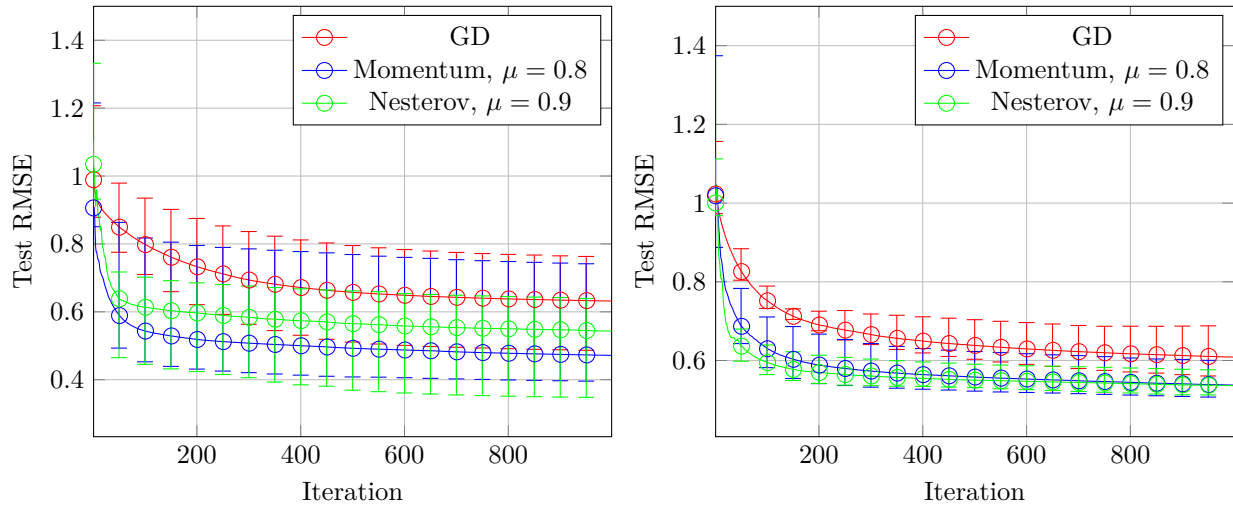


Figure 2: The effect of Momentum for Kernel Regression in the **Boston Housing** dataset (*Left*) and **Concrete** dataset (*Right*). We observe faster convergence. In both experiments, we use the Radial Basis Function.

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] Steinbrunn William Pfisterer Matthias Janosi, Andras and Robert Detrano. Heart Disease. UCI Machine Learning Repository, 1988. DOI: <https://doi.org/10.24432/C52P4X>.
- [10] 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.
- [11] 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.
- [12] Ashish Khetan, Zachary C. Lipton, and Anima Anandkumar. Learning from noisy singly-labeled data. In *International Conference on Learning Representations*, 2018.
- [13] 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.
- [14] Tian Li, Ahmad Beirami, Maziar Sanjabi, and Virginia Smith. Tilted empirical risk minimization. In *International Conference on Learning Representations*, 2021.
- [15] 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.
- [16] Jean-Jacques Moreau. Proximité et dualité dans un espace hilbertien. *Bulletin de la Société mathématique de France*, 93:273–299, 1965.

- [17] 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.
- [18] Yurii Evgen’evich Nesterov. A method of solving a convex programming problem with convergence rate $\mathcal{O}(\frac{1}{k^2})$. In *Doklady Akademii Nauk*, volume 269, pages 543–547. Russian Academy of Sciences, 1983.
- [19] 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.
- [20] Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *Ussr computational mathematics and mathematical physics*, 4(5):1–17, 1964.
- [21] Ning Qian. On the momentum term in gradient descent learning algorithms. *Neural networks*, 12(1):145–151, 1999.
- [22] 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.
- [23] R Tyrrell Rockafellar. *Convex analysis*. 2015.
- [24] Ralph Tyrrell Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, 1970.
- [25] 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.
- [26] Markus Schneider. Probability inequalities for kernel embeddings in sampling without replacement. In *International Conference on Artificial Intelligence and Statistics*, 2016.
- [27] Jonathan Richard Shewchuk et al. An introduction to the conjugate gradient method without the agonizing pain, 1994.
- [28] 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.
- [29] 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.

A Kernel Embedding Inequalities

Lemma 29. (*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 16 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 [26], Theorem 2

Lemma 30. 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:

$$\left\| \sum_{i=1}^n \xi_i \right\| \leq C \text{ with high probability} \quad (48)$$

First, we can note,

$$\sum_{i=1}^n \xi_i \sim \mathcal{N}(n\mu, n\sigma^2) \quad (49)$$

With Hoeffding's Inequality, we have:

$$\mathbb{P} \left\{ \left\| \frac{1}{n} \sum_{i=1}^n \xi_i - \mu \right\| \geq \epsilon \right\} \leq 2 \exp \left(-\frac{2n^2\epsilon^2}{\sum_{i=1}^n \|\xi_i\|_{\psi_2}^2} \right) \quad (50)$$

A.1 Proof of Lemma 17

Proof.

$$\left\| \frac{1}{np} \sum_{i \in S \cap P} \mathbf{k}_i \right\| = \frac{1}{np} \left\| \sum_{i \in S \cap P} \sum_{j \in X} \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) \right\| \quad (51)$$

■

A.2 Proof of Lemma 18

Proof.

■

B Proofs for Section 3

In this section we give some deferred proofs.

B.1 Proof of Lemma 10

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

$$\|\nabla_{\mathbf{w}} g(t, \mathbf{w})\|_{\mathcal{H}} = \left\| \frac{2}{np} \sum_{i=1}^n \mathbb{1}_{t \geq (\mathbf{k}_i^\top \mathbf{w} - y_i)^2} (\mathbf{k}_i (\mathbf{k}_i^\top \mathbf{w} - y_i)) \right\|_{\mathcal{H}} \quad (52)$$

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 $t \geq (\mathbf{k}_i^\top \mathbf{w} - y_i)^2$.

$$= \left\| \frac{2}{np} \sum_{i=1}^m \mathbf{k}_i (\mathbf{k}_i^\top \mathbf{w} - y_i) \right\|_{\mathcal{H}} \quad (53)$$

$$\stackrel{(a)}{\leq} \frac{2}{np} \left(\left\| \sum_{i=1}^m \mathbf{k}_i (\mathbf{k}_i^\top \mathbf{w}) \right\|_{\mathcal{H}} + \left\| \sum_{i=1}^m \mathbf{k}_i y_i \right\|_{\mathcal{H}} \right) \quad (54)$$

$$\stackrel{(b)}{\leq} \frac{2}{np} \left(\left\| \sum_{i=1}^m \mathbf{k}_i \right\|_{\mathcal{H}} \left\| \sum_{i=1}^m \mathbf{k}_i^\top \mathbf{w} \right\|_{\mathcal{H}} + \left\| \sum_{i=1}^m \mathbf{k}_i (\mathbf{k}_i^\top \mathbf{w}^* + \eta_i) \right\|_{\mathcal{H}} \right) \quad (55)$$

$$\stackrel{(c)}{\leq} \frac{2}{np} \left(\left\| \sum_{i=1}^m \mathbf{k}_i \right\|_{\mathcal{H}}^2 \|\mathbf{w}\|_{\mathcal{H}} + \left\| \sum_{i=1}^m \mathbf{k}_i \right\|_{\mathcal{H}} \left\| \sum_{i=1}^m \mathbf{k}_i^\top \mathbf{w}^* \right\|_{\mathcal{H}} + \left\| \sum_{i=1}^m \mathbf{k}_i \eta_i \right\|_{\mathcal{H}} \right) \quad (56)$$

$$\stackrel{(d)}{\leq} \frac{4R}{np} \left\| \sum_{i=1}^{np} \mathbf{k}_i \right\|_{\mathcal{H}}^2 + \frac{2}{np} \left\| \sum_{i=1}^{np} \mathbf{k}_i \eta_i \right\|_{\mathcal{H}} \quad (57)$$

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

B.2 Proof of Lemma 12

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

$$\|\nabla_{\mathbf{w}} g(t, \mathbf{w})\|_{\mathcal{H}} = \left\| \frac{1}{np} \sum_{i=1}^n \mathbb{1}_{t \geq (1-y_i) \log(\mathbf{w}^\top \mathbf{k}_i)} (y_i - \sigma(\mathbf{w}^\top \mathbf{k}_i)) \mathbf{k}_i \right\|_{\mathcal{H}} \quad (58)$$

$$\leq \frac{1}{np} \left\| \sum_{i \in S} y_i \mathbf{k}_i \right\|_{\mathcal{H}} + \frac{1}{np} \left\| \sum_{i \in S} \sigma(\mathbf{w}^\top \mathbf{k}_i) \mathbf{k}_i \right\|_{\mathcal{H}} \quad (59)$$

$$\stackrel{(a)}{\leq} \frac{1}{np} \left(\sum_{i \in S} y_i^2 \right)^{1/2} \left(\sum_{i \in S} \|\mathbf{k}_i\|_{\mathcal{H}}^2 \right)^{1/2} + \frac{1}{np} \left(\sum_{i \in S} \sigma(\mathbf{w}^\top \mathbf{k}_i)^2 \right)^{1/2} \left(\sum_{i \in S} \|\mathbf{k}_i\|_{\mathcal{H}}^2 \right)^{1/2} \quad (60)$$

$$\leq \frac{2\sqrt{np}}{np} \left(\sum_{i=1}^n \|\mathbf{k}_i\|_{\mathcal{H}} \right) \quad (61)$$

(a) follows from the fact that $y_i \in \{0, 1\}$ and $\text{range}(\sigma) \in [0, 1]$. This completes the proof. \blacksquare

B.3 Proof of Lemma 15

Proof. We use the spectral norm of the gradient to bound L from above. Let S be denoted as the subquantile set.

$$\|\nabla \mathbf{w}g(t, \mathbf{W})\|_2 = \left\| \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^{|\mathcal{C}|} \mathbf{k}_i \left(\mathbb{1}_{y_i=j} - \frac{\exp(\mathbf{W}_j^\top \mathbf{k}_i)}{\sum_{h=1}^{|\mathcal{C}|} \exp(\mathbf{W}_h^\top \mathbf{k}_i)} \right) \right\| \quad (62)$$

$$\leq \frac{1}{np} \left\| \sum_{i=1}^n \mathbf{k}_i \right\| + \frac{1}{np} \left\| \sum_{i=1}^n \sum_{j=1}^{|\mathcal{C}|} \mathbf{k}_i \left(\frac{\exp(\mathbf{W}_j^\top \mathbf{k}_i)}{\sum_{h=1}^{|\mathcal{C}|} \exp(\mathbf{W}_h^\top \mathbf{k}_i)} \right) \right\| \quad (63)$$

$$\leq \frac{2}{np} \sum_{i=1}^n \|\mathbf{k}_i\| \quad (64)$$

■

B.4 Proof of Lemma 20

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

$$\Phi(\mathbf{w}) - \Phi(\mathbf{w}^*) = \sum_{i \in S} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 - \sum_{j \in P} (\mathbf{w}^{*\top} \mathbf{k}_j - y_j)^2 \quad (65)$$

$$= \sum_{i \in S \cap P} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 + \sum_{i \in S \cap Q} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 - \sum_{j \in P} (\mathbf{w}^{*\top} \mathbf{k}_j - y_j)^2 \quad (66)$$

$$\stackrel{(a)}{\leq} \sum_{i \in S \cap P} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 + \sum_{i \in P \setminus S} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 - \sum_{j \in P} (\mathbf{w}^{*\top} \mathbf{k}_j - y_j)^2 \quad (67)$$

$$= \sum_{i \in P} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 - \sum_{j \in P} (\mathbf{w}^{*\top} \mathbf{k}_j - y_j)^2 \quad (68)$$

$$= \sum_{i \in P} (\mathbf{w}^\top \mathbf{k}_i - \mathbf{w}^{*\top} \mathbf{k}_i - \eta_i)^2 - \sum_{j \in P} \eta_j^2 \quad (69)$$

$$= \sum_{i \in P} \left((\mathbf{w} - \mathbf{w}^*)^\top \mathbf{k}_i - \eta_i \right)^2 - \sum_{j \in P} \eta_j^2 \quad (70)$$

$$= \sum_{i \in P} \left\| (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{k}_i \right\|^2 - 2\eta_i (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{k}_i + \eta_i^2 - \sum_{j \in P} \eta_j^2 \quad (71)$$

$$\stackrel{\text{lem. 31}}{\leq} \frac{\|\mathbf{w} - \mathbf{w}^*\|_{\mathcal{H}}^2}{\sqrt{\sigma_{\min}(\mathbf{K})}} \sigma_{\max} \left(\sum_{i \in P} \mathbf{k}_i \mathbf{k}_i^\top \right) - 2 \frac{\|\mathbf{w} - \mathbf{w}^*\|_{\mathcal{H}}}{\sqrt{\sigma_{\max}(\mathbf{K})}} \left\| \sum_{i \in P} \eta_i \mathbf{k}_i \right\| \quad (72)$$

$$\leq \eta^2 \frac{\sigma_{\max}(\sum_{i \in P} \mathbf{k}_i \mathbf{k}_i^\top)}{\sqrt{\sigma_{\min}(\mathbf{K})}} - 2\eta \frac{\left\| \sum_{i \in P} \eta_i \mathbf{k}_i \right\|}{\sqrt{\sigma_{\max}(\mathbf{K})}} \quad (73)$$

(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 (8).

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

We thus have \mathbf{w} is a stationary point if

$$\eta^2 \frac{\sigma_{\max}(\sum_{i \in P} \mathbf{k}_i \mathbf{k}_i^\top) - 2\sqrt{\sigma_{\min}(\mathbf{K})}}{2\lambda\sqrt{\sigma_{\min}(\mathbf{K})}} - 2\eta \frac{\left\| \sum_{i \in P} \eta_i \mathbf{k}_i \right\|}{\sqrt{\sigma_{\max}(\mathbf{K})}} \leq 0 \quad (74)$$

After rearranging we end up with

$$\eta \leq \frac{2\lambda\sqrt{\sigma_{\min}(\mathbf{K})}\|\sum_{i \in P} \eta_i \mathbf{k}_i\|}{\left(\sigma_{\max}(\sum_{i \in P} \mathbf{k}_i \mathbf{k}_i^\top) - 2\sqrt{\sigma_{\min}(\mathbf{K})}\right)\sqrt{\sigma_{\max}(\mathbf{K})}} \quad (75)$$

■

B.5 Proof of Lemma 21

Proof.

Let S be the set containing the points with the minimum error from X w.r.t to the weights vector \mathbf{w} . Define $\eta_i \triangleq (\mathbf{w}^{*\top} \mathbf{k}_i - y_i)$ where $i \in P$.

$$\Phi(\mathbf{w}) - \Phi(\mathbf{w}^*) = \sum_{i \in S} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 - \sum_{j \in P} (\mathbf{w}^{*\top} \mathbf{k}_j - y_j)^2 \quad (76)$$

$$= \sum_{i \in S \cap P} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 + \sum_{i \in S \cap Q} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 - \sum_{j \in P} (\mathbf{w}^{*\top} \mathbf{k}_j - y_j)^2 \quad (77)$$

$$\geq \sum_{i \in S \cap P} (\mathbf{w}^\top \mathbf{k}_i - y_i)^2 - \sum_{j \in P} (\mathbf{w}^{*\top} \mathbf{k}_j - y_j)^2 \quad (78)$$

$$= \sum_{i \in S \cap P} (\mathbf{w}^\top \mathbf{k}_i - \mathbf{w}^{*\top} \mathbf{k}_i - \eta_i)^2 - \sum_{j \in P} \eta_j^2 \quad (79)$$

$$= \sum_{i \in S \cap P} \left((\mathbf{w} - \mathbf{w}^*)^\top \mathbf{k}_i - \eta_i \right)^2 - \sum_{j \in P} \eta_j^2 \quad (80)$$

$$= \sum_{i \in S \cap P} \left\| (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{k}_i \right\|^2 - 2\eta_i (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{k}_i - \sum_{j \in P \setminus S} \eta_j^2 \quad (81)$$

$$\stackrel{\text{lem. 31}}{\geq} \underbrace{\sum_{i \in S \cap P} \left\| (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{k}_i \right\|^2}_{\Xi} - 2 \frac{\eta}{\sqrt{\sigma_{\min}(\mathbf{K})}} \underbrace{\left\| \sum_{i \in S \cap P} \eta_i \mathbf{k}_i \right\|}_{\beta} - \underbrace{\sum_{j \in P \setminus S} \eta_j^2}_{\phi} \quad (82)$$

Now we will upper bound Ξ .

$$\sum_{i \in S \cap P} \left\| (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{k}_i \right\|^2 \stackrel{(a)}{=} \sum_{i \in S \cap P} (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{k}_i \mathbf{k}_i^\top (\mathbf{w} - \mathbf{w}^*) \quad (83)$$

$$= (\mathbf{w} - \mathbf{w}^*)^\top \left(\sum_{i \in S \cap P} \mathbf{k}_i \mathbf{k}_i^\top \right) (\mathbf{w} - \mathbf{w}^*) \quad (84)$$

$$\geq \|\mathbf{w} - \mathbf{w}^*\|^2 \sigma_{\min} \left(\sum_{i \in S \cap P} \mathbf{k}_i \mathbf{k}_i^\top \right) \quad (85)$$

$$\stackrel{\text{lem. 31}}{\geq} \frac{1}{\sqrt{\sigma_{\max}(\mathbf{K})}} \|\mathbf{w} - \mathbf{w}^*\|_{\mathcal{H}}^2 \underbrace{\sigma_{\min} \left(\sum_{i \in S \cap P} \mathbf{k}_i \mathbf{k}_i^\top \right)}_{\xi} \quad (86)$$

We will now lower bound ξ , from [1], this should be non-zero. Define μ such that $\sum_{i \in P \cap S} \mathbb{E} [\mathbf{k}_i \mathbf{k}_i^\top] \preceq \mu \mathbf{I}$, then

$$\sigma_{\min} \left(\sum_{i \in S \cap P} \mathbf{k}_i \mathbf{k}_i^\top \right) = \sigma_{\min} \left(\left(\sum_{i \in S \cap P} \mathbf{k}_i \mathbf{k}_i^\top + \mu \mathbf{I} \right) - \mu \mathbf{I} \right) \quad (87)$$

$$\geq \mu - \sigma_{\max} \left(\sum_{i \in S \cap P} \mathbf{k}_i \mathbf{k}_i^\top - \mu \mathbf{I} \right) \quad (88)$$

We thus have

$$\Phi(\mathbf{w}) - \Phi(\mathbf{w}^*) \geq \frac{\eta^2 \sigma_{\min}(\sum_{i \in S \cap P} \mathbf{k}_i \mathbf{k}_i^\top)}{\sqrt{\sigma_{\max}(\mathbf{K})}} - 2 \frac{\eta \|\sum_{i \in S \cap P} \eta_i \mathbf{k}_i\|}{\sqrt{\sigma_{\min}(\mathbf{K})}} - \sum_{j \in P \setminus S} \eta_j^2 \quad (89)$$

This completes the proof. \blacksquare

B.6 Proof of Theorem 22

Proof. First,

$$\|\nabla \Phi_\lambda(\mathbf{w})\|_{\mathcal{H}} = \left\| \frac{1}{\lambda} (\mathbf{w} - \text{prox}_{\lambda, \Phi}(\mathbf{w})) \right\|_{\mathcal{H}} = \left\| \frac{1}{\lambda} \left(\mathbf{w} - \arg \min_{\hat{\mathbf{w}} \in \mathcal{K}} \left(\Phi(\hat{\mathbf{w}}) + \frac{1}{2\lambda} \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}}^2 \right) \right) \right\|_{\mathcal{H}} = 0 \quad (90)$$

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

$$\Phi(\hat{\mathbf{w}}) < \Phi(\tilde{\mathbf{w}}) + \frac{1}{2\lambda} \|\tilde{\mathbf{w}} - \hat{\mathbf{w}}\|_{\mathcal{H}}^2 \quad (91)$$

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

$$\Phi(\hat{\mathbf{w}}) - \Phi(\mathbf{w}^*) \leq \frac{1}{2\lambda} \|\hat{\mathbf{w}} - \mathbf{w}^*\|_{\mathcal{H}}^2 \quad (92)$$

We proceed by proof by contradiction. Assume $\|\hat{\mathbf{w}} - \mathbf{w}^*\| > \eta$, then if $\Phi(\hat{\mathbf{w}}) - \Phi(\mathbf{w}^*) > \frac{1}{2}\eta^2$, then we will have $\hat{\mathbf{w}}$ is not a stationary point, which will imply $\|\hat{\mathbf{w}} - \mathbf{w}^*\|_{\mathcal{H}} \leq \eta$. Therefore, we attempt to find the minimum value for η . From Lemma 21, we have we have

$$\Phi(\mathbf{w}) - \Phi(\mathbf{w}^*) \geq \frac{\eta^2 \sigma_{\min}(\sum_{i \in S \cap P} \mathbf{k}_i \mathbf{k}_i^\top)}{\sqrt{\sigma_{\max}(\mathbf{K})}} - 2 \frac{\eta \|\sum_{i \in S \cap P} \eta_i \mathbf{k}_i\|}{\sqrt{\sigma_{\min}(\mathbf{K})}} - \sum_{j \in P \setminus S} \eta_j^2 \quad (93)$$

$$\triangleq \frac{\eta^2 \xi}{\sqrt{\sigma_{\max}(\mathbf{K})}} - 2 \frac{\eta \beta}{\sqrt{\sigma_{\min}(\mathbf{K})}} - \phi \quad (94)$$

We thus have,

$$\frac{\eta^2 \xi}{\sqrt{\sigma_{\max}(\mathbf{K})}} - 2 \frac{\eta \beta}{\sqrt{\sigma_{\min}(\mathbf{K})}} - \phi \geq \frac{1}{2\lambda} \eta^2 \quad (95)$$

$$\eta^2 \left(\frac{\xi}{\sqrt{\sigma_{\max}(\mathbf{K})}} - \frac{1}{2\lambda} \right) - \eta \left(\frac{2\beta}{\sqrt{\sigma_{\min}(\mathbf{K})}} \right) - \phi \geq 0 \quad (96)$$

$$\eta^2 \underbrace{\left(\frac{2\xi\lambda - \sqrt{\sigma_{\max}(\mathbf{K})}}{2\lambda\sqrt{\sigma_{\max}(\mathbf{K})}} \right)}_A - \eta \underbrace{\left(\frac{2\beta}{\sqrt{\sigma_{\min}(\mathbf{K})}} \right)}_B - \phi \geq 0 \quad (97)$$

Here we can note from the definition of ϕ in Equation (82), $\phi \geq 0$, by the positive solution of the quadratic

$$\eta \geq \frac{-B + \sqrt{B^2 + 4A\phi}}{2A} \geq \frac{B\sqrt{1 + \frac{4A\phi}{B^2}} - B}{2A} \quad (98)$$

Therefore, when Equation (98) holds, we have a contradiction. It thus follows

$$\eta \leq \frac{B\sqrt{1 + \frac{4A\phi}{B^2}} - B}{2A} \leq \frac{\sqrt{4A\phi}}{2A} = \left(\sum_{j \in P \setminus S} \eta_j^2 \right)^{1/2} \left(\frac{2\xi\lambda - \sqrt{\sigma_{\max}(\mathbf{K})}}{2\lambda\sqrt{\sigma_{\max}(\mathbf{K})}} \right)^{1/2} \leq \left(\sum_{j \in P \setminus S} \eta_j^2 \right)^{1/2} \sqrt{\xi} \quad (99)$$

This completes the proof. \blacksquare

C Proofs for Section 4

In this section we give the optimization results from § section 4.

C.1 Proof of Theorem 28

D Auxiliary Lemmas

Lemma 31. Let \mathbf{K} be the kernel matrix s.t. $K_{i,j} \triangleq k(\mathbf{x}_i, \mathbf{x}_j)$, where $k : (\mathcal{X}, \mathcal{X}) \rightarrow \mathbb{R}$ represents the kernel function. Let \mathbf{w} and $\hat{\mathbf{w}}$ be two functions in the RKHS of k and $\{\mathbf{x}_i\}_{i \in X}$, then it follows

$$\frac{1}{\sqrt{\sigma_{\max}(\mathbf{K})}} \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}} \leq \|\mathbf{w} - \hat{\mathbf{w}}\|_2 \leq \frac{1}{\sqrt{\sigma_{\min}(\mathbf{K})}} \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}} \quad (100)$$

Proof. We will start with the definition of the RKHS norm.

$$\|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}}^2 = (\mathbf{w} - \hat{\mathbf{w}})^\top \mathbf{K} (\mathbf{w} - \hat{\mathbf{w}}) \quad (101)$$

$$\leq \sigma_{\max}(\mathbf{K}) \|\mathbf{w} - \hat{\mathbf{w}}\|_2^2 \quad (102)$$

Rearranging the inequality we have

$$\|\mathbf{w} - \hat{\mathbf{w}}\|_2 \geq \frac{\|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}}}{\sqrt{\sigma_{\max}(\mathbf{K})}} \quad (103)$$

Similarly, we have

$$\|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}}^2 \geq \sigma_{\min}(\mathbf{K}) \|\mathbf{w} - \hat{\mathbf{w}}\|_2^2 \quad (104)$$

Similarly rearrange inequality and the proof is complete. \blacksquare

Lemma 32. Let $a, b, c \in \mathbb{R}_+$, be positive real numbers, it then follows

E Base Learner Algorithm

In this section, we note that by using a base learner similar to Sever, [4], in each iteration, we obtain the Trimmed Maximum Likelihood Estimator from [1].

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 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  $\mathbf{w}_{(T)}$ 

```

The analysis for linear regression is in [1].

F Experimental Details

In this section we give details on datasets and hyperparameters.

F.1 Kernel Regression

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 6: 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.

F.2 Kernel Binary Classification

Dataset	Dimension d	Sample Size n	Source
Heart Disease	13	303	[9]
Breast Cancer	32	569	Kaggle

Table 7: Datasets for Kernel Binary Classification.

G Detailed Related Works

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

G.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

$$(\nabla_w F(w, u))^\top (\tilde{w} - w) \geq 0 \quad \forall \tilde{w} \in K \quad (105)$$

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.

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

First we will give the algorithm

$$S^{(t)} = \arg \min_{T \subset S^{(0)}: |T|=(1-2\epsilon)n} \sum_{i \in T} -\log f(y_i | \langle \beta^{(t)}, \mathbf{x}_i \rangle) \quad (106)$$

$$\beta^{(t+1)} = \arg \min_{\beta, \|\beta\| \leq R} \sum_{i \in S^{(t)}} -\log f(y_i | \langle \beta^{(t)}, \mathbf{x}_i \rangle) \quad (107)$$

In Equation (106), the algorithm chooses the $(1 - 2\epsilon)n$ points giving the least error and put this in the set $S^{(t)}$. Next, in Equation (107), 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].

$$\frac{1}{n} \sum_{i \in S} \nabla_\beta \log f(y_i | \langle \beta, \mathbf{x}_i \rangle)^\top \frac{(\beta^* - \beta)}{\|\beta^* - \beta\|} \leq \gamma \quad (108)$$

We see Equation (108) is an upper bound, instead of a lower bound, of Equation (105). 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.