

# Subquantile Minimization for Kernel Learning in the Huber $\epsilon$ -Contamination Model

Arvind Rathnashyam\*      Alex Gittens†

October 14, 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 [19, 33]. 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  $\epsilon$  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](mailto:rathna@rpi.edu)

†CS, Rensselaer Polytechnic Institute, [gittes@rpi.edu](mailto:gittes@rpi.edu)

# 1 Introduction

There has been extensive study of algorithms to learn the target distribution from a Huber  $\epsilon$ -Contaminated Model for a Generalized Linear Model (GLM), [5, 1, 20, 24, 9] as well as for linear regression [2, 22]. Robust Statistics has been studied extensively [6] for problems such as high-dimensional mean estimation [26, 3] and Robust Covariance Estimation [4, 8]. Recently, there has been an interest in solving robust machine learning problems by gradient descent [27, 5]. Subquantile minimization aims to address the shortcomings of standard ERM in applications of noisy/corrupted data [17, 15]. 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  $\epsilon$ -Contamination Model [13]).** 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  $\epsilon$ -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  $\epsilon n$  are arbitrarily corrupted by an adversary. Let  $K$  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:

$$\|\mathbf{f}_{\hat{\mathbf{w}}} - \mathbf{f}_{\mathbf{w}^*}\|_{\mathcal{H}} \left( \frac{2L \left( \sum_{j \in P \setminus S} \eta_j^2 \right)}{2L\sigma_{\min} \left( \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top \right) - 1} \right)^{1/2} + \frac{4L \left| \sum_{i \in S \cap P} \eta_i \sqrt{K(\mathbf{x}_i, \mathbf{x}_i)} \right|}{\sigma_{\min} \left( \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top \right)} + L\epsilon \quad (1)$$

Linear Regression:

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

Kernel Binary Classification:

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

Kernel Multi-Class Classification:

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

We will now state our main contributions clearly.

## Contributions

1. We propose a gradient-descent based algorithm for robust kernel learning in the Huber  $\epsilon$ -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  $\epsilon$ -Contamination Model and works in minimax optimization that will be relevant to our theoretical analysis.

### Robust Algorithms

[5] 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 data settings.

[20] 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 and lack of theoretical error bounds due to the nonlinearity induced by the exponential.

[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

[16] 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.

[37] 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  $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$  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

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

Given a convex objective function,  $\ell$ , the kernelized learning problem becomes:

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

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

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

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

**Claim 3.** *The function  $g(t, \mathbf{f}_{\mathbf{w}})$  defined in Equation (7) is non-convex-concave, i.e. it is not convex with respect to  $\mathbf{f}_{\mathbf{w}}$  and is concave with respect to  $t$ . Furthermore,  $g(t, \mathbf{f}_{\mathbf{w}})$  is  $\rho$ -weakly convex where  $\rho$  is the  $\beta$ -smoothness factor of  $g(t, \mathbf{f}_{\mathbf{w}})$  w.r.t  $\mathbf{f}_{\mathbf{w}}$ .*

**Proof.** We will show  $-g(t, \mathbf{f}_{\mathbf{w}})$  is convex with respect to  $t$ . Let  $\nu_i \triangleq (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2$ .

$$-g(\lambda t_1 + (1 - \lambda)t_2, \mathbf{f}_{\mathbf{w}}) = -\lambda t_1 - (1 - \lambda)t_2 + \sum_{i=1}^n (\lambda t_1 + (1 - \lambda)t_2 - \nu_i)^+ \quad (10)$$

$$\leq -\lambda t_1 - (1 - \lambda)t_2 + \sum_{i=1}^n \lambda(t_1 - \nu_i)^+ + (1 - \lambda)(t_2 - \nu_i)^+ \quad (11)$$

$$= -\lambda g(t_1, \mathbf{f}_{\mathbf{w}}) - (1 - \lambda)g(t_2, \mathbf{f}_{\mathbf{w}}) \quad (12)$$

Therefore we have  $g(t, \mathbf{f}_{\mathbf{w}})$  is concave in  $t$ . Next we will prove  $g(t, \mathbf{f}_{\mathbf{w}})$  is not convex in  $\mathbf{f}_{\mathbf{w}}$ . ■

**Claim 4.** *The function  $g(t, \mathbf{f}_{\mathbf{w}})$  defined in Equation (7) is  $L$ -weakly convex in  $\mathbf{f}_{\mathbf{w}}$ , where  $L$  is lipschitz constant of the gradient of  $g(t, \mathbf{f}_{\mathbf{w}})$  w.r.t  $\mathbf{f}_{\mathbf{w}}$ . This is true for Conditional Value at Risk, [? ].*

**Proof.** Here we note that if we add  $\max_{i \in [n]} |k(\mathbf{x}_i, \mathbf{x}_i) + \mathbf{f}_{\mathbf{w}}(\mathbf{x}_i) - y_i|$  to each of the second derivatives, then we are pushing the trace to be negative. Note this is the  $L$ -lipschitz gradient. This is equivalent to  $g(\mathbf{f}_{\mathbf{w}}, t) + \frac{L}{2} \|\mathbf{f}_{\mathbf{w}}\|_{\mathcal{H}}^2$ . ■

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:  
 $X, (n \times d), n \gg d$ ; Learning schedule:  
 $\alpha_1, \dots, \alpha_T$ ; Ridge parameter:  $\lambda$

**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:    $S_{(k)} \leftarrow \text{SUBQUANTILE}(\mathbf{w}^{(k)}, X)$ 
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:  $X, (n \times d)$ ,  
Convex Loss Function  $f$

**Output:** Subquantile Matrix  $S$

```
1:  $\hat{\nu}_i \leftarrow \ell(\mathbf{x}_i; \mathbf{f}_{\mathbf{w}}, y_i)$  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  
    $\ell(\mathbf{x}_i; \mathbf{f}_{\mathbf{w}}, y_i) \leq t$ 
4:  $S \leftarrow (\mathbf{x}_1^\top \quad \dots \quad \mathbf{x}_{np}^\top)^\top$ 
5: return  $S$ 
```

---

### 3 Structural Results

To consider theoretical guarantees of Subquantile Minimization, we first analyze the inner and outer optimization problems. We first analyze kernel learning in the presence of corrupted data. Next, we provide error bounds for the two most important kernel learning problems, kernel ridge regression, and kernel classification. Now we will give our first result regarding kernel learning in the Huber  $\epsilon$ -contamination model. Now we will analyze the two-step minimax optimization steps described in Equations (8) and (9).

**Lemma 5.** *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  $g$ .

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

$$= -1 + \frac{1}{np} \sum_{i=1}^{np} \begin{cases} 1 & \text{if } t > \hat{\nu}_i \\ 0 & \text{if } t < \hat{\nu}_i \\ [0, 1] & \text{if } t = \hat{\nu}_i \end{cases} \quad (14)$$

$$= 0 \text{ when } t = \hat{\nu}_{np} \quad (15)$$

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

It therefore follows,

$$\sum_{i=1}^n \mathbb{I} \left\{ \hat{\nu}_{np} \geq \left( f_{\mathbf{w}}^{(k)}(\mathbf{x}_i) - y_i \right)^2 \right\} \left( f_{\mathbf{w}}^{(k)}(\mathbf{x}_i) - y_i \right)^2 \in \max_{t \in \mathbb{R}} g(t, \mathbf{f}_{\mathbf{w}}^{(k)}) \quad (16)$$

**Interpretation 6.** From lemma 5, 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 7.** *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 5, it then follows  $\nabla_{\mathbf{w}} g(t^{(k)}, \mathbf{f}_{\mathbf{w}}^{(k)}) = \frac{1}{np} \sum_{i=1}^{np} \nabla f(\mathbf{x}_i; \mathbf{f}_{\mathbf{w}}^{(k)}, y_i)$*

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

$$\nabla_{\mathbf{f}} g(t^{(k+1)}, \mathbf{f}_{\mathbf{w}}^{(k)}) = \nabla_{\mathbf{f}} \left( \hat{\nu}_{np} - \frac{1}{np} \sum_{i=1}^n \left( \hat{\nu}_{np} - \ell(\mathbf{x}_i; \mathbf{f}_{\mathbf{w}}^{(k)}, y_i) \right)^+ \right) \quad (17)$$

$$= -\frac{1}{np} \sum_{i=1}^{np} \nabla_{\mathbf{f}} \left( \hat{\nu}_{np} - \ell(\mathbf{x}_i; \mathbf{f}_{\mathbf{w}}^{(k)}, y_i) \right)^+ \quad (18)$$

$$= \frac{1}{np} \sum_{i=1}^n \nabla_{\mathbf{f}} \ell(\mathbf{x}_i; \mathbf{f}_{\mathbf{w}}^{(k)}, y_i) \begin{cases} 1 & \text{if } t > \hat{\nu}_i \\ 0 & \text{if } t < \hat{\nu}_i \\ [0, 1] & \text{if } t = \hat{\nu}_i \end{cases} \quad (19)$$

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

$$\nabla_{\mathbf{f}} g(t^{(k+1)}, \mathbf{f}_{\mathbf{w}}^{(k)}) = \frac{1}{np} \sum_{i=1}^{np} \nabla_{\mathbf{f}} \ell(\mathbf{x}_i; \mathbf{f}_{\mathbf{w}}^{(k)}, y_i) \quad (20)$$

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.

### 3.1 On the Softplus Approximation

It is clear our objective function is non-smooth. Thus we propose to use the Softplus approximation to smooth the function.

$$\zeta_\gamma(x) = \frac{1}{\lambda} \log(1 + e^{\lambda x}) \quad (21)$$

We then have the approximation of  $g$  as

$$\tilde{g}_\lambda(t, f_{\mathbf{w}}) \triangleq t - \sum_{i=1}^n \zeta_\lambda(t - \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i)) \quad (22)$$

$$= t - \frac{1}{np} \sum_{i=1}^n \frac{1}{\lambda} \log(1 + \exp(\lambda(t - \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i)))) \quad (23)$$

Now we compute the derivatives,

$$\nabla_t \tilde{g}_\lambda(t, f_{\mathbf{w}}) = \nabla_t \left( t - \frac{1}{np} \sum_{i=1}^n \frac{1}{\lambda} \ln(1 + \exp(\lambda(t - \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i)))) \right) \quad (24)$$

$$= 1 - \frac{1}{np} \sum_{i=1}^n \sigma(\lambda(t - \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i))) \quad (25)$$

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

$$\lim_{\lambda \rightarrow \infty} \nabla_t \tilde{g}_\lambda(t, f_{\mathbf{w}}) = 1 - \frac{1}{np} \sum_{i=1}^n \mathbb{I}\{t - \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i)\} \quad (26)$$

$$\nabla_f \tilde{g}_\lambda(t, f_{\mathbf{w}}) = \nabla_f \left( t - \frac{1}{np} \sum_{i=1}^n \frac{1}{\lambda} \ln(1 + \exp(\lambda(t - \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i)))) \right) \quad (27)$$

$$= \frac{1}{np} \sum_{i=1}^n \nabla_f \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i) \sigma(\lambda(t - \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i))) \quad (28)$$

We therefore similarly have,

$$\lim_{\lambda \rightarrow \infty} \nabla_f \tilde{g}_\lambda(t, f_{\mathbf{w}}) = \frac{1}{np} \sum_{i=1}^n \mathbb{I}\{t - \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i)\} \nabla_f \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i) \quad (29)$$

Then the second derivative is given by

$$[\nabla_f^2 \tilde{g}_\lambda(t, f_{\mathbf{w}})]_{i,j} = [\nabla_f]_j \left( \frac{1}{np} \sum_{i=1}^n [\nabla_f]_i \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i) \sigma(\lambda(t - \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i))) \right) \quad (30)$$

$$\begin{aligned} &= \frac{1}{np} \sum_{i=1}^n \left( [\nabla_f^2]_{i,j} \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i) \sigma(\lambda(t - \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i))) \right. \\ &\quad \left. - [\nabla_f \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i)]_j [\nabla_f \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i)]_i \sigma(\lambda(t - \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i))) (1 - \sigma(\lambda(t - \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i)))) \right) \end{aligned} \quad (31)$$

We then similarly have,

$$\lim_{\lambda \rightarrow \infty} [\nabla_f^2 \tilde{g}_\lambda(t, f_{\mathbf{w}})]_{j,k} = \frac{1}{np} \sum_{i=1}^n \mathbb{I}\{t - \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i)\} [\nabla_f^2 \ell(f_{\mathbf{w}}; \mathbf{x}_i, y_i)]_{j,k} \quad (32)$$

### 3.2 Weakly Convex Concave Optimization Theory

With our smoothed function, we are now able to use the weakly-convex concave minimization literature to analyze  $g$ .

**Definition 8. (Moreau Envelope, [21]).** Let  $f$  be proper lower semi-continuous convex function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , then the Moreau Envelope is defined as:

$$f_\lambda(\mathbf{x}) \triangleq \inf_{\hat{\mathbf{x}} \in \mathcal{X}} \left( f(\hat{\mathbf{x}}) + \frac{1}{2\lambda} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \right) \quad (33)$$

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

**Definition 9. (Moreau Envelope on closed, convex set, [3]).** Let  $f$  be proper lower semi-continuous convex function  $\ell : \mathcal{K} \rightarrow \mathbb{R}$ , where  $\mathcal{K} \subset \mathcal{X}$  is a closed and convex set, then the Moreau Envelope is defined as:

$$M_{\lambda\ell}(\mathbf{f}_\mathbf{w}) \triangleq \inf_{\mathbf{f}_{\hat{\mathbf{w}}} \in \mathcal{K}} \left\{ \ell(\mathbf{f}_{\hat{\mathbf{w}}}) + \frac{1}{2\lambda} \|\mathbf{f}_\mathbf{w} - \mathbf{f}_{\hat{\mathbf{w}}}\|_{\mathcal{H}}^2 \right\} \quad (34)$$

**Lemma 10.** Let  $\ell$  be a proper convex function,  $\ell : \mathcal{K} \rightarrow \mathbb{R}$ , where  $\mathcal{K} \subset \mathcal{H}$  is a closed and convex set. Define a new function,  $\bar{\ell} : \mathcal{H} \rightarrow \mathbb{R}$ ,

$$\bar{\ell}(\mathbf{f}_\mathbf{w}) = \begin{cases} \ell(\mathbf{f}_\mathbf{w}) & \mathbf{f}_\mathbf{w} \in \mathcal{K} \\ \infty & \mathbf{f}_\mathbf{w} \notin \mathcal{K} \end{cases} = \ell(\mathbf{f}_\mathbf{w}) + \mathbb{I}\{\mathbf{f}_\mathbf{w} \in \mathcal{K}\} \quad (35)$$

where  $\mathbb{I}\{\cdot\}$  is the convex set indicator function. Then the derivative of the Moreau Envelope defined in Definition 9 is given by

$$\nabla M_{\lambda\bar{\ell}}(\mathbf{f}_\mathbf{w}) = \frac{1}{\lambda} \left( \mathbf{f}_\mathbf{w} - \arg \min_{\mathbf{f}_{\hat{\mathbf{w}}} \in \mathcal{K}} \left\{ \bar{\ell}(\mathbf{f}_{\hat{\mathbf{w}}}) + \frac{1}{2\lambda} \|\mathbf{f}_\mathbf{w} - \mathbf{f}_{\hat{\mathbf{w}}}\|_{\mathcal{H}}^2 \right\} \right) \quad (36)$$

**Proof.** It follows

$$M_{\lambda\bar{\ell}}(\mathbf{f}_\mathbf{w}) = \inf_{\mathbf{f}_{\hat{\mathbf{w}}} \in \mathcal{H}} \left\{ \bar{\ell}(\mathbf{f}_{\hat{\mathbf{w}}}) + \frac{1}{2\lambda} \|\mathbf{f}_\mathbf{w} - \mathbf{f}_{\hat{\mathbf{w}}}\|_{\mathcal{H}}^2 \right\} = \inf_{\mathbf{f}_{\hat{\mathbf{w}}} \in \mathcal{K}} \left\{ \ell(\mathbf{f}_{\hat{\mathbf{w}}}) + \frac{1}{2\lambda} \|\mathbf{f}_\mathbf{w} - \mathbf{f}_{\hat{\mathbf{w}}}\|_{\mathcal{H}}^2 \right\} \quad (37)$$

since we are assuming  $\mathcal{K}$  is non-empty. Then, from standard results in convex analysis, we have

$$\nabla M_{\lambda\bar{\ell}}(\mathbf{f}_\mathbf{w}) = \frac{1}{\lambda} (\mathbf{f}_\mathbf{w} - \text{prox}_{\lambda\bar{\ell}}(\mathbf{f}_\mathbf{w})) = \frac{1}{\lambda} \left( \mathbf{f}_\mathbf{w} - \arg \min_{\mathbf{f}_{\hat{\mathbf{w}}} \in \mathcal{H}} \left\{ \bar{\ell}(\mathbf{f}_{\hat{\mathbf{w}}}) + \frac{1}{2\lambda} \|\mathbf{f}_\mathbf{w} - \mathbf{f}_{\hat{\mathbf{w}}}\|_{\mathcal{H}}^2 \right\} \right) \quad (38)$$

$$= \frac{1}{\lambda} \left( \mathbf{f}_\mathbf{w} - \arg \min_{\mathbf{f}_{\hat{\mathbf{w}}} \in \mathcal{K}} \left\{ \ell(\mathbf{f}_{\hat{\mathbf{w}}}) + \frac{1}{2\lambda} \|\mathbf{f}_\mathbf{w} - \mathbf{f}_{\hat{\mathbf{w}}}\|_{\mathcal{H}}^2 \right\} \right) \quad (39)$$

This concludes the proof. ■

With Lemma 10, we now have a finite Lipschitz Constant for the Kernel Regression which will allow us to use theory in Moreau Envelopes for convergence analysis. Our analysis will continue by looking at the function  $\bar{\ell}$  rather than  $\ell$ , although since we are doing a projection on to the convex set, so  $\bar{\ell}$  does not change our gradient descent.

**Remark 11.** Define the function  $\Phi(\mathbf{f}_\mathbf{w}) \triangleq \max_{t \in \mathbb{R}} g(t, \mathbf{f}_\mathbf{w})$ . This function is a  $L$ -weakly convex function in  $\mathcal{K}$ , i.e.,  $\Phi(\mathbf{f}_\mathbf{w}) + \frac{L}{2} \|\mathbf{f}_\mathbf{w}\|_{\mathcal{H}}^2$  is a convex function over  $\mathbf{w}$  in the convex and compact set  $\mathcal{K}$ .

**Assumption 12.** Define  $\Phi(\cdot)$  as the function in remark 11. Then it follows  $\arg \min_{\mathbf{w} \in \mathbb{R}^d} \Phi(\mathbf{w}) = \mathbf{w}^*$



### 3.3 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

$$\ell(\mathbf{x}, y_i; \mathbf{f}_{\mathbf{w}}) = (\mathbf{f}_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 \quad (40)$$

For our theory, we need the  $L$ -lipschitz constant and  $\beta$ -smoothness constant.

**Lemma 13.** ( *$L$ -Lipschitz of  $g(t, \mathbf{f}_{\mathbf{w}})$  w.r.t  $\mathbf{f}_{\mathbf{w}}$* ). Let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ , represent the data vectors. It then follows for any  $\mathbf{f}_{\mathbf{w}}, \mathbf{f}_{\hat{\mathbf{w}}} \in \mathcal{K}$ :

$$|g(t, \mathbf{f}_{\mathbf{w}}) - g(t, \mathbf{f}_{\hat{\mathbf{w}}})| \leq L \|\mathbf{f}_{\mathbf{w}} - \mathbf{f}_{\hat{\mathbf{w}}}\|_{\mathcal{H}} \quad (41)$$

where  $L = \frac{2R}{np} (\sum_{i=1}^n \sqrt{K(\mathbf{x}_i, \mathbf{x}_i)})^2 + \frac{2\|\mathbf{y}\|}{np} (\sum_{i=1}^n \sqrt{K(\mathbf{x}_i, \mathbf{x}_i)})$

**Lemma 14.** ( *$\beta$ -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:

$$\|\nabla_{\mathbf{w}} g(t, \mathbf{f}_{\mathbf{w}}) - \nabla_{\mathbf{f}} g(t, \mathbf{f}_{\hat{\mathbf{w}}})\| \leq \beta \|\mathbf{f}_{\mathbf{w}} - \mathbf{f}_{\hat{\mathbf{w}}}\|_{\mathcal{H}} \quad (42)$$

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  $t \geq (\mathbf{f}_{\mathbf{w}}(\mathbf{x}) - y)^2$ . Since  $g$  is twice differentiable, we will analyze the Hessian.

$$\|\nabla_{\mathbf{f}}^2 g(t, \mathbf{f}_{\mathbf{w}})\|_{\mathcal{H}} = \quad (43)$$

This concludes the proof.  $\blacksquare$

Similar results for the Lipschitz Constant for non-kernelized learning algorithms can be seen in [38].

### 3.4 Kernel Binary Classification

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

$$\ell(\mathbf{x}_i, y_i; \mathbf{f}_{\mathbf{w}}) = -(y_i \log(\sigma(\mathbf{f}_{\mathbf{w}}(\mathbf{x}_i))) + (1 - y_i) \log(1 - \sigma(\mathbf{f}_{\mathbf{w}}(\mathbf{x}_i)))) \quad (44)$$

Similar to § section 3.3, we require the  $L$ -Lipschitz constant and  $\beta$ -smoothness constant.

**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:

$$|g(t, \mathbf{f}_{\mathbf{w}}) - g(t, \mathbf{f}_{\hat{\mathbf{w}}})| \leq L \|\mathbf{f}_{\mathbf{w}} - \mathbf{f}_{\hat{\mathbf{w}}}\|_{\mathcal{H}} \quad (45)$$

where  $L =$

**Lemma 16.** ( *$\beta$ -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:

$$\|\nabla_{\mathbf{f}} g(t, \mathbf{f}_{\mathbf{w}}) - \nabla_{\mathbf{f}} g(t, \mathbf{f}_{\hat{\mathbf{w}}})\| \leq \beta \|\mathbf{f}_{\mathbf{w}} - \mathbf{f}_{\hat{\mathbf{w}}}\|_{\mathcal{H}} \quad (46)$$

where  $\beta =$

### 3.5 Kernel Multi-Class Classification

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

$$\ell(\mathbf{x}_i, y_i; \mathbf{W}) = - \sum_{j=1}^{|C|} \mathbb{I}\{y_i = j\} \log \left( \frac{\exp(\mathbf{W}_k^\top \mathbf{k}_i)}{\sum_{h=1}^{|C|} \exp(\mathbf{W}_h^\top \mathbf{k}_i)} \right) \quad (47)$$

**Lemma 17.** ( *$L$ -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:

$$|g(t, \mathbf{w}) - g(t, \hat{\mathbf{w}})| \leq L \|\mathbf{w} - \hat{\mathbf{w}}\|_{\mathcal{H}} \quad (48)$$

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

### 3.6 Necessary Kernel Inequalities

We will first extend the idea of Resilience [34] to kernel learning.

**Definition 18. (Resilience)** from [34]. Let  $\mathcal{H}$  represent the RKHS associated with the proper kernel  $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ , then given the feature mapping  $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$ , and the set  $X = \{\mathbf{x}_i\}_{i=1}^n = P \cup Q$ , such that  $|P| = n(1 - \epsilon)$  and  $|Q| = n\epsilon$ , It then follows for any subset  $T \subseteq P$  such that  $|T| = (1 - 2\epsilon)n$ ,

$$\left\| \frac{1}{|T|} \sum_{i \in T} \phi(\mathbf{x}_i) - \mu_{\mathbb{P}} \right\| \leq \tau$$

where  $\mu_{\mathbb{P}} = \mathbb{E}_{\mathbf{x} \sim \mathbb{P}}[\phi(\mathbf{x})]$  is the kernel mean embedding for the distribution,  $\mathbb{P}$ . We say the set  $X$  has  $(\epsilon, \tau)$ -resilience in the Reproducing Kernel Hilbert Space.

Without the idea of resilience defined in definition 18, we will be unable to put error bounds on our algorithm.

**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

$$\|\nabla \Phi_{\lambda}(\mathbf{w})\|_{\mathcal{H}} = 0 \quad (49)$$

i.e.

$$\mathbf{f}_{\mathbf{w}} = \arg \min_{\mathbf{f}_{\hat{\mathbf{w}}} \in \mathcal{K}} \left( \Phi(\mathbf{f}_{\hat{\mathbf{w}}}) + \frac{1}{2\lambda} \|\mathbf{f}_{\mathbf{w}} - \mathbf{f}_{\hat{\mathbf{w}}}\|_{\mathcal{H}}^2 \right) \quad (50)$$

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

$$\Phi(\mathbf{f}_{\mathbf{w}}) - \Phi(\mathbf{f}_{\mathbf{w}^*}) \geq \eta^2 \sigma_{\min} \left( \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^{\top} \right) - 2\eta \left( \sum_{i \in S \cap P} K(\mathbf{x}_i, \mathbf{x}_i) \right)^{1/2} \left( \sum_{i \in S \cap P} \eta_i^2 \right)^{1/2} - \sum_{j \in P \setminus S} \eta_j^2 \quad (51)$$

The proof is deferred to § ??.

**Theorem 21.** Let  $\mathbf{f}_{\hat{\mathbf{w}}}$  be a stationary point defined in ?? for the function  $\Phi$  defined in ??. Then,

$$\eta \leq \left( \frac{2L \left( \sum_{j \in P \setminus S} \eta_j^2 \right)}{2L \sigma_{\min} \left( \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^{\top} \right) - 1} \right)^{1/2} + \frac{4L \left| \sum_{i \in S \cap P} \eta_i \sqrt{K(\mathbf{x}_i, \mathbf{x}_i)} \right|}{\sigma_{\min} \left( \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^{\top} \right)} \quad (52)$$

where  $L$  is the Lipschitz Constant given in Lemma 13.

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

**Lemma 22.** ([32, 31? ]). Assume the function  $\Phi$  is  $\ell$ -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:

$$\|\hat{\mathbf{w}} - \mathbf{w}\| = \lambda \epsilon \text{ and } \min_{\mathbf{g} \in \partial \Phi(\hat{\mathbf{w}}) + \partial \mathcal{I}_{\mathcal{K}}(\hat{\mathbf{w}})} \|\mathbf{g}\| \leq \epsilon \quad (53)$$

*How to extend this to Hilbert Space Norm?*

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

$$\partial \mathcal{I}_{\mathcal{K}} = \mathcal{N}(\hat{\mathbf{w}}) = \{\mathbf{f}_{\hat{\mathbf{w}}} \in \mathbb{R}^n \mid \langle \hat{\mathbf{w}}, \bar{\mathbf{w}} - \hat{\mathbf{w}} \rangle \leq \forall \bar{\mathbf{w}} \in \mathcal{K}\} \quad (54)$$

We will define a convex cone.

**Definition 23.** A set  $\Omega$  is a cone if  $\lambda x \in \Omega$  whenever  $x \in \Omega$  and  $\lambda \geq 0$ , if  $\Omega$  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 24.** *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*

$$\|\hat{\mathbf{w}} - \mathbf{f}_{\mathbf{w}^*}\|_{\mathcal{H}} \leq \Xi \quad (55)$$

**Definition 25. (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 8. Then we say that a point  $\mathbf{w}$  is a  $\rho$ -approximate stationary point if  $\|f_\beta(\mathbf{w})\|_2 \leq \rho$ .

The approximate stationary point in definition 25 is used in the analysis of the minimax algorithm in [3]. First, if you can prove a stationary point is good, theorem 21, then using lemma 22, 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  $\Phi$  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  $\mu$ -approximate stationary point. However, since our methods are in the kernelized setting. The 2-norm,  $\|\mathbf{w} - \mathbf{w}^*\|$  is not sufficient, we want  $\|\mathbf{w} - \mathbf{w}^*\|_{\mathcal{H}}$  to be close, as the RKHS being small indicates the function  $f(\mathbf{x})$  and  $f^*(\mathbf{x})$  will be close.

### 4.1 Optimization in the Reproducing Kernel Hilbert Space

In this section, we will discuss and give necessary optimization results in the RKHS norm.

### 4.2 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.

#### 4.2.1 Momentum Accelerated Gradient Descent

In this section we study Momentum Accelerated Gradient Descent [28, 25] with our non-convex-concave optimization algorithm.

$$\mathbf{b}^{(t)} = \mu \mathbf{b}^{(t-1)} + \nabla_{\mathbf{f}} \Phi \left( \mathbf{f}_{\mathbf{w}}^{(t-1)} \right) \quad (56)$$

$$\mathbf{f}_{\mathbf{w}}^{(t)} = \mathbf{f}_{\mathbf{w}}^{(t-1)} - \alpha \mathbf{b}^{(t)} \quad (57)$$

**Theorem 26.** *Momentum Accelerated Gradient Descent given in Equations (56) and (57) reaches a  $\eta$ -approximate stationary point. Algorithm 1 reaches a  $\eta$ -approximate stationary point in a polynomial number of iterations.*

$$\mathbb{E} \left[ \left\| \nabla \Phi_{1/2\ell}(\mathbf{f}_{\mathbf{w}}) \right\|_{\mathcal{H}}^2 \right] \leq \quad (58)$$

#### 4.2.2 Nesterov Accelerated Gradient Descent

In this section we study Nesterov Accelerated Gradient Descent [23] with our non-convex-concave optimization algorithm.

$$\mathbf{b}^{(t+1)} = (1 + \mu) \mathbf{f}_{\mathbf{w}}^{(t)} - \mu \mathbf{f}_{\mathbf{w}}^{(t-1)} \quad (59)$$

$$\mathbf{f}_{\mathbf{w}}^{(t+1)} = \mathbf{b}^{(t+1)} - \alpha \nabla_{\mathbf{f}} \Phi \left( \mathbf{f}_{\mathbf{w}}^{(t)} \right) \quad (60)$$

**Theorem 27.** *Nesterov Accelerated Gradient Descent given in Equations (59) and (60) reaches a  $\eta$ -approximate stationary point. Algorithm 1 reaches a  $\eta$ -approximate stationary point in a polynomial number of iterations.*

$$\mathbb{E} \left[ \left\| \nabla \Phi_{1/2\ell}(\mathbf{f}_{\mathbf{w}}) \right\|_{\mathcal{H}}^2 \right] \leq \quad (61)$$

## 5 Experiments

We perform numerical experiments on state of the art datasets comparing with other state of the art methods. We initialize the weights parameterizing  $f_{\mathbf{w}}$  with the Glorot Initialization Scheme [10].

---

### Algorithm 3: SUBQUANTILE-KERNEL

---

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

**Output:** Trained Parameters:  $\mathbf{f}_{\mathbf{w}}^{(T)}$

```

1:  $\mathbf{w}_i^{(0)} \leftarrow \text{Unif} \left[ -\sqrt{\frac{6}{n}}, \sqrt{\frac{6}{n}} \right], \forall i \in [n]$  ▷ Base Learner
2: for  $k = 1, 2, \dots, T$  do
3:    $S^{(k)} \leftarrow \text{SUBQUANTILE}(\mathbf{f}_{\mathbf{w}}^{(k)}, \mathbf{X})$  ▷ Algorithm 2
4:    $\nabla_{\mathbf{f}} g \left( t^{(k+1)}, \mathbf{f}_{\mathbf{w}}^{(k)} \right) \leftarrow 2 \sum_{i \in S^{(k)}} \left( \mathbf{f}_{\mathbf{w}}^{(k)}(\mathbf{x}_i) - y_i \right) \cdot K(\mathbf{x}_i, \cdot)$  ▷ Regression
5:    $\nabla_{\mathbf{f}} g \left( t^{(k+1)}, \mathbf{f}_{\mathbf{w}}^{(k)} \right) \leftarrow \sum_{i \in S^{(k)}} \left( \sigma \left( \mathbf{f}_{\mathbf{w}}^{(k)}(\mathbf{x}_i) \right) - y_i \right) \cdot K(\mathbf{x}_i, \cdot)$  ▷ Binary Classification
6:    $\mathbf{f}_{\mathbf{w}}^{(k+1)} \leftarrow \mathbf{f}_{\mathbf{w}}^{(k)} - \alpha_{(k)} \nabla_{\mathbf{f}} g \left( t^{(k+1)}, \mathbf{f}_{\mathbf{w}}^{(k)} \right)$  ▷  $\mathbf{w}$ -update in eqn. (9)
7: end
8: return  $\mathbf{f}_{\mathbf{w}}^{(T)}$ 

```

---

Algorithms	Test RMSE							
	Concrete		Wine Quality		Boston Housing		Drug	
	$\epsilon = 0.2(\downarrow)$	$\epsilon = 0.4(\downarrow)$	$\epsilon = 0.2(\downarrow)$	$\epsilon = 0.4(\downarrow)$	$\epsilon = 0.2(\downarrow)$	$\epsilon = 0.4(\downarrow)$	$\epsilon = 0.2(\downarrow)$	$\epsilon = 0.4(\downarrow)$
KRR	1.355 <sub>(0.0934)</sub>	2.282 <sub>(0.2063)</sub>	1.437 <sub>(0.0979)</sub>	2.272 <sub>(0.1088)</sub>	1.285 <sub>(0.0896)</sub>	2.266 <sub>(0.0686)</sub>	1.478 <sub>(0.0533)</sub>	2.381 <sub>(0.0203)</sub>
TERM [20]	0.829 <sub>(0.0422)</sub>	0.928 <sub>(0.0197)</sub>	1.854 <sub>(0.7437)</sub>	1.069 <sub>(0.1001)</sub>	0.879 <sub>(0.0178)</sub>	0.875 <sub>(0.0711)</sub>	$\infty$	$\infty$
SEVER [5]	<u>0.533</u> <sub>(0.0347)</sub>	<u>0.592</u> <sub>(0.0548)</sub>	<u>0.915</u> <sub>(0.0343)</sub>	<u>0.841</u> <sub>(0.0413)</sub>	<u>0.526</u> <sub>(0.0287)</sub>	<u>0.720</u> <sub>(0.1147)</sub>	<b>1.172</b> <sub>(0.0542)</sub>	<b>1.215</b> <sub>(0.0536)</sub>
SUBQUANTILE	<b>0.519</b> <sub>(0.0134)</sub>	<b>0.547</b> <sub>(0.0174)</sub>	<b>0.808</b> <sub>(0.0389)</sub>	<b>0.827</b> <sub>(0.0216)</sub>	<b>0.468</b> <sub>(0.0896)</sub>	<b>0.458</b> <sub>(0.0662)</sub>	<u>1.280</u> <sub>(0.0568)</sub>	<u>1.372</u> <sub>(0.0294)</sub>
Genie ERM	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$

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

In 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 ?? and figures fig. 1, subquantile minimization has near optimal performance in the Polynomial Regression Synthetic Dataset.

### 5.1 Linear Regression

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

### 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.

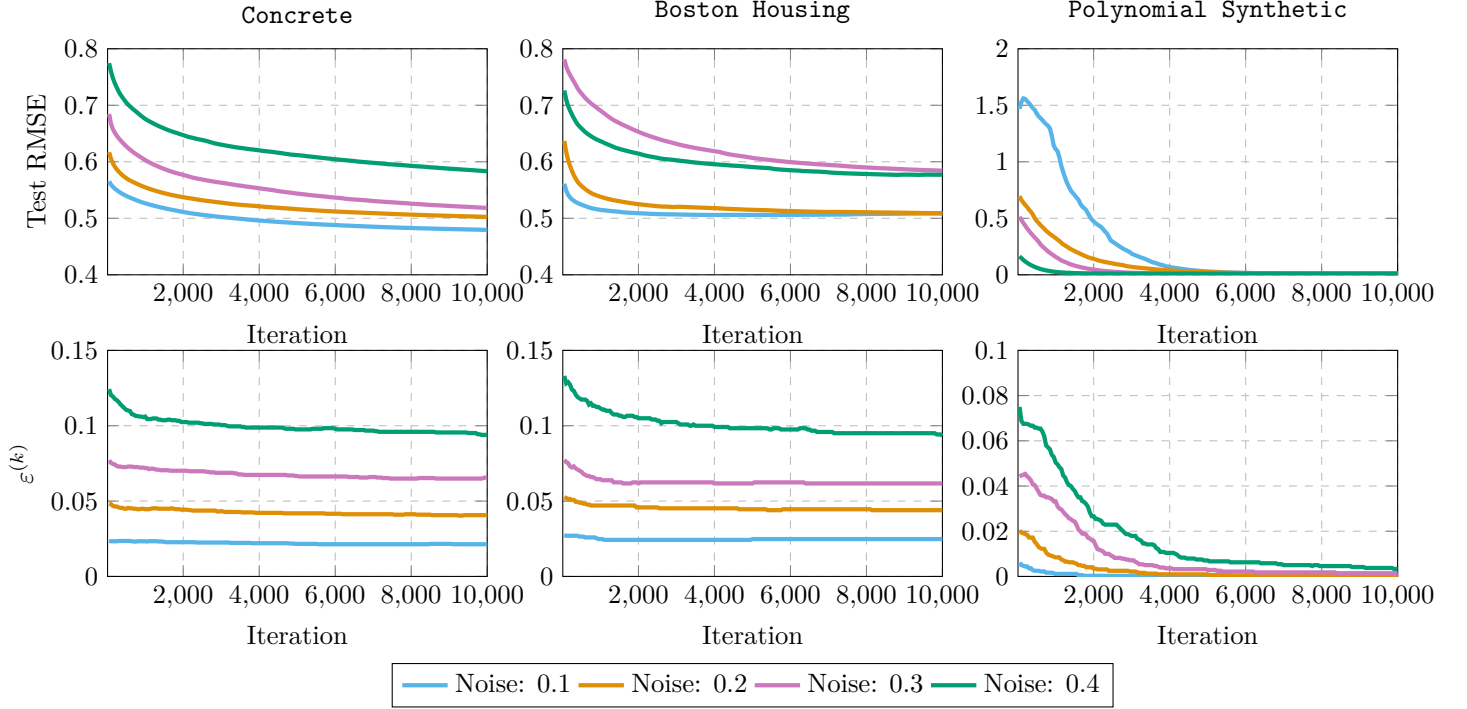


Figure 1: Test RMSE over the iterations in **Concrete**, **Boston Housing**, and **Polynomial** Datasets for SUBQUANTILE at different noise levels

Algorithms	Test RMSE							
	Boston Housing		Wine Quality		Concrete		Drug	
	Label(↓)	Label+Feature	Label	Label+Feature	Label	Label+Feature	Label	Label+Feature
KRR	0.907 <sub>(0.2724)</sub>	90.799 <sub>(5.7170)</sub>	0.894 <sub>(0.0404)</sub>	62.913 <sub>(7.4959)</sub>	0.747 <sub>(0.0465)</sub>	77.383 <sub>(5.5692)</sub>	2.679 <sub>(0.1286)</sub>	141.690 <sub>(3.5297)</sub>
RANSAC [9]	1.167 <sub>(0.6710)</sub>	22.460 <sub>(19.1987)</sub>	1.489 <sub>(0.2730)</sub>	39.630 <sub>(13.0294)</sub>	0.870 <sub>(0.2308)</sub>	23.629 <sub>(16.1023)</sub>	2.801 <sub>(0.2004)</sub>	117.389 <sub>(8.3915)</sub>
CRR [2]	0.636 <sub>(0.0905)</sub>	88.626 <sub>(5.7380)</sub>	<b>0.814</b> <sub>(0.0481)</sub>	58.488 <sub>(3.5612)</sub>	0.710 <sub>(0.0919)</sub>	73.932 <sub>(4.7867)</sub>	1.887 <sub>(0.1463)</sub>	152.827 <sub>(6.6038)</sub>
STIR [22]	<u>0.562</u> <sub>(0.0626)</sub>	78.878 <sub>(8.0164)</sub>	0.828 <sub>(0.0293)</sub>	58.352 <sub>(4.6700)</sub>	<b>0.684</b> <sub>(0.0245)</sub>	76.555 <sub>(4.5927)</sub>	1.721 <sub>(0.1520)</sub>	144.975 <sub>(5.4953)</sub>
SEVER [5]	0.601 <sub>(0.0979)</sub>	5.980 <sub>(8.2603)</sub>	<u>0.818</u> <sub>(0.0373)</sub>	9.065 <sub>(13.7632)</sub>	<b>0.684</b> <sub>(0.0438)</sub>	4.119 <sub>(8.2436)</sub>	1.469 <sub>(0.1162)</sub>	156.043 <sub>(4.5543)</sub>
TERM [20]	0.608 <sub>(0.1357)</sub>	<u>0.569</u> <sub>(0.0620)</sub>	0.840 <sub>(0.0563)</sub>	<u>0.827</u> <sub>(0.0255)</sub>	0.830 <sub>(0.0934)</sub>	<u>0.808</u> <sub>(0.0726)</sub>	<b>1.185</b> <sub>(0.1077)</sub>	<b>1.147</b> <sub>(0.1258)</sub>
SUBQUANTILE	<b>0.503</b> <sub>(0.0470)</sub>	<b>0.560</b> <sub>(0.0373)</sub>	0.838 <sub>(0.0260)</sub>	<b>0.821</b> <sub>(0.0305)</sub>	0.757 <sub>(0.1174)</sub>	<b>0.630</b> <sub>(0.0269)</sub>	<u>1.244</u> <sub>(0.1091)</sub>	<u>2.413</u> <sub>(0.6737)</sub>
Genie ERM	0.630 <sub>(0.1015)</sub>	0.665 <sub>(0.1134)</sub>	0.838 <sub>(0.0130)</sub>	0.865 <sub>(0.0222)</sub>	0.763 <sub>(0.0390)</sub>	0.768 <sub>(0.0181)</sub>	0.988 <sub>(0.0823)</sub>	0.985 <sub>(0.0838)</sub>

Table 2: For only Label Noise,  $y_{\text{noisy}} \sim \mathcal{N}(5, 5)$ . For Label and Feature Noise  $\mathbf{x}_{\text{noisy}} = 100\mathbf{x}_{\text{original}}$  and  $y_{\text{noisy}} = 10000y_{\text{original}}$ .

### 5.3 Kernel Multi-Class Classification

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

### 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.

Algorithms	Test Accuracy							
	Heart Disease				Breast Cancer			
	Label		Label+Feature		Label		Label+Feature	
	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$	$\epsilon = 0.2(\uparrow)$	$\epsilon = 0.4(\uparrow)$
SVM	<u>0.826</u> <sub>(0.0482)</sub>	<b>0.625</b> <sub>(0.0475)</sub>	0.549 <sub>(0.0763)</sub>	0.513 <sub>(0.0344)</sub>	<u>0.940</u> <sub>(0.0140)</sub>	0.781 <sub>(0.0496)</sub>	0.651 <sub>(0.0370)</sub>	0.633 <sub>(0.0428)</sub>
SEVER [5]	0.728 <sub>(0.1134)</sub>	0.531 <sub>(0.0321)</sub>	0.807 <sub>(0.0557)</sub>	<b>0.677</b> <sub>(0.1552)</sub>	0.918 <sub>(0.0439)</sub>	0.575 <sub>(0.1456)</sub>	0.954 <sub>(0.0146)</sub>	0.809 <sub>(0.0947)</sub>
TERM [20]	0.790 <sub>(0.0420)</sub>	<u>0.610</u> <sub>(0.0667)</sub>	<u>0.816</u> <sub>(0.0485)</sub>	<u>0.607</u> <sub>(0.1318)</sub>	0.937 <sub>(0.0151)</sub>	<b>0.793</b> <sub>(0.0691)</sub>	<b>0.972</b> <sub>(0.0151)</sub>	<b>0.973</b> <sub>(0.0190)</sub>
SUBQUANTILE	<b>0.846</b> <sub>(0.0382)</sub>	0.597 <sub>(0.0794)</sub>	<b>0.830</b> <sub>(0.0471)</sub>	0.516 <sub>(0.0623)</sub>	<b>0.956</b> <sub>(0.0055)</sub>	<u>0.785</u> <sub>(0.1059)</sub>	<u>0.955</u> <sub>(0.0159)</sub>	<u>0.823</u> <sub>(0.1445)</sub>
Genie ERM	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$

Table 3: Heart Disease and Breast Cancer Dataset. Label Noise:  $y_{\text{noise}} = \mathbb{I}\{y_{\text{original}} = 0\}$ . Feature Noise:  $\mathbf{x}_{\text{noise}} = 100\mathbf{x}_{\text{original}}$ . The Radial Basis Function is used in all experiments.

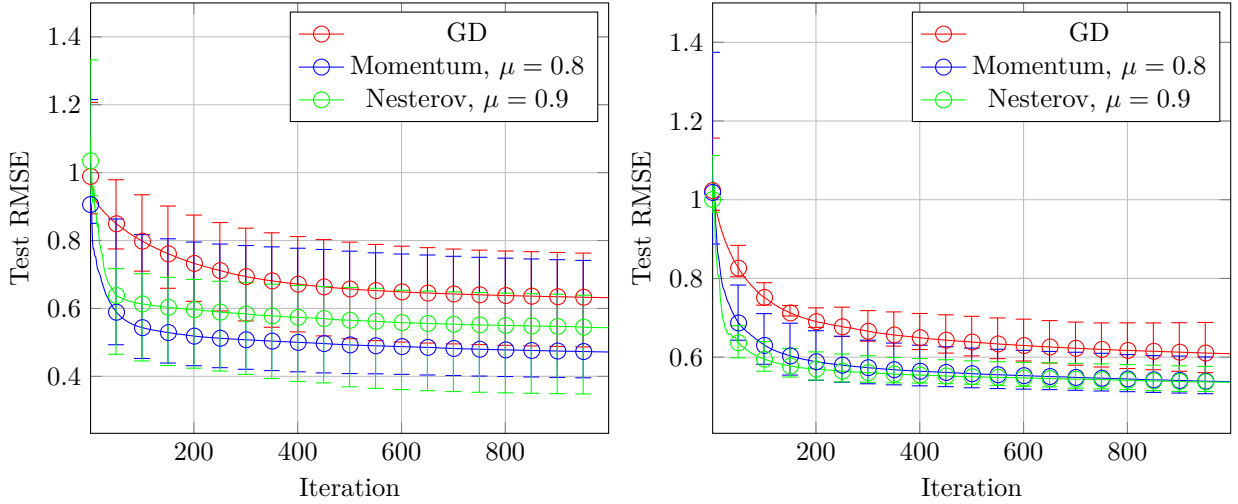


Figure 2: The effect of Momentum for Kernel Regression in the Boston Housing dataset (Left) and Concrete dataset (Right). We observe faster convergence. In both experiments, we use the Radial Basis Function.

## 6 Discussion

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

**Theory.** From our theoretical bounds, we see the more good data we obtain, the closer the resultant function is to the optimal function in the RKHS, an important idea in learning theory. We see we obtain an approximately close function to the optimal with high probability. Furthermore, we extend minimax optimization with maximizing oracle gradient descent literature with novel bounds for accelerated gradient descent.

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

**Interpretability.** One of the strengths in Subquantile Optimization is the high interpretability. Once training is finished, we can see the  $n(1 - p)$  points with highest error to find the outliers. 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.

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



## 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] Yu Cheng, Ilias Diakonikolas, Rong Ge, and David P Woodruff. Faster algorithms for high-dimensional robust covariance estimation. In *Conference on Learning Theory*, pages 727–757. PMLR, 2019.
- [5] 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.
- [6] Ilias Diakonikolas and Daniel M Kane. *Algorithmic high-dimensional robust statistics*. Cambridge University Press, 2023.
- [7] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
- [8] Jianqing Fan, Weichen Wang, and Yiqiao Zhong. An  $\infty$  eigenvector perturbation bound and its application to robust covariance estimation. *Journal of Machine Learning Research*, 18(207):1–42, 2018.
- [9] 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.
- [10] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010.
- [11] Arthur Gretton. Introduction to rkhs, and some simple kernel algorithms. *Adv. Top. Mach. Learn. Lecture Conducted from University College London*, 16(5-3):2, 2013.
- [12] David Gross. Recovering low-rank matrices from few coefficients in any basis. *IEEE Transactions on Information Theory*, 57(3):1548–1566, 2011.
- [13] Peter J. Huber and Elvezio Ronchetti. *Robust statistics*. Wiley series in probability and statistics. Wiley, Hoboken, N.J., 2nd ed. edition, 2009.
- [14] Steinbrunn William Pfisterer Matthias Janosi, Andras and Robert Detrano. Heart Disease. UCI Machine Learning Repository, 1988. DOI: <https://doi.org/10.24432/C52P4X>.
- [15] 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.
- [16] 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.
- [17] Ashish Khetan, Zachary C. Lipton, and Anima Anandkumar. Learning from noisy singly-labeled data. In *International Conference on Learning Representations*, 2018.

- [18] Jonas Moritz Kohler and Aurelien Lucchi. Sub-sampled cubic regularization for non-convex optimization. In *International Conference on Machine Learning*, pages 1895–1904. PMLR, 2017.
- [19] Yassine Laguel, Krishna Pillutla, Jérôme Malick, and Zaid Harchaoui. Superquantiles at work: Machine learning applications and efficient subgradient computation. *Set-Valued and Variational Analysis*, 29(4):967–996, Dec 2021.
- [20] Tian Li, Ahmad Beirami, Maziar Sanjabi, and Virginia Smith. Tilted empirical risk minimization. In *International Conference on Learning Representations*, 2021.
- [21] Jean-Jacques Moreau. Proximité et dualité dans un espace hilbertien. *Bulletin de la Société mathématique de France*, 93:273–299, 1965.
- [22] 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.
- [23] Yurii Evgen’evich Nesterov. A method of solving a convex programming problem with convergence rate  $O(k^{-2})$ . In *Doklady Akademii Nauk*, volume 269, pages 543–547. Russian Academy of Sciences, 1983.
- [24] 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.
- [25] Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *Ussr computational mathematics and mathematical physics*, 4(5):1–17, 1964.
- [26] Adarsh Prasad, Sivaraman Balakrishnan, and Pradeep Ravikumar. A unified approach to robust mean estimation. *arXiv preprint arXiv:1907.00927*, 2019.
- [27] Adarsh Prasad, Arun Sai Suggala, Sivaraman Balakrishnan, and Pradeep Ravikumar. Robust estimation via robust gradient estimation. *arXiv preprint arXiv:1802.06485*, 2018.
- [28] Ning Qian. On the momentum term in gradient descent learning algorithms. *Neural networks*, 12(1):145–151, 1999.
- [29] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc., 2007.
- [30] 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.
- [31] R Tyrrell Rockafellar. *Convex analysis*. 2015.
- [32] Ralph Tyrell Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, 1970.
- [33] 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.
- [34] 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.

- [35] Joel A Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of computational mathematics*, 12:389–434, 2012.
- [36] Joel A Tropp et al. An introduction to matrix concentration inequalities. *Foundations and Trends® in Machine Learning*, 8(1-2):1–230, 2015.
- [37] 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.
- [38] Rahul Yedida, Snehanthu Saha, and Tejas Prashanth. Lipschitzlr: Using theoretically computed adaptive learning rates for fast convergence. *Applied Intelligence*, 51:1460–1478, 2021.

<b>A</b>	<b>Concentration Inequalities</b>	<b>21</b>
A.1	Linear Kernel . . . . .	21
A.2	Polynomial Kernel . . . . .	25
A.3	Gaussian Kernel . . . . .	25
<b>B</b>	<b>Proofs for Section 3</b>	<b>27</b>
B.1	Proof of Lemma 13 . . . . .	27
B.2	Proof of Lemma 15 . . . . .	28
B.3	Proof of Lemma 17 . . . . .	28
B.4	Proof of Lemma 20 . . . . .	28
B.5	Proof of Theorem 21 . . . . .	29
<b>C</b>	<b>Proofs for Section 4</b>	<b>31</b>
C.1	Proof of Theorem 26 . . . . .	31
C.2	Proof of Theorem 27 . . . . .	31
<b>D</b>	<b>Necessary Lemmas</b>	<b>32</b>
<b>E</b>	<b>Experimental Details</b>	<b>33</b>
E.1	Kernel Regression . . . . .	33
E.2	Kernel Binary Classification . . . . .	33
E.3	Kernel Multi-Class Classification . . . . .	33
E.4	Linear Regression . . . . .	33
<b>F</b>	<b>Detailed Related Works</b>	<b>34</b>
F.1	High-dimensional Robust Mean Estimation via Gradient Descent <a href="#">[3]</a> . . . . .	34
F.2	Trimmed Maximum Likelihood Estimation for Robust Generalized Linear Model <a href="#">[1]</a> . . . . .	34

## A Concentration Inequalities

In this section we will give various concentration inequalities on the inlier data. We first restate our assumptions.

**Assumption 28.** (Sub-Gaussian Design of Covariates). We assume each covariate is drawn i.i.d from a zero-mean covariance  $\Sigma$  sub-Gaussian distribution with sub-Gaussian norm  $K_1 \in \mathbb{R}_{++}$ .

$$\mathbb{E}[\mathbf{x}_i] = \mathbf{0} \quad (62)$$

$$\mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top] = \Sigma \quad (63)$$

For all  $i \in [n]$ ,

$$\mathbb{P}\{\mathbf{v}^\top \mathbf{x}_i \geq t\} \leq \exp\left(-\frac{t^2}{K_1^2}\right) \quad (64)$$

Is this equivalent to

$$\mathbb{P}\{\mathbf{v}^\top \mathbf{x}_i \geq t\} \leq \exp\left(-\frac{t^2}{K_1^2}\right) \quad (65)$$

**Assumption 29.** (Sub-Gaussian Design of Optimal Residuals). Recall the residual is defined as  $\eta_i \triangleq \mathbf{f}_{\mathbf{w}}^*(\mathbf{x}_i) - y_i$ . Then we assume for some  $K_2 \in \mathbb{R}_{++}$

$$\mathbb{E}[\eta_i] = 0 \quad (66)$$

$$\mathbb{P}\{|\eta_i| \geq t\} \leq 2 \exp\left(-\frac{t^2}{K_2^2}\right) \quad (67)$$

### A.1 Linear Kernel

We will first begin with the Linear Kernel Case.

**Assumption 30.** Let the data have dimension  $d$  such that Assumption 28 holds. Then assume

$$n > \frac{1}{18} \left( \sqrt{(18 - 24dK)^2 - 36d^2 K^2 \log(2d)} - (18 - 24dK) \log(2d) \right) = \Omega(dK \log(2d)) \quad (68)$$

**Lemma 31.** Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  have sub-Gaussian design described in Assumption 28 and  $\eta_1, \dots, \eta_n$  have sub-Gaussian design described in Assumption 29. It then follows

$$\mathbb{E} \left\| \sum_{i=1}^n \eta_i \mathbf{x}_i \right\| \leq \frac{1}{2} K_1 K_2 e^{1/4} \sqrt{8\pi n d} = \mathcal{O}(K_1 K_2 \sqrt{nd}) \quad (69)$$

and with probability at least 0.99 we have

$$\left\| \sum_{i=1}^n \eta_i \mathbf{x}_i \right\|_2 \leq 7K_1 K_2 \sqrt{nd} \quad (70)$$

**Proof.** We will first upper bound the event  $\|\sum_{i=1}^n \eta_i \mathbf{x}_i\| \geq t$ . To this we utilize the vector Bernstein Inequality.

$$\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^d x_i^2} \leq \sqrt{dK_1^2} \leq \sqrt{d} K_1 \quad (71)$$

$$\mathbb{E}[\|\eta_i \mathbf{x}_i\|^2] = \mathbb{E}[\eta_i^2 \|\mathbf{x}_i\|^2] = \mathbb{E}[\eta_i^2] \mathbb{E}\left[\sum_{i=1}^d x_i^2\right] \leq dK_1^2 K_2^2 \quad (72)$$

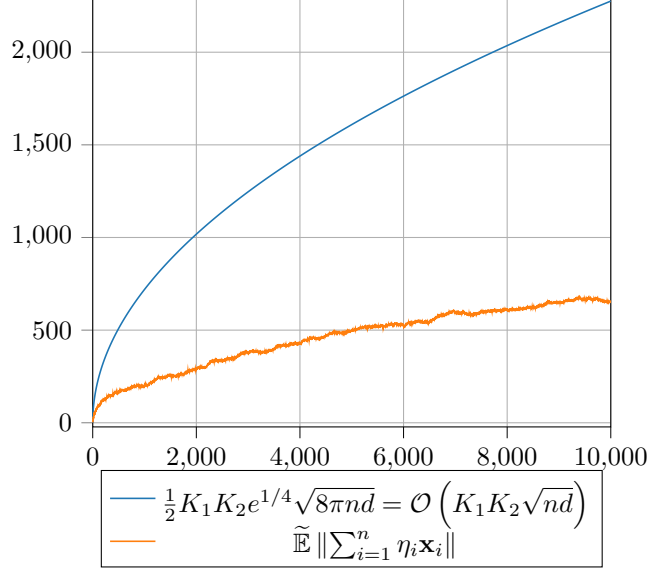


Figure 3: The bound for ?? compared to an average over 10 trials. We set  $\Sigma = \mathbf{I}$ ,  $K = 1$ ,  $d = 20$ .

$$\mathbb{E} \left[ \left\| \sum_{i=1}^n \eta_i \mathbf{x}_i \right\| \right] = \int_0^\infty \mathbb{P} \left\{ \left\| \sum_{i=1}^n \eta_i \mathbf{x}_i \right\| > t \right\} dt \quad (73)$$

$$\leq \int_0^\infty \exp \left( -\frac{t^2}{8ndK_1^2K_2^2} + \frac{1}{4} \right) dt \quad (74)$$

$$= \frac{1}{2} K_1 K_2 e^{1/4} \sqrt{8\pi nd} \quad (75)$$

In probability, we apply the Vector Bernstein Inequality and obtain with probability at least  $1 - \delta$

$$\left\| \sum_{i=1}^n \eta_i \mathbf{x}_i \right\|_2 \leq K_1 K_2 \sqrt{nd} + 2K_1 K_2 \sqrt{nd \ln \left( \frac{1}{\delta} \right)} \quad (76)$$

This completes the proof  $\blacksquare$

**Lemma 32.** Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  have sub-Gaussian design described in Assumption 28 such that  $n > d \log(d)$ . It then follows

$$\mathbb{E} \left[ \sigma_{\min} \left( \sum_{k \in [n]} \mathbf{x}_k \mathbf{x}_k^\top \right) \right] \geq \sigma_{\min}(\Sigma) \left( n - \left( \sqrt{2(dK-1)n \log(2d)} + \frac{1}{3} dK \log(2d) \right) \right) = \Omega(\sigma_{\min}(\Sigma) n) \quad (77)$$

and with probability  $(1 - \delta)$

$$\sigma_{\min} \left( \sum_{k \in [n]} \mathbf{x}_k \mathbf{x}_k^\top \right) \geq \frac{1}{1-\delta} \sigma_{\min}(\Sigma) \left( n - \left( \sqrt{2(dK-1)n \log(2d)} + \frac{1}{3} dK \log(2d) \right) \right) \quad (78)$$

**Proof.** In expectation we have

$$\mu_{\min} \triangleq \sigma_{\min} \left( \sum_{i=1}^n \mathbb{E} [\mathbf{x}_i \mathbf{x}_i^\top] \right) = \sigma_{\min} \left( \sum_{i=1}^n \Sigma \right) = n \sigma_{\min}(\Sigma) \quad (79)$$

Then we have by the sub-Gaussian design in Assumption 28.

$$\sigma_{\min} \left( \sum_{k \in [n]} \mathbf{x}_k \mathbf{x}_k^\top \right) = \sigma_{\min} \left( n\Sigma + \sum_{k \in [n]} \mathbf{x}_k \mathbf{x}_k^\top - n\Sigma \right) \quad (80)$$

$$\stackrel{\text{Weyl's}}{\geq} \sigma_{\min}(n\mathbf{\Sigma}) - \sigma_{\max}\left(\sum_{k \in [n]} \mathbf{x}_k \mathbf{x}_k^\top - n\mathbf{\Sigma}\right) \quad (81)$$

$$= \mu_{\min} - \underbrace{\left\| \sum_{k \in [n]} \mathbf{x}_k \mathbf{x}_k^\top - \mathbf{\Sigma} \right\|_2}_A \quad (82)$$

We assume  $\mathbf{\Sigma} \succcurlyeq \mathbf{0}$  and symmetric, therefore we can represent  $\mathbf{\Sigma} = (\mathbf{\Sigma}^{1/2})^\top \mathbf{\Sigma}^{1/2} = \mathbf{\Sigma}^{1/2}(\mathbf{\Sigma}^{1/2})^\top$ . Next we will upper bound A.

$$A \triangleq \left\| \sum_{k \in [n]} \mathbf{x}_k \mathbf{x}_k^\top - \mathbf{\Sigma} \right\|_2 \quad (83)$$

$$= \left\| \sum_{k \in [n]} \left( \mathbf{\Sigma}^{1/2} \hat{\mathbf{x}}_k \right) \left( \mathbf{\Sigma}^{1/2} \hat{\mathbf{x}}_k \right)^\top - \mathbf{\Sigma} \right\|_2 \quad (84)$$

$$= \left\| \sum_{k \in [n]} \mathbf{\Sigma}^{1/2} \hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top \left( \mathbf{\Sigma}^{1/2} \right)^\top - \mathbf{\Sigma} \right\|_2 \quad (85)$$

$$= \left\| \mathbf{\Sigma}^{1/2} \left( \sum_{k \in [n]} \mathbf{x}_k \mathbf{x}_k^\top - \mathbf{I} \right) \left( \mathbf{\Sigma}^{1/2} \right)^\top \right\|_2 \quad (86)$$

$$\leq \underbrace{\left\| \sum_{k \in [n]} \hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top - \mathbf{I} \right\|}_B \|\mathbf{\Sigma}\| \quad (87)$$

It then suffices to upper bound B. First, let  $\hat{\mathbf{x}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and define  $K \geq 1$  s.t.  $\hat{\mathbf{x}} \leq dK$ , then

$$\mathbb{V}(\hat{\mathbf{x}}\hat{\mathbf{x}}^\top - \mathbf{I}) = \mathbb{E} \left[ (\hat{\mathbf{x}}\hat{\mathbf{x}}^\top - \mathbf{I})^\top (\hat{\mathbf{x}}\hat{\mathbf{x}}^\top - \mathbf{I}) \right] \quad (88)$$

$$= \mathbb{E} [\hat{\mathbf{x}}\hat{\mathbf{x}}^\top \hat{\mathbf{x}}\hat{\mathbf{x}}^\top] - 2\mathbb{E} [\hat{\mathbf{x}}\hat{\mathbf{x}}^\top] + \mathbf{I} \quad (89)$$

$$= \mathbb{E} [\hat{\mathbf{x}}\hat{\mathbf{x}}^\top \hat{\mathbf{x}}\hat{\mathbf{x}}^\top] - \mathbf{I} \quad (90)$$

$$= \mathbb{E} [\|\hat{\mathbf{x}}\|^2 \hat{\mathbf{x}}\hat{\mathbf{x}}^\top] - \mathbf{I} \quad (91)$$

$$\leq dK \mathbb{E} [\hat{\mathbf{x}}\hat{\mathbf{x}}^\top] - \mathbf{I} \quad (92)$$

$$= (dK - 1) \mathbf{I} \quad (93)$$

Then from Matrix Bernstein, we have

$$\mathbb{E}[B] \leq \sqrt{2(dK - 1)n \log(2d)} + \frac{1}{3}dK \log(2d) \quad (94)$$

Then we have in expectation

$$\mathbb{E} \left[ \sigma_{\min} \left( \sum_{k \in [n]} \mathbf{x}_k \mathbf{x}_k^\top \right) \right] \geq n \sigma_{\min}(\mathbf{\Sigma}) - \|\mathbf{\Sigma}\| \left( \sqrt{2(dK - 1)n \log(2d)} + \frac{1}{3}dK \log(2d) \right) \quad (95)$$

Assume the spectrum of  $\mathbf{\Sigma}$  is flat, then we have

$$\mathbb{E} \left[ \sigma_{\min} \left( \sum_{k \in [n]} \mathbf{x}_k \mathbf{x}_k^\top \right) \right] \geq \sigma_{\min}(\mathbf{\Sigma}) \left( n - \left( \sqrt{2(dK - 1)n \log(2d)} + \frac{1}{3}dK \log(2d) \right) \right) \quad (96)$$

Similarly from the Matrix Bernstein Inequality in [36], we have

$$\mathbb{P} \left\{ \sigma_{\min} \left( \sum_{k \in [n]} \mathbf{x}_k \mathbf{x}_k^\top \right) > t \right\} \leq 2d \exp \left( \frac{-t^2}{2n(dK-1) + Kt/3} \right) \quad (97)$$

We also like to note if the singular values have linear decay rate th ■

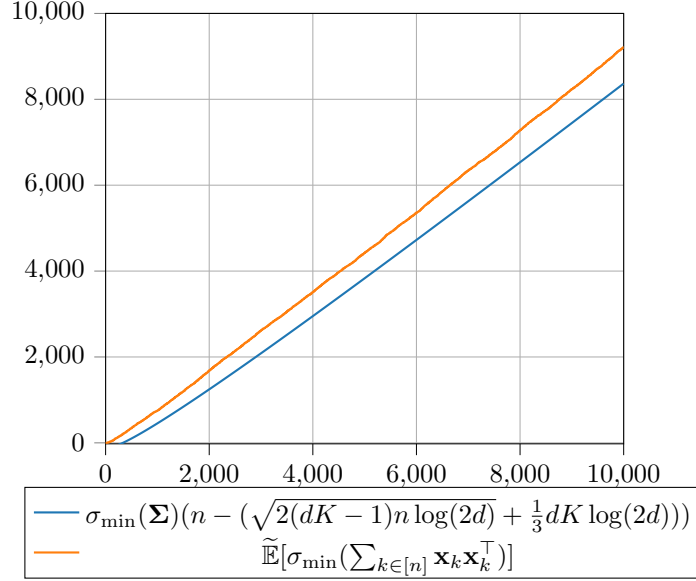


Figure 4: The bound for Lemma 32 compared to an average over 10 trials. We set  $\Sigma = \mathbf{I}$ ,  $K = 4$ ,  $d = 20$ .

**Lemma 33.** Let  $\eta_i \in P \setminus S$  be defined as in Assumption 29, then it follows

$$\mathbb{E} \left[ \sum_{i \in P \setminus S} \eta_i^2 \right] \leq eK_3^2 n(1 - 2\varepsilon) \quad (98)$$

**Proof.** Note  $|P \setminus S| = (1 - 2\varepsilon)n$ , then

$$\mathbb{P} \left\{ \sum_{i \in P \setminus S} \eta_i^2 \geq t \right\} = \mathbb{P} \left\{ \sum_{i \in P \setminus S} (\eta_i^2 - \mathbb{V}(\eta_i)) + \mathbb{V}(\eta_i) \geq t \right\} \quad (99)$$

$$= \mathbb{P} \left\{ \sum_{i \in P \setminus S} (\eta_i^2 - \mathbb{E}[\eta_i^2]) \geq t - n(1 - 2\varepsilon)\mathbb{E}[\eta_i^2] \right\} \quad (100)$$

$$\stackrel{\text{Bernstein}}{\leq} \exp \left( - \frac{(t - n(1 - 2\varepsilon)\mathbb{E}[\eta_i^2])^2}{2 \left( \sum_{i \in P \setminus S} \mathbb{E}[(\eta_i^2 - \mathbb{E}[\eta_i^2])^2] + \frac{1}{3}Mt \right)} \right) \quad (101)$$

**Alternative Proof.** We will utilize Chernoff's Inequality.

$$\mathbb{P} \left\{ \sum_{i \in P \setminus S} \eta_i^2 \geq t \right\} \stackrel{\zeta_1}{\leq} \mathbb{P} \left\{ \exp \left( \lambda \sum_{i \in P \setminus S} \eta_i^2 \right) \geq \exp(\lambda t) \right\} \quad (102)$$

$$\stackrel{\text{Markov}}{\leq} \exp(-\lambda t) \mathbb{E} \left[ \exp \left( \lambda \sum_{i \in P \setminus S} \eta_i^2 \right) \right] \quad (103)$$



$$= \exp(-\lambda t) \mathbb{E} \left[ \prod_{i \in P \setminus S} \exp(\lambda \eta_i^2) \right] \quad (104)$$

$$= \exp(-\lambda t) \prod_{i \in P \setminus S} \mathbb{E} [\exp(\lambda \eta_i^2)] \quad (105)$$

$$\stackrel{\zeta_2}{\leq} \exp(-\lambda t) \prod_{i \in P \setminus S} \exp(\lambda K_3^2) \quad (106)$$

$$= \exp(-\lambda t + n(1 - 2\varepsilon)\lambda K_3^2) \quad (107)$$

$$\stackrel{\zeta_3}{=} \exp\left(-\frac{t}{K_3^2 n(1 - 2\varepsilon)} + 1\right) \quad (108)$$

$(\zeta_1)$  holds for any  $\lambda > 0$ .  $(\zeta_2)$  holds for any  $\lambda \leq K_3$ . Then we can calculate the expectation by the following,

$$\mathbb{E} \left[ \sum_{i \in P \setminus S} \eta_i^2 \geq t \right] = \int_0^\infty \mathbb{P} \left\{ \sum_{i \in P \setminus S} \eta_i^2 \geq t \right\} dt \quad (109)$$

$$\leq \int_0^\infty \exp\left(-\frac{t}{K_3^2 n(1 - 2\varepsilon)} + 1\right) dt \quad (110)$$

$$= e K_3^2 n(1 - 2\varepsilon) \quad (111)$$

This concludes the proof. ■

## A.2 Polynomial Kernel

The polynomial kernel is given by

$$K(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^\top \mathbf{x}_2 + C)^p \quad (112)$$

The feature map for the polynomial kernel is given as

$$\phi_{\text{poly}}(\mathbf{x}) = [x_1, \dots, x_d, x_1^2, \dots, x_d^2, \dots, x_1^p, \dots, x_d^p, x_1 x_2, \dots, x_{d-1} x_d] \in \mathbb{R}^{d^p} \quad (113)$$

**Lemma 34.** *Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  have sub-Gaussian design described in Assumption 28 and  $\eta_1, \dots, \eta_n$  have sub-Gaussian design described in Assumption 29. It then follows*

$$\mathbb{E} \left\| \sum_{i=1}^n \eta_i \phi_{\text{poly}}(\mathbf{x}_i) \right\| \leq \Xi \quad (114)$$

**Lemma 35.** *Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  have sub-Gaussian design described in Assumption 28 such that  $n > d \log(d)$ . It then follows*

$$\sigma_{\min} \left( \sum_{i=1}^n \phi_{\text{poly}}(\mathbf{x}_i) \phi_{\text{poly}}(\mathbf{x}_i)^\top \right) \geq \Xi \quad (115)$$

## A.3 Gaussian Kernel

The gaussian kernel is given for  $\gamma > 0$  by

$$K(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2\right) \quad (116)$$

**Lemma 36.** *Let  $\eta_1, \dots, \eta_n \sim \mathcal{N}(0, \sigma^2)$  be sub-Gaussian. Then it follows*

$$\mathbb{E} \left| \sum_{i=1}^n \eta_i \right| \leq \Xi \quad (117)$$

We will utilize the Random Fourier Features (RFF) representation [29]. Let  $\mathbf{w}_1, \dots, \mathbf{w}_d \sim \mathcal{N}_d(\mathbf{0}, \frac{2}{\gamma} \mathbf{I})$ , then the RFF representation is given by

$$\phi_{\text{RFF}}(\mathbf{x}) = \frac{1}{\sqrt{d}} [\cos(\mathbf{w}_1^\top \mathbf{x}), \sin(\mathbf{w}_1^\top \mathbf{x}), \dots, \cos(\mathbf{w}_d^\top \mathbf{x}), \sin(\mathbf{w}_d^\top \mathbf{x})] \quad (118)$$

The key idea behind the Randomized Fourier Features is

$$\mathbb{E} [\phi_{\text{RFF}}(\mathbf{x}_1)^\top \phi_{\text{RFF}}(\mathbf{x}_2)] = \exp \left( -\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 \right) \quad (119)$$

We will first derive a relation between  $\mathbb{E} [\phi_{\text{RFF}}(\mathbf{x}_1)^\top \phi_{\text{RFF}}(\mathbf{x}_2)]$  and  $\mathbb{E} [\phi_{\text{RFF}}(\mathbf{x}_1)]^\top \mathbb{E} [\phi_{\text{RFF}}(\mathbf{x}_2)]$ .

## B Proofs for Section 3

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

### B.1 Proof of Lemma 13

**Proof.** For any  $\mathbf{f}_{\mathbf{w}_1}, \mathbf{f}_{\mathbf{w}_2} \in \mathcal{K}$ , we want to make sure the gradient is bounded.

$$|g(t, \mathbf{f}_{\mathbf{w}_1}) - g(t, \mathbf{f}_{\mathbf{w}_2})| = \left| \int_0^1 \nabla_{\mathbf{f}} g(t, (1-\lambda)\mathbf{f}_{\mathbf{w}_1} + \lambda\mathbf{f}_{\mathbf{w}_2})(\mathbf{f}_{\mathbf{w}_1} - \mathbf{f}_{\mathbf{w}_2}) d\lambda \right| \quad (120)$$

$$\leq \|\mathbf{f}_{\mathbf{w}_1} - \mathbf{f}_{\mathbf{w}_2}\|_{\mathcal{H}} \left| \int_0^1 \nabla_{\mathbf{f}} g(t, (1-\lambda)\mathbf{f}_{\mathbf{w}_1} + \lambda\mathbf{f}_{\mathbf{w}_2}) d\lambda \right| \quad (121)$$

$$\stackrel{(a)}{\leq} \|\mathbf{f}_{\mathbf{w}_1} - \mathbf{f}_{\mathbf{w}_2}\|_{\mathcal{H}} \max_{\mathbf{w} \in \mathcal{K}} \|\nabla_{\mathbf{f}} g(t, \mathbf{f}_{\mathbf{w}})\|_{\mathcal{H}} \quad (122)$$

In (a), we note that since  $\mathcal{K}$  is convex, then by definition as  $\mathbf{f}_{\mathbf{w}_1}, \mathbf{f}_{\mathbf{w}_2} \in \mathcal{K}$ , we have for  $\lambda \in [0, 1]$ , the convex combination  $(1-\lambda)\mathbf{f}_{\mathbf{w}_1} + \lambda\mathbf{f}_{\mathbf{w}_2} \in \mathcal{K}$ . We use the  $\mathcal{H}$  norm of the gradient to bound  $L$  from above for an element in the convex closed set  $\mathcal{K}$ .

$$\|\nabla g(t, \mathbf{f}_{\mathbf{w}})\|_{\mathcal{H}} = \left\| \frac{2}{np} \sum_{i=1}^n \mathbb{I}\{t \geq (\mathbf{f}_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2\} (\mathbf{f}_{\mathbf{w}}(\mathbf{x}_i) - y_i) \cdot K(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}} \quad (123)$$

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  $t \geq (\mathbf{f}_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2$ .

$$= \left\| \frac{2}{np} \sum_{i=1}^m (\mathbf{f}_{\mathbf{w}}(\mathbf{x}_i) - y_i) \cdot K(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}} \quad (124)$$

$$\leq \frac{2}{np} \left( \left\| \sum_{i=1}^m \mathbf{f}_{\mathbf{w}}(\mathbf{x}_i) \cdot K(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}} + \left\| \sum_{i=1}^m y_i K(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}} \right) \quad (125)$$

$$\leq \frac{2}{np} \left( \left\| \sum_{i=1}^m \left( \sum_{j=1}^n w_j K(\mathbf{x}_i, \mathbf{x}_j) \right) K(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}} + \left\| \sum_{i=1}^m y_i \left\| \sum_{i=1}^m K(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}} \right\|_{\mathcal{H}} \right) \quad (126)$$

$$\stackrel{(a)}{=} \frac{2}{np} \left( \left\| \sum_{i=1}^m \left\langle \sum_{j=1}^n w_j K(\mathbf{x}_j, \cdot), K(\mathbf{x}_i, \cdot) \right\rangle_{\mathcal{H}} \cdot K(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}} + \left\| \sum_{i=1}^m y_i \left\| \sum_{i=1}^m \sqrt{K(\mathbf{x}_i, \mathbf{x}_i)} \right\|_{\mathcal{H}} \right) \quad (127)$$

$$\leq \frac{2}{np} \left( \left\| \left\langle \sum_{j=1}^n w_j K(\mathbf{x}_j, \cdot), \sum_{i=1}^m K(\mathbf{x}_i, \cdot) \right\rangle_{\mathcal{H}} \right\|_{\mathcal{H}} \left\| \sum_{i=1}^m K(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}} + \left\| \sum_{i=1}^m y_i \left\| \sum_{i=1}^m \sqrt{K(\mathbf{x}_i, \mathbf{x}_i)} \right\|_{\mathcal{H}} \right) \quad (128)$$

$$\leq \frac{2}{np} \left( \left\| \sum_{j=1}^n w_j K(\mathbf{x}_j, \cdot) \right\|_{\mathcal{H}} \left( \sum_{i=1}^m \sqrt{K(\mathbf{x}_i, \mathbf{x}_i)} \right)^2 + \left\| \sum_{i=1}^n y_i \left\| \sum_{i=1}^m \sqrt{K(\mathbf{x}_i, \mathbf{x}_i)} \right\|_{\mathcal{H}} \right) \quad (129)$$

$$\leq \frac{2}{np} \left( \|\mathbf{f}_{\mathbf{w}}\|_{\mathcal{H}} \left( \sum_{i=1}^m \sqrt{K(\mathbf{x}_i, \mathbf{x}_i)} \right)^2 + \|\mathbf{y}\|_1 \left( \sum_{i=1}^n \sqrt{K(\mathbf{x}_i, \mathbf{x}_i)} \right) \right) \quad (130)$$

$$\leq \frac{2R}{np} \left( \sum_{i=1}^n \sqrt{K(\mathbf{x}_i, \mathbf{x}_i)} \right)^2 + \frac{2\|\mathbf{y}\|_1}{np} \left( \sum_{i=1}^n \sqrt{K(\mathbf{x}_i, \mathbf{x}_i)} \right) \quad (131)$$

(a) follows from the reproducing property for RKHS [11]. If we have a normalized kernel such as the Gaussian Kernel, then we have the Lipschitz Constant is finite. Furthermore, if the adversary introduces label corruption that tends to  $\infty$ , then these points will not be in the Subquantile as  $\mathbf{f}_{\mathbf{w}}$  has bounded norm, so it will have infinite error. This concludes the proof.  $\blacksquare$

## B.2 Proof of Lemma 15

**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  $\sigma(x) = \frac{1}{1+e^{-x}}$ .

$$\|\nabla_{\mathbf{f}} g(t, \mathbf{f}_{\mathbf{w}})\|_{\mathcal{H}} = \left\| \frac{1}{np} \sum_{i=1}^n \mathbb{I}\{t \geq (1 - y_i) \log(\mathbf{f}_{\mathbf{w}}(\mathbf{x}_i))\} (y_i - \sigma(\mathbf{f}_{\mathbf{w}}(\mathbf{x}_i))) \cdot K(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}} \quad (132)$$

$$\leq \frac{1}{np} \left\| \sum_{i \in S} (y_i - \sigma(\mathbf{f}_{\mathbf{w}}(\mathbf{x}_i))) \right\| \left\| \sum_{i \in S} K(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}} \quad (133)$$

$$\stackrel{(a)}{\leq} \sum_{i \in S} \sqrt{K(\mathbf{x}_i, \mathbf{x}_i)} \quad (134)$$

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

## B.3 Proof of Lemma 17

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

$$\|\nabla_{\mathbf{W}} g(t, \mathbf{W})\|_2 = \quad (135)$$

■

## B.4 Proof of Lemma 20

**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  $\eta_i \triangleq (\mathbf{f}_{\mathbf{w}^*}(\mathbf{x}_i) - y_i)$  where  $i \in P$ .

$$\Phi(\mathbf{f}_{\mathbf{w}}) - \Phi(\mathbf{f}_{\mathbf{w}^*}) = \sum_{i \in S} (\mathbf{f}_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 - \sum_{j \in P} (\mathbf{f}_{\mathbf{w}^*}(\mathbf{x}_j) - y_j)^2 \quad (136)$$

$$= \sum_{i \in S \cap P} (\mathbf{f}_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 + \sum_{i \in S \cap Q} (\mathbf{f}_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 - \sum_{j \in P} (\mathbf{f}_{\mathbf{w}^*}(\mathbf{x}_j) - y_j)^2 \quad (137)$$

$$\geq \sum_{i \in S \cap P} (\mathbf{f}_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 - \sum_{j \in P} (\mathbf{f}_{\mathbf{w}^*}(\mathbf{x}_j) - y_j)^2 \quad (138)$$

$$= \sum_{i \in S \cap P} (\mathbf{f}_{\mathbf{w}}(\mathbf{x}_i) - \mathbf{f}_{\mathbf{w}^*}(\mathbf{x}_i) - \eta_i)^2 - \sum_{j \in P} \eta_j^2 \quad (139)$$

$$= \sum_{i \in S \cap P} ((\mathbf{f}_{\mathbf{w}} - \mathbf{f}_{\mathbf{w}^*})(\mathbf{x}_i) - \eta_i)^2 - \sum_{j \in P} \eta_j^2 \quad (140)$$

$$\geq \sum_{i \in S \cap P} \underbrace{((\mathbf{f}_{\mathbf{w}} - \mathbf{f}_{\mathbf{w}^*})(\mathbf{x}_i))^2}_{A_1} - 2 \underbrace{\sum_{i \in S \cap P} \eta_i (\mathbf{f}_{\mathbf{w}} - \mathbf{f}_{\mathbf{w}^*})(\mathbf{x}_i)}_{A_2} - \underbrace{\sum_{j \in P \setminus S} \eta_j^2}_{A_3} \quad (141)$$

Now we will upper bound  $A_1$ .

$$\sum_{i \in S \cap P} ((\mathbf{f}_{\mathbf{w}} - \mathbf{f}_{\mathbf{w}^*})(\mathbf{x}_i))^2 \stackrel{(a)}{=} \sum_{i \in S \cap P} \left\langle \sum_{j \in X} (w_j - w_j^*) K(\mathbf{x}_j, \cdot), K(\mathbf{x}_i, \cdot) \right\rangle_{\mathcal{H}}^2 \quad (142)$$

$$= \sum_{i \in S \cap P} \left\langle \sum_{j \in X} (w_j - w_j^*) \phi(\mathbf{x}_j), \phi(\mathbf{x}_i) \right\rangle \left\langle \phi(\mathbf{x}_i), \sum_{j \in X} (w_j - w_j^*) \phi(\mathbf{x}_j) \right\rangle \quad (143)$$

$$\geq \|\mathbf{f}_{\mathbf{w}} - \mathbf{f}_{\mathbf{w}^*}\|_{\mathcal{H}}^2 \underbrace{\sigma_{\min} \left( \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top \right)}_{B_1} \quad (144)$$

Next we will upper bound  $A_2$ ,

$$A_2 = \sum_{i \in S \cap P} \eta_i (\mathbf{f}_{\mathbf{w}} - \mathbf{f}_{\mathbf{w}^*})(\mathbf{x}_i) \quad (145)$$

$$= \sum_{i \in S \cap P} \left\langle \sum_{j \in X} (w_j - w_j^*) K(\mathbf{x}_j, \cdot), \eta_i K(\mathbf{x}_i, \cdot) \right\rangle_{\mathcal{H}} \quad (146)$$

$$= \left\langle \sum_{j \in X} (w_j - w_j^*) K(\mathbf{x}_j, \cdot), \sum_{i \in S \cap P} \eta_i K(\mathbf{x}_i, \cdot) \right\rangle_{\mathcal{H}} \quad (147)$$

$$\leq \|\mathbf{f}_{\mathbf{w}} - \mathbf{f}_{\mathbf{w}^*}\|_{\mathcal{H}} \left\| \sum_{i \in S \cap P} \eta_i K(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}} = \|\mathbf{f}_{\mathbf{w}} - \mathbf{f}_{\mathbf{w}^*}\|_{\mathcal{H}} \left\| \sum_{i \in S \cap P} \eta_i \phi(\mathbf{x}_i) \right\|_2 \quad (148)$$

We will now lower bound  $B_1$ . For the linear kernel, for  $B_1$  to be greater than 0, than if  $\mathbf{x} \in \mathbb{R}^d$ , then we must have  $\frac{d}{1-2\varepsilon} < n$ , otherwise we will have a rank-deficient matrix which will thus have singular values of value 0. For the polynomial kernel, for  $B_1$  to be greater than 0, than  $n > d^r$  where  $r$  is the polynomial degree. We thus have

$$\Phi(\mathbf{f}_{\mathbf{w}}) - \Phi(\mathbf{f}_{\mathbf{w}^*}) \geq \eta^2 \sigma_{\min} \left( \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top \right) - 2\eta \left\| \sum_{i \in S \cap P} \eta_i \phi(\mathbf{x}_i) \right\| - \sum_{j \in P \setminus S} \eta_j^2 \quad (149)$$

This completes the proof. ■

## B.5 Proof of Theorem 21

**Proof.** First,

$$\|\nabla \Phi_\lambda(\mathbf{f}_{\mathbf{w}})\|_{\mathcal{H}} = \left\| \frac{1}{\lambda} \left( \mathbf{f}_{\mathbf{w}} - \arg \min_{\mathbf{f}_{\tilde{\mathbf{w}}} \in \mathcal{K}} \left( \Phi(\mathbf{f}_{\tilde{\mathbf{w}}}) + \frac{1}{2\lambda} \|\mathbf{f}_{\mathbf{w}} - \mathbf{f}_{\tilde{\mathbf{w}}}\|_{\mathcal{H}}^2 \right) \right) \right\|_{\mathcal{H}} = 0 \quad (150)$$

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

$$\Phi(\mathbf{f}_{\tilde{\mathbf{w}}}) < \Phi(\mathbf{f}_{\tilde{\mathbf{w}}}) + \frac{1}{2\lambda} \|\mathbf{f}_{\tilde{\mathbf{w}}} - \mathbf{f}_{\tilde{\mathbf{w}}}\|_{\mathcal{H}}^2 \quad (151)$$

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

$$\Phi(\mathbf{f}_{\tilde{\mathbf{w}}}) - \Phi(\mathbf{f}_{\mathbf{w}^*}) \leq \frac{1}{2\lambda} \|\mathbf{f}_{\tilde{\mathbf{w}}} - \mathbf{f}_{\mathbf{w}^*}\|_{\mathcal{H}}^2 \quad (152)$$

We proceed by proof by contradiction. Assume  $\|\mathbf{f}_{\tilde{\mathbf{w}}} - \mathbf{w}^*\| > \eta$ , then if  $\Phi(\mathbf{f}_{\tilde{\mathbf{w}}}) - \Phi(\mathbf{w}^*) > \frac{1}{2}\eta^2$ , then we will have  $\mathbf{f}_{\tilde{\mathbf{w}}}$  is not a stationary point, which will imply  $\|\mathbf{f}_{\tilde{\mathbf{w}}} - \mathbf{w}^*\|_{\mathcal{H}} \leq \eta$ . Therefore, we attempt to find the minimum value for  $\eta$ . From Lemma 20, we have we have

$$\Phi(\mathbf{w}) - \Phi(\mathbf{w}^*) \geq \eta^2 \sigma_{\min} \left( \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top \right) - 2\eta \left\| \sum_{i \in S \cap P} \eta_i \phi(\mathbf{x}_i) \right\| - \sum_{j \in P \setminus S} \eta_j^2 \quad (153)$$

From the definition of stationary point, we have

$$\eta^2 \left( \sigma_{\min} \left( \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top \right) - \frac{1}{2}L \right) - 2\eta \left\| \sum_{i \in S \cap P} \eta_i \phi(\mathbf{x}_i) \right\| - \sum_{j \in P \setminus S} \eta_j^2 \geq 0 \quad (154)$$

From lower bounding the positive solution of the quadratic equation, we have

$$\eta \geq \Omega \left( \frac{2L \left( \sum_{j \in P \setminus S} \eta_j^2 \right)}{2L\sigma_{\min} \left( \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top \right) - 1} \right)^{1/2} + \frac{4L \left\| \sum_{i \in S \cap P} \eta_i \phi(\mathbf{x}_i) \right\|}{\sigma_{\min} \left( \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top \right)} \quad (155)$$

Therefore, when ?? holds, we have a contradiction. It thus follows

$$\eta \leq \left( \frac{2L \left( \sum_{j \in P \setminus S} \eta_j^2 \right)}{2L\sigma_{\min} \left( \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top \right) - 1} \right)^{1/2} + \frac{4L \left\| \sum_{i \in S \cap P} \eta_i \phi(\mathbf{x}_i) \right\|}{\sigma_{\min} \left( \sum_{i \in S \cap P} \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top \right)} \quad (156)$$

This completes the proof. ■

## C Proofs for Section 4

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

### C.1 Proof of Theorem 26

We will start with the definition of the Moreau Envelope.

$$\bar{\Phi}_\lambda(\mathbf{f}_\mathbf{w}) = \frac{1}{\lambda} \left( \mathbf{f}_\mathbf{w} - \arg \min_{\mathbf{f}_{\hat{\mathbf{w}}} \in \mathcal{K}} \left\{ \Phi(\mathbf{f}_{\hat{\mathbf{w}}}) + \frac{1}{2\lambda} \|\mathbf{f}_\mathbf{w} - \mathbf{f}_{\hat{\mathbf{w}}}\|_{\mathcal{H}}^2 \right\} \right) \quad (157)$$

Note  $g(t, \mathbf{f}_\mathbf{w})$  is  $L$ -lipschitz in  $\mathbf{f}_\mathbf{w}$ . Let  $\mathbf{f}_{\hat{\mathbf{w}}}^{(t)} = \arg \min_{\mathbf{f}_{\hat{\mathbf{w}}} \in \mathcal{K}} \{\bar{\Phi}(\mathbf{f}_\mathbf{w}) + \|\mathbf{f}_{\hat{\mathbf{w}}} - \mathbf{f}_\mathbf{w}\|_{\mathcal{H}}^2\}$ . Then we have,

$$\bar{\Phi}_\lambda(\mathbf{f}_{\hat{\mathbf{w}}}^{(t+1)}) \leq \Phi(\mathbf{f}_{\hat{\mathbf{w}}}^{(t)}) + \left\| \mathbf{f}_{\hat{\mathbf{w}}}^{(t+1)} - \mathbf{f}_{\hat{\mathbf{w}}}^{(t)} \right\|_{\mathcal{H}}^2 \quad (158)$$

$$= \Phi(\mathbf{f}_{\hat{\mathbf{w}}}^{(t)}) + \left\| \mathbf{f}_{\hat{\mathbf{w}}}^{(t)} - \Pi_{\mathcal{K}} \left( \mathbf{f}_{\hat{\mathbf{w}}}^{(t)} - \alpha \left( \mu \mathbf{b}^{(t)} \right) \right) \right\|_{\mathcal{H}}^2 \quad (159)$$

$$= \Phi(\mathbf{f}_{\hat{\mathbf{w}}}^{(t)}) + \left\| \mathbf{f}_{\hat{\mathbf{w}}}^{(t)} - \Pi_{\mathcal{K}} \left( \mathbf{f}_{\hat{\mathbf{w}}}^{(t)} - \alpha \left( \mu \left( \mu \mathbf{b}^{(t-1)} + \nabla_{\mathbf{f}} \Phi(\mathbf{f}_{\hat{\mathbf{w}}}^{(t-1)}) \right) \right) \right) \right\|_{\mathcal{H}}^2 \quad (160)$$

$$= \Phi(\mathbf{f}_{\hat{\mathbf{w}}}^{(t)}) + \left\| \mathbf{f}_{\hat{\mathbf{w}}}^{(t)} - \Pi_{\mathcal{K}} \left( \mathbf{f}_{\hat{\mathbf{w}}}^{(t)} - \alpha \left( \sum_{i=k}^t \mu^k (1 - \mu) \nabla_{\mathbf{f}} \mathbf{f}_{\hat{\mathbf{w}}}^{(t-k)} \right) \right) \right\|_{\mathcal{H}}^2 \quad (161)$$

$$\leq \Phi(\mathbf{f}_{\hat{\mathbf{w}}}^{(t)}) + \left\| \mathbf{f}_{\hat{\mathbf{w}}}^{(t)} - \mathbf{f}_{\hat{\mathbf{w}}}^{(t)} - \alpha \left( \sum_{k=1}^t \mu^k (1 - \mu) \nabla_{\mathbf{f}} \mathbf{f}_{\hat{\mathbf{w}}}^{(t-k)} \right) \right\|_{\mathcal{H}}^2 \quad (162)$$

### C.2 Proof of Theorem 27

## D Necessary Lemmas

**Lemma 37** (MGF of Sub-Exponential Random Variable). *Let  $x$  be a Sub-exponential variable, then we have*

$$\mathbb{E}[\exp(\lambda x)] \leq \exp(C\lambda^2) \quad (163)$$

**Theorem 38** (Matrix Chernoff, [35]). *Let  $X_k$  be a sequence of independent, random, self-adjoint matrices with dimension  $d$  s.t.*

$$X_k \succcurlyeq \mathbf{0} \quad \text{and} \quad \lambda_{\max}(X_k) \leq R \quad \text{almost surely} \quad (164)$$

*Define*

$$\mu_{\min} \triangleq \lambda_{\min}\left(\sum_k \mathbb{E}X_k\right) \quad (165)$$

*Then for  $\delta \in [0, 1]$*

$$\mathbb{P}\left\{\lambda_{\min}\left(\sum_k X_k\right) \geq (1-\delta)\mu_{\min}\right\} \leq d \cdot \left[\frac{e^{-\delta}}{(1-\delta)^{1-\delta}}\right]^{\mu_{\min}/R} \quad (166)$$

**Theorem 39** (Matrix Bernstein, [36]). *Let  $X_k$  be a sequence of independent, random, self-adjoint matrices with dimension  $d$  s.t.*

$$\mathbb{E}X_k = \mathbf{0} \quad (167)$$

**Theorem 40** (Vector Bernstein, [12, 18]). *Let  $\mathbf{x}_k$  be a sequence of independent, random vectors such that*

$$\mathbb{E}\mathbf{x}_k = \mathbf{0} \quad \forall k \quad (168)$$



## E Experimental Details

In this section we give details on datasets and hyperparameters.

### E.1 Kernel Regression

Our datasets are synthetic and are sourced from [7]

Dataset	Dimension $d$	Sample Size $n$	Source
Polynomial	3	1000	Ours
Boston Housing	13	506	[7]
Concrete Data	8	1030	[7]
Wine Quality	11	1599	[7]

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.2 Kernel Binary Classification

Dataset	Dimension $d$	Sample Size $n$	Source
Heart Disease	13	303	[14]
Breast Cancer	32	569	Kaggle

Table 5: Datasets for Kernel Binary Classification.

### E.3 Kernel Multi-Class Classification

Dataset	Dimension $d$	Sample Size $n$	Source
---------	---------------	-----------------	--------

Table 6: Datasets for Kernel Multi-Class Classification.

### E.4 Linear Regression

Dataset	Dimension $d$	Sample Size $n$	Source
Boston Housing	14	506	Kaggle
Wine Quality	11	1599	[7]
Concrete	8	1030	[7]
Drug			

Table 7: Datasets for Linear Regression.

## F Detailed Related Works

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

### F.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  $\epsilon$ -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

$$(\nabla_w F(w, u))^\top (\tilde{w} - w) \geq 0 \quad \forall \tilde{w} \in K \quad (169)$$

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  $\beta$ -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.

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

First we will give the algorithm

$$S^{(t)} = \arg \min_{T \subset S^{(0)}: |T|=(1-2\epsilon)n} \sum_{i \in T} -\log f(y_i | \langle \beta^{(t)}, \mathbf{x}_i \rangle) \quad (170)$$

$$\beta^{(t+1)} = \arg \min_{\beta, \|\beta\| \leq R} \sum_{i \in S^{(t)}} -\log f(y_i | \langle \beta^{(t)}, \mathbf{x}_i \rangle) \quad (171)$$

In Equation (170), the algorithm chooses the  $(1 - 2\epsilon)n$  points giving the least error and put this in the set  $S^{(t)}$ . Next, in Equation (171), the algorithm then finds  $\beta$  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].

$$\frac{1}{n} \sum_{i \in S} \nabla_\beta \log f(y_i | \langle \beta, \mathbf{x}_i \rangle)^\top \frac{(\beta^* - \beta)}{\|\beta^* - \beta\|} \leq \gamma \quad (172)$$

We see Equation (172) is an upper bound, instead of a lower bound, of Equation (169). Next, they prove their algorithm reaches a  $\eta$  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  $\eta$  improvement.