

Subquantile Minimization: Theory and Applications for Machine Learning

Arvind Rathnashyam* Alex Gittens†

September 1, 2023

Abstract

In this paper we propose Subquantile Minimization for adversarial corruption in the training set. We study the Huber Contamination Problem for the Kernel Learning Problem where the distribution is formed as, $\hat{\mathbb{P}} = (1 - \varepsilon)\mathbb{P} + \varepsilon\mathbb{Q}$, and we want to find the function $\inf_f \mathbb{E}_{\mathbf{x} \in \mathbb{P}} [\ell_f(\mathbf{x})]$, from the noisy distribution, $\hat{\mathbb{P}}$. We assume the adversary has knowledge of the true distribution of \mathbb{P} , and is able to corrupt the covariates and the labels of ε samples. To our knowledge, we are the first to study the problem of general kernel learning in the Huber Contamination Model. We theoretically analyze Kernel Ridge Regression and Kernel Classification and empirically show the strength of Subquantile Minimization. Furthermore, we run experiments on various datasets and compare with the state-of-the-art algorithms to show the superior performance of Subquantile Minimization.

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

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

1	Introduction	3
1.1	Related Work	3
1.2	Notation	3
2	Subquantile Minimization	4
3	Theory	4
3.1	Necessary Kernel Inequalities	6
3.2	Kernel Ridge Regression	6
3.3	Kernel Classification	9
4	Experiments	9
4.1	Kernel Ridge Regression	9
4.2	Kernel Classification	10
5	Discussion	11
A	Kernel Embedding Inequalities	14
B	Proofs	15
B.1	Proof of theorem 15	15
B.2	Proof of lemma 7	15
B.3	Proof of theorem 25	16
C	Robust Linear Regression	18
D	Base Learner Algorithm	19
E	Experimental Details	20

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), (Diakonikolas u. a., 2019), (Awasthi u. a., 2022), (Li u. a., 2021). 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). 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 $n(1-\epsilon)$ samples arbitrarily. The algorithm is then given the ϵ -corrupted data matrix \mathbf{X} and \mathbf{y} as training data.

Contributions

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

1.1 Related Work

Robust Algorithms

(Diakonikolas u. a., 2019) 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.

(Li u. a., 2021) 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.

(Awasthi u. a., 2022) 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.

(Cheng u. a., 2020) 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.

Minimax Optimization

(Jin u. a., 2020) 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.

(Yang u. a., 2022) 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: $k(\mathbf{x}, \mathbf{y}) = \exp\left(-\gamma \|\mathbf{x} - \mathbf{y}\|_2^2\right)$.

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 often denote X as the set of all data vectors, $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$. Uppercase bold $(\mathbf{X}, \mathbf{Y}, \mathbf{Z}, \dots)$ are matrices. Uppercase Roman are sets (X, S, P, Q) . Lowercase bold are vectors $(\mathbf{x}, \mathbf{y}, \mathbf{z}, \dots)$.

We denote $\mathbf{I}_{k \times k}$ as the $k \times k$ identity matrix. The spectral norm of \mathbf{A} is $\|\mathbf{A}\|_2 = \max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\| = \sigma_{\max}(\mathbf{A})$.

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.

2 Subquantile Minimization

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

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

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 (Razaviyayn u. a., 2020).

$$t^{(k+1)} = \arg \max_{t \in \mathbb{R}} g(\mathbf{w}^{(k)}, t) \quad (2)$$

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha \nabla g(\mathbf{w}^{(k)}, t^{(k+1)}) \quad (3)$$

We note this is a non-convex concave minimax optimization problem.

Remark 2. 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} .

We provide an algorithm for Subquantile Minimization of the ridge regression kernel algorithm. We solve the kernel ridge regression with Functional Gradient Descent.

Algorithm 1: SUBQ-GRADIENT

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

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

```

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

```

Algorithm 2: SUBQUANTILE

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

Output: Subquantile Matrix \mathbf{S}

```

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

```

3 Theory

To consider theoretical guarantees of Subquantile Minimization, we first analyze the inner and outer optimization problems. We first analyze kernel learning in the presence of corrupted data. Next, we provide error bounds for the two most important kernel learning problems, kernel ridge regression, and kernel classification.

Assumption 3. We represent the data matrix $\mathbf{X} = (\hat{\mathbf{P}}^\top \quad \hat{\mathbf{Q}}^\top)^\top$ and the labels vector as $\mathbf{y} = (\hat{\mathbf{y}}_P^\top \quad \hat{\mathbf{y}}_Q^\top)^\top$ where the rows of $\hat{\mathbf{P}}$ and $\hat{\mathbf{Q}}$ are sampled from the distribution, $\mathcal{N}_d(\mathbf{0}, \mathbf{I})$

Lemma 4. 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 .

$$\partial(-g(t, \mathbf{w})) = \partial \left(-t + \frac{1}{np} \sum_{i=1}^n (t - \hat{\nu}_i) \right) \quad (4)$$

$$= -1 + \frac{1}{np} \sum_{i=1}^{np} \begin{cases} 1 & \text{if } t > \hat{\nu}_i \\ 0 & \text{if } t < \hat{\nu}_i \\ [0, 1] & \text{if } t = \hat{\nu}_i \end{cases} \quad (5)$$

$$= 0 \text{ when } t = \hat{\nu}_{np} \quad (6)$$

This is equivalent to the p -quantile of the Risk. \blacksquare

Interpretation 5. From lemma 4, 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 6. 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 4, 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:

$$\nabla_{\mathbf{w}} g(t^{(k+1)}, \mathbf{w}^{(k)}) = \nabla_{\mathbf{w}} \left(\hat{\nu}_{np} - \frac{1}{np} \sum_{i=1}^n (\hat{\nu}_{np} - f(\mathbf{x}_i; \mathbf{w}, y_i))^+ \right) \quad (7)$$

$$= -\frac{1}{np} \sum_{i=1}^{np} \nabla_{\mathbf{w}} (\hat{\nu}_{np} - f(\mathbf{x}_i; \mathbf{w}, y_i))^+ \quad (8)$$

$$= \frac{1}{np} \sum_{i=1}^n \nabla_{\mathbf{w}} f(\mathbf{x}_i; \mathbf{w}^{(k)}, y_i) \begin{cases} 1 & \text{if } t > \hat{\nu}_i \\ 0 & \text{if } t < \hat{\nu}_i \\ [0, 1] & \text{if } t = \hat{\nu}_i \end{cases} \quad (9)$$

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

$$\nabla_{\mathbf{w}} g(t^{(k+1)}, \mathbf{w}^{(k)}) = \frac{1}{np} \sum_{i=1}^{np} \nabla_{\mathbf{w}} f(\mathbf{x}_i; \mathbf{w}, y_i) \quad (10)$$

This concludes the proof. \blacksquare

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

$$|g(t, \mathbf{w}) - g(t, \hat{\mathbf{w}})| \leq L \|\mathbf{w} - \hat{\mathbf{w}}\| \quad (11)$$

$$\text{where } L = \frac{K}{np} \sigma_{\max}^2 \left(\sum_{i=1}^n (\mathbf{k}_i \mathbf{k}_i^\top) \right) + \frac{2}{np} \sigma_{\max} \left(\sum_{i=1}^n \mathbf{k}_i \right) \|\mathbf{y}\|_2$$

Lemma 8. (β -Smoothness of $g(t, \mathbf{w})$ w.r.t \mathbf{w}). Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ represent the rows of the data matrix \mathbf{X} . It then follows:

$$\|\nabla_{\mathbf{w}} g(t, \mathbf{w}) - \nabla_{\mathbf{w}} g(t, \hat{\mathbf{w}})\| \leq \beta \|\mathbf{w} - \hat{\mathbf{w}}\| \quad (12)$$

$$\text{where } \beta = \frac{2}{np} \sigma_{\max}(\mathbf{X})$$

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

$$\|\nabla_{\mathbf{w}}^2 g(t, \mathbf{w})\|_2 = \left\| \frac{2}{np} \sum_{\mathbf{x} \in S} \mathbf{k}_x \mathbf{k}_x^\top \right\|_2 \stackrel{\text{eqn. (50)}}{\leq} \left\| \frac{2}{np} \sum_{\mathbf{x} \in X} \mathbf{k}_x \mathbf{k}_x^\top \right\|_2 = \frac{2}{np} \sigma_{\max} \left(\sum_{\mathbf{x} \in X} \mathbf{k}_x \mathbf{k}_x^\top \right) \stackrel{\text{lem. 13}}{\leq} \quad (13)$$

This concludes the proof. \blacksquare

3.1 Necessary Kernel Inequalities

We will first extend the idea of Resilience [Steinhardt u. a. \(2018\)](#) to kernel learning.

Definition 9. (Resilience) from ([Steinhardt u. a., 2018](#)). Let \mathcal{H} represent a RKHS, then given the feature mapping $\phi : \mathcal{X} \rightarrow \mathcal{H}$, and the set $X = \{\mathbf{x}_i\}_{i=1}^n = \widehat{P} \cup \widehat{Q}$, such that $|\widehat{P}| = n(1 - \epsilon)$ and $|\widehat{Q}| = n\epsilon$, 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 Hilbert Space.

Without the idea of resilience defined in [definition 9](#), we will be unable to put error bounds on our algorithm.

Lemma 10. *Let S be the set of elements in the subquantile, then*

$$\left\| \frac{1}{np} \sum_{i \in S \cap \widehat{P}} \mathbf{k}_i \right\| \leq \mathcal{O}() \quad (14)$$

Proof.

$$\left\| \frac{1}{np} \sum_{i \in S \cap \widehat{P}} \mathbf{k}_i \right\| = \frac{1}{np} \left\| \sum_{i \in S \cap \widehat{P}} \sum_{j \in X} \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) \right\| \quad (15)$$

■

Lemma 11. *Under the same setting as [lemma 10](#),*

$$\left\| \frac{1}{np} \sum_{i \in S \cap \widehat{Q}} \mathbf{k}_i \right\| \leq \mathcal{O}() \quad (16)$$

Lemma 12. *Under the same setting as [lemma 10](#),*

$$\left\| \frac{1}{np} \sum_{i \in S \cap \widehat{P}} \xi_i \mathbf{k}_i \right\| \leq \mathcal{O}() \quad (17)$$

Lemma 13. *Under the same setting as [lemma 10](#),*

$$\left\| \frac{1}{np} \sum_{i \in S} \mathbf{k}_i \mathbf{k}_i^\top \right\| \leq \mathcal{O}() \quad (18)$$

3.2 Kernel Ridge Regression

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: $f(\mathbf{x}^*; \mathbf{w}) = \sum_{i=1}^n \mathbf{w}_i \kappa(\mathbf{x}_i, \mathbf{x}^*)$

First let us define the kernels:

RBK Kernel: $\mathbf{k}_\mathbf{x} \triangleq \left(\exp(-\gamma \|\mathbf{x} - \mathbf{x}_1\|_2^2) \quad \cdots \quad \exp(-\gamma \|\mathbf{x} - \mathbf{x}_n\|_2^2) \right)^\top$ where $\gamma \triangleq \frac{1}{d}$

Linear Kernel: $\mathbf{k}_\mathbf{x} \triangleq (\mathbf{x}^\top \mathbf{x}_1 \quad \cdots \quad \mathbf{x}^\top \mathbf{x}_n)$

Assumption 14. (Gauss-Markov Assumption) Let $\widehat{P} = \{(\mathbf{p}_1, y_1), (\mathbf{p}_2, y_2), \dots, (\mathbf{p}_m, y_m)\}$ and $K(\mathbf{x}, \mathbf{y})$ be a proper kernel. If $\mathbf{p} \in \widehat{P}$, then define $\mathbf{k}_{\mathbf{p}} \triangleq (K(\mathbf{p}, \mathbf{p}_1) \ K(\mathbf{p}, \mathbf{p}_2) \ \dots \ K(\mathbf{p}, \mathbf{p}_m))^{\top}$: we first assume y can be modeled as a non-linear function of the data to smooth the landscape of the function.

$$y = f(\mathbf{p}) + \widehat{\xi}_P \text{ where } \widehat{\xi}_P \sim \mathcal{N}(0, \widehat{\sigma}_P^2) \quad (19)$$

Furthermore, we assume this function can be represented by the Reproducing Kernel Hilbert Space:

$$y = (\mathbf{k}_{\mathbf{p}}^{\top} \mathbf{w}^* - f(\mathbf{p})) + \xi_P \text{ where } \xi_P \sim \mathcal{N}(0, \sigma_P^2) \quad (20)$$

Theorem 15. (Monotonically Decreasing). Let f be a convex loss function and \mathbf{X} follows [assumption 3](#) with learning schedule $\alpha_{(1)}, \alpha_{(2)}, \dots, \alpha_{(T)}$. Then it follows at any iteration $k \in \mathbb{N}$:

$$g(t^{(k+1)}, \mathbf{w}^{(k+1)}) \leq g(t^{(k)}, \mathbf{w}^{(k)}) \quad (21)$$

From our definition of $S^{(k)}$ in [theorem 15](#), we are interested in as $k \rightarrow \infty$ the quantities: $|\mathbf{x} \in S^{(k)} \cap \widehat{P}|$ and $|\mathbf{x} \in S^{(k)} \cap \widehat{Q}|$, where the latter cardinality represents the number of corrupted points in the subquantile set.

Definition 16. (Moreau Envelope). Let f be proper lower semi-continuous convex function $f: \mathcal{X} \rightarrow \mathbb{R}$, then the Moreau Envelope is defined as:

$$f_{\lambda}(\mathbf{x}) \triangleq \inf_{\widehat{\mathbf{x}} \in \mathcal{X}} \left(f(\widehat{\mathbf{x}}) + \frac{1}{2\lambda} \|\mathbf{x} - \widehat{\mathbf{x}}\|_2^2 \right) \quad (22)$$

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

Definition 17. (First-Order Stationary Point). For any proper lower semi-continuous function f and closed convex set \mathcal{K} consider its associated Moreau envelope $f_{\beta}(\mathbf{w})$ in [definition 16](#). Then we say that a point \mathbf{w} is a first-order stationary point if $\|\Phi_{\beta}(\mathbf{w})\|_2 = 0$

Definition 18. (First-Order Stationary Point). Let $\Phi(\mathbf{w}) = \max_t g(t, \mathbf{w})$. Then \mathbf{w} is a first-order stationary point if

$$(\nabla_{\mathbf{w}} \Phi(\mathbf{w}))^{\top} (\widetilde{\mathbf{w}} - \mathbf{w}) \geq 0 \ \forall \widetilde{\mathbf{w}} \in \mathcal{K} \quad (23)$$

The idea of the First-Order Stationary Point will have significance in our analysis of the base learner algorithm in [§ appendix D](#).

Assumption 19. Define $\Phi(\cdot)$ as the function in [remark 2](#). Then it follows $\arg \min_{\mathbf{w} \in \mathbb{R}^d} \Phi(\mathbf{w}) = \mathbf{w}^*$

Theorem 20. Let $\widehat{\mathbf{w}}$ be a First-Order Stationary Point defined in [definition 17](#), it then follows:

Proof. We will define the function Φ as in [remark 2](#). The derivative of the Moreau Envelope is well known, $\nabla \Phi_{\lambda}(\mathbf{w}) = \mathbf{w} - \frac{1}{2\lambda} \text{prox}_{\lambda \Phi}(\mathbf{w})$. Thus, if $\|\nabla \Phi_{1/2\ell}(\mathbf{w})\| = 0$, then it follows $\mathbf{w} = \ell \text{prox}_{1/2\ell \Phi}(\mathbf{w}) \stackrel{\text{def}}{=} \arg \min_{\widehat{\mathbf{w}}} \left(\Phi(\mathbf{w}) + \ell \|\widehat{\mathbf{w}} - \mathbf{w}\|^2 \right)$. Thus it follows for any $\widehat{\mathbf{w}}$ s.t. $\Phi(\widehat{\mathbf{w}}) < \Phi(\mathbf{w})$, then it holds $\ell \|\mathbf{w} - \widehat{\mathbf{w}}\|^2 > \Phi(\widehat{\mathbf{w}}) - \Phi(\mathbf{w})$. Therefore, if we let $\mathbf{w} = \mathbf{w}^*$, then we find $\Phi(\widehat{\mathbf{w}}) < \Phi(\mathbf{w}^*) + \ell \|\mathbf{w}^* - \widehat{\mathbf{w}}\|^2$.

Upper bound of $\|\mathbf{w}^* - \widehat{\mathbf{w}}\|^2$ The idea here is if $\|\widehat{\mathbf{w}} - \mathbf{w}^*\|$ is large, then the norm of the derivative must be correspondingly large. However, without any distribution assumptions on \widehat{Q} , how to show $\widehat{\mathbf{w}}$ is not a good weight vector for \widehat{Q} . ■

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 21. ((Rockafellar, 1970)). Assume the function Φ is ℓ -weakly convex. Let $\lambda < \frac{1}{\ell}$, and denote $\widehat{\mathbf{w}} = \arg \min_{\mathbf{w}'} \left(\Phi(\mathbf{w}') + \frac{1}{2\lambda} \|\mathbf{w} - \mathbf{w}'\|^2 \right)$, $\|\nabla \Phi_{\lambda}\| \leq \epsilon$ implies:

$$\|\widehat{\mathbf{w}} - \mathbf{w}\| = \lambda \epsilon \text{ and } \min_{\mathbf{g} \in \partial \Phi(\widehat{\mathbf{w}})} \|\mathbf{g}\| \leq \epsilon \quad (24)$$

Definition 22. (Approximate First-Order Stationary Point) from (Cheng u. a., 2020). For any function f and closed convex set \mathcal{K} consider its associated Moreau envelope $f_\beta(\mathbf{w})$ in definition 16. 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 22 is used in the analysis of the minimax algorithm in (Cheng u. a., 2020). First, if you can prove a stationary point is good, theorem 20, then using lemma 21, you can show an approximate stationary point is good.

Definition 23. (Approximate First-Order Stationary Point) from (Awasthi u. a., 2022). A regression coefficients vector, \mathbf{w} is an approximate stationary point if:

$$\left(\sum_{\mathbf{x} \in S} \mathbf{k}_x (\mathbf{k}_x^\top \mathbf{w} - y) \right)^\top \left(\frac{(\mathbf{w} - \mathbf{w}^*)}{\|\mathbf{w} - \mathbf{w}^*\|} \right) \leq \eta \quad (25)$$

We adopt the proof strategy of (Awasthi u. a., 2022) and (Cheng u. a., 2020), and have a two-part proof strategy. First we show an approximate stationary point is close to the true distribution of \mathbb{P} . Then, we analyze the optimization to show ?? 1 converges to an approximate stationary point in a polynomial number of iterations.

Theorem 24. Let $\hat{\mathbf{w}}$ be the vector returned from ?? 1 after T iterations with $\alpha = \frac{1}{\beta}$, then we reach an approximation stationary point defined in ??,

$$\left(\sum_{\mathbf{x} \in S} \mathbf{k}_x (\mathbf{k}_x^\top \mathbf{w} - y) \right)^\top \left(\frac{(\mathbf{w} - \mathbf{w}^*)}{\|\mathbf{w} - \mathbf{w}^*\|} \right) \leq \eta \quad (26)$$

Theorem 25. (Approximate Stationary Point is Good). Let $\hat{\mathbf{w}}$ be a η -stationary point as defined in definition 23. It then follows:

$$\|\hat{\mathbf{w}} - \mathbf{w}^*\| \leq \frac{\|\sum_{\mathbf{x} \in S \cap \hat{P}} \xi_x \mathbf{k}_x\| + \sqrt{\frac{p}{1-p}} (\|\sum_{i \in \hat{P} \cap S} \xi_i\| \|\sum_{\mathbf{x} \in S \cap \hat{P}} \mathbf{k}_x\|)}{\|\sum_{\mathbf{x} \in S \cap \hat{P}} \mathbf{k}_x\|^2 - \sqrt{\frac{p}{1-p}} \|\sum_{\mathbf{x} \in \hat{P} \setminus S} \mathbf{k}_x\|} \quad (27)$$

where $S^{(T)}$ represents the np set of points in the subquantile.

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

Lemma 26. content...

Proof. content... ■

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 (Lin u. a., 2020) Theorem 31 and (Cheng u. a., 2020) Lemma 4.2, it follows:

$$\mathbb{E} \left[\|\nabla \Phi_{1/2\ell}(\bar{\mathbf{w}})\|^2 \right] \leq 2 \cdot \frac{(\Phi_{1/2\ell}(\mathbf{w}_0) - \min \Phi(\mathbf{w})) + \ell \beta^2 \gamma^2}{\gamma \sqrt{T+1}} \quad (28)$$

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

3.3 Kernel Classification

4 Experiments

4.1 Kernel Ridge 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. (3)
6: end
7: return  $\mathbf{w}_{(T)}$ 

```

Objectives	Test RMSE (Polynomial Regression (Degree = 3))			
	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon = 0.4$
KRR	0.460 _(0.2143)	1.171 _(0.7809)	0.950 _(0.3053)	1.230 _(0.4678)
TERM (Li u. a., 2021)	∞	∞	∞	∞
SEVER (Diakonikolas u. a., 2019)	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)			
	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon = 0.4$
KRR	0.544 _(0.0712)	0.862 _(0.1199)	0.865 _(0.0811)	1.049 _(0.2249)
TERM (Li u. a., 2021)	0.888 _(0.1360)	0.891 _(0.1699)	1.023 _(0.1329)	0.931 _(0.0433)
SEVER (Diakonikolas u. a., 2019)	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 $\lambda = 2$

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.

Objectives	Test RMSE (Concrete Data Regression)			
	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon = 0.4$
KRR	0.802 _(0.0324)	0.929 _(0.0209)	0.993 _(0.0441)	0.775 _(0.0514)
TERM (Li u. a., 2021)	0.874 _(0.0205)	0.916 _(0.0421)	0.840 _(0.0249)	0.878 _(0.0749)
SEVER (Diakonikolas u. a., 2019)	0.532 _(0.0134)	0.516 _(0.0340)	0.526 _(0.0217)	0.552 _(0.0444)
SUBQUANTILE($p = (1 - \epsilon)$)	0.468_(0.0220)	0.491_(0.0271)	0.555_(0.0391)	0.566_(0.0405)
Genie ERM	∞	∞	∞	∞

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

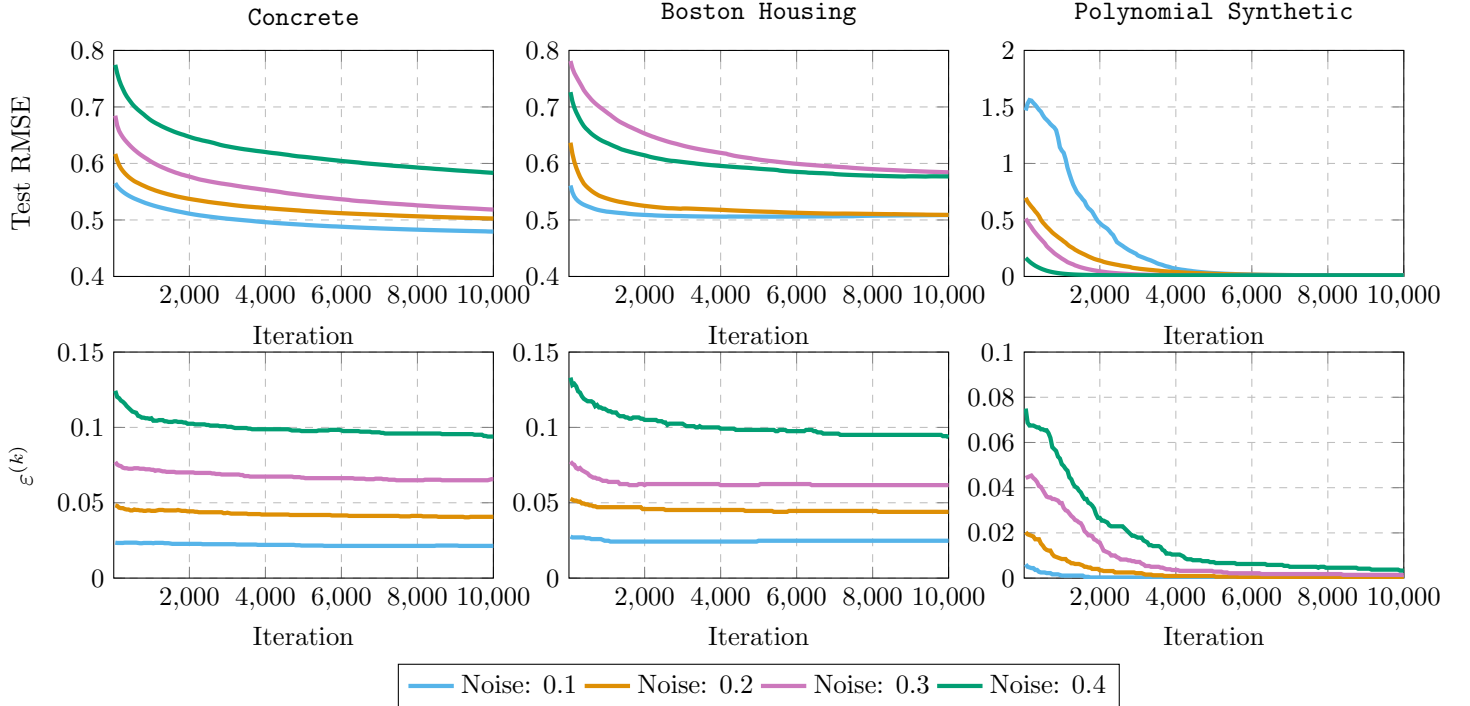


Figure 1: Test RMSE over the iterations in Concrete, Boston Housing, and Polynomial Datasets for SUBQUANTILE at different noise levels

4.2 Kernel Classification

Algorithm 4: SUBQ-KERNEL-CLASSIFICATION

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

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

- 1: $\mathbf{w}_{(0)} \leftarrow (\mathbf{K}^\top \mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{K}^\top \mathbf{y}$ \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 eqn. (3)
 - 6: **end**
 - 7: **return** $\mathbf{w}_{(T)}$
-

5 Discussion

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

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

- [Awasthi u. a. 2022] AWASTHI, Pranjal ; DAS, Abhimanyu ; KONG, Weihao ; SEN, Rajat: Trimmed Maximum Likelihood Estimation for Robust Generalized Linear Model. In: OH, Alice H. (Hrsg.) ; AGARWAL, Alekh (Hrsg.) ; BELGRAVE, Danielle (Hrsg.) ; CHO, Kyunghyun (Hrsg.): *Advances in Neural Information Processing Systems*, URL https://openreview.net/forum?id=VHmdFPy4U_u, 2022
- [Cheng u. a. 2020] CHENG, Yu ; DIAKONIKOLAS, Ilias ; GE, Rong ; SOLTANOLKOTABI, Mahdi: High-dimensional Robust Mean Estimation via Gradient Descent. In: III, Hal D. (Hrsg.) ; SINGH, Aarti (Hrsg.): *Proceedings of the 37th International Conference on Machine Learning* Bd. 119, PMLR, 13–18 Jul 2020, S. 1768–1778. – URL <https://proceedings.mlr.press/v119/cheng20a.html>
- [Diakonikolas u. a. 2019] DIAKONIKOLAS, Ilias ; KAMATH, Gautam ; KANE, Daniel M. ; LI, Jerry ; STEINHARDT, Jacob ; STEWART, Alistair: Sever: A Robust Meta-Algorithm for Stochastic Optimization. In: *Proceedings of the 36th International Conference on Machine Learning*, JMLR, Inc., 2019 (ICML '19), S. 1596–1606
- [Dua und Graff 2017] DUA, Dheeru ; GRAFF, Casey: *UCI Machine Learning Repository*. 2017. – URL <http://archive.ics.uci.edu/ml>
- [Jambulapati u. a. 2020] JAMBULAPATI, Arun ; LI, Jerry ; TIAN, Kevin: Robust Sub-Gaussian Principal Component Analysis and Width-Independent Schatten Packing. In: LAROCHELLE, H. (Hrsg.) ; RANZATO, M. (Hrsg.) ; HADSELL, R. (Hrsg.) ; BALCAN, M.F. (Hrsg.) ; LIN, H. (Hrsg.): *Advances in Neural Information Processing Systems* Bd. 33, Curran Associates, Inc., 2020, S. 15689–15701. – URL https://proceedings.neurips.cc/paper_files/paper/2020/file/b58144d7e90b5a43edcce1ca9e642882-Paper.pdf
- [Jin u. a. 2020] JIN, Chi ; NETRAPALLI, Praneeth ; JORDAN, Michael: What is Local Optimality in Nonconvex-Nonconcave Minimax Optimization? In: III, Hal D. (Hrsg.) ; SINGH, Aarti (Hrsg.): *Proceedings of the 37th International Conference on Machine Learning* Bd. 119, PMLR, 13–18 Jul 2020, S. 4880–4889. – URL <https://proceedings.mlr.press/v119/jin20e.html>
- [Li u. a. 2021] LI, Tian ; BEIRAMI, Ahmad ; SANJABI, Maziar ; SMITH, Virginia: Tilted Empirical Risk Minimization. In: *International Conference on Learning Representations*, URL <https://openreview.net/forum?id=K5YasWXZT30>, 2021
- [Liao u. a. 2020] LIAO, Zhenyu ; COUILLET, Romain ; MAHONEY, Michael W.: A random matrix analysis of random Fourier features: beyond the Gaussian kernel, a precise phase transition, and the corresponding double descent. In: LAROCHELLE, H. (Hrsg.) ; RANZATO, M. (Hrsg.) ; HADSELL, R. (Hrsg.) ; BALCAN, M.F. (Hrsg.) ; LIN, H. (Hrsg.): *Advances in Neural Information Processing Systems* Bd. 33, Curran Associates, Inc., 2020, S. 13939–13950. – URL https://proceedings.neurips.cc/paper_files/paper/2020/file/a03fa30821986dff10fc66647c84c9c3-Paper.pdf
- [Lin u. a. 2020] LIN, Tianyi ; JIN, Chi ; JORDAN, Michael I.: Near-Optimal Algorithms for Minimax Optimization. In: ABERNETHY, Jacob (Hrsg.) ; AGARWAL, Shivani (Hrsg.): *Proceedings of Thirty Third Conference on Learning Theory* Bd. 125, PMLR, 09–12 Jul 2020, S. 2738–2779. – URL <https://proceedings.mlr.press/v125/lin20a.html>
- [Rahimi und Recht 2007] RAHIMI, Ali ; RECHT, Benjamin: Random Features for Large-Scale Kernel Machines. In: PLATT, J. (Hrsg.) ; KOLLER, D. (Hrsg.) ; SINGER, Y. (Hrsg.) ; ROWEIS, S. (Hrsg.): *Advances in Neural Information Processing Systems* Bd. 20, Curran Associates, Inc., 2007. – URL https://proceedings.neurips.cc/paper_files/paper/2007/file/013a006f03dbc5392effeb8f18fda755-Paper.pdf
- [Razaviyayn u. a. 2020] RAZAVIYAYN, Meisam ; HUANG, Tianjian ; LU, Songtao ; NOUIEHED, Maher ; SANJABI, Maziar ; HONG, Mingyi: Nonconvex Min-Max Optimization: Applications, Challenges, and Recent Theoretical Advances. In: *IEEE Signal Processing Magazine* 37 (2020), Nr. 5, S. 55–66

- [Rockafellar 1970] ROCKAFELLAR, Ralph T.: *Convex Analysis*. Princeton : Princeton University Press, 1970. – URL <https://doi.org/10.1515/9781400873173>. – ISBN 9781400873173
- [Schneider 2016] SCHNEIDER, Markus: Probability Inequalities for Kernel Embeddings in Sampling without Replacement. In: *International Conference on Artificial Intelligence and Statistics*, 2016
- [Staib und Jegelka 2019] STAIB, Matthew ; JEGELKA, Stefanie: Distributionally Robust Optimization and Generalization in Kernel Methods. In: WALLACH, H. (Hrsg.) ; LAROCHELLE, H. (Hrsg.) ; BEYGELZIMER, A. (Hrsg.) ; ALCHÉ-BUC, F. d' (Hrsg.) ; FOX, E. (Hrsg.) ; GARNETT, R. (Hrsg.): *Advances in Neural Information Processing Systems* Bd. 32, Curran Associates, Inc., 2019. – URL https://proceedings.neurips.cc/paper_files/paper/2019/file/1770ae9e1b6bc9f5fd2841f141557ffb-Paper.pdf
- [Steinhardt u. a. 2018] STEINHARDT, Jacob ; CHARIKAR, Moses ; VALIANT, Gregory: Resilience: A Criterion for Learning in the Presence of Arbitrary Outliers. In: KARLIN, Anna R. (Hrsg.): *9th Innovations in Theoretical Computer Science Conference, ITCS 2018, January 11-14, 2018, Cambridge, MA, USA* Bd. 94, Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2018, S. 45:1–45:21. – URL <https://doi.org/10.4230/LIPIcs.ITCS.2018.45>
- [Yang u. a. 2022] YANG, Junchi ; ORVIETO, Antonio ; LUCCHI, Aurelien ; HE, Niao: Faster Single-loop Algorithms for Minimax Optimization without Strong Concavity. In: CAMPS-VALLS, Gustau (Hrsg.) ; RUIZ, Francisco J. R. (Hrsg.) ; VALERA, Isabel (Hrsg.): *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics* Bd. 151, PMLR, 28–30 Mar 2022, S. 5485–5517. – URL <https://proceedings.mlr.press/v151/yang22b.html>

A Kernel Embedding Inequalities

Lemma 28. (*Resilience on Inlier Samples*). Let $X = \{\mathbf{x}_i\}_{i=1}^n$ and $\hat{P} = \{\mathbf{x}_i\}_{i=1}^{np}$, and $[\mathbf{k}_i]_j = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$. If the conditions in assumption [definition 9](#) it then follows:

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

where $\|\phi(\mathbf{x})\| \leq d$ for any $\mathbf{x} \in \mathcal{X}$ a.s. A similar inequality can be found in ([Schneider, 2016](#)), Theorem 2

Lemma 29. 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:

$$\left\| \sum_{i=1}^n \xi_i \right\| \leq C \text{ with high probability} \quad (29)$$

First, we can note,

$$\sum_{i=1}^n \xi_i \sim \mathcal{N}(n\mu, n\sigma^2) \quad (30)$$

With Hoeffding's Inequality, we have:

$$\mathbb{P} \left\{ \left\| \frac{1}{n} \sum_{i=1}^n \xi_i - \mu \right\| \geq \epsilon \right\} \leq 2 \exp \left(-\frac{2n^2\epsilon^2}{\sum_{i=1}^n \|\xi_i\|_{\psi_2}^2} \right) \quad (31)$$

Lemma 30. (*Resilience of corrupted sample*). Let $X = \{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^n$. It then follows for any subset $S \subset X$, s.t. $|S| \geq np\epsilon^{(k)}$.

B Proofs

B.1 Proof of **theorem 15**

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

Note this is an equivalent characterization of g from **lemma 4**.

$$F(\mathbf{w}^{(k+1)}, S^{(k+1)}) \leq F(\mathbf{w}^{(k)}, S^{(k)}) \quad (32)$$

$$F(\mathbf{w}^{(k+1)}, S^{(k+1)}) - F(\mathbf{w}^{(k)}, S^{(k+1)}) \leq F(\mathbf{w}^{(k)}, S^{(k)}) - F(\mathbf{w}^{(k)}, S^{(k+1)}) \quad (33)$$

Upper Bound of LHS

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

$$F(\mathbf{w}^{(k+1)}, S^{(k+1)}) - F(\mathbf{w}^{(k)}, S^{(k+1)}) \leq \langle \nabla_{\mathbf{w}} F(\mathbf{w}^{(k)}, S^{(k+1)}), \mathbf{w}^{(k+1)} - \mathbf{w}^{(k)} \rangle + \frac{L}{2} \left\| \mathbf{w}^{(k+1)} - \mathbf{w}^{(k)} \right\|_2^2 \quad (34)$$

$$= \langle \nabla_{\mathbf{w}} F(\mathbf{w}^{(k)}, S^{(k+1)}), -\alpha^{(k)} \nabla_{\mathbf{w}} F(\mathbf{w}^{(k)}, S^{(k+1)}) \rangle + \frac{L}{2} \left\| \mathbf{w}^{(k+1)} - \mathbf{w}^{(k)} \right\|_2^2 \quad (35)$$

$$= -\alpha^{(k)} \left\| \nabla_{\mathbf{w}} F(\mathbf{w}^{(k)}, S^{(k+1)}) \right\|_2^2 + \frac{L}{2} \left\| \mathbf{w}^{(k+1)} - \mathbf{w}^{(k)} \right\|_2^2 \quad (36)$$

$$= -\alpha^{(k)} \left\| \nabla_{\mathbf{w}} F(\mathbf{w}^{(k)}, S^{(k+1)}) \right\|_2^2 + \frac{L\alpha_{(k)}^2}{2} \left\| \nabla_{\mathbf{w}} F(\mathbf{w}^{(k)}, S^{(k+1)}) \right\|_2^2 \quad (37)$$

$$= \left(\frac{L\alpha_{(k)}^2}{2} - \alpha_{(k)} \right) \left\| \nabla_{\mathbf{w}} F(\mathbf{w}^{(k)}, S^{(k+1)}) \right\|_2^2 \quad (38)$$

Lower Bound of RHS

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

$$F(\mathbf{w}^{(k)}, S^{(k)}) - F(\mathbf{w}^{(k)}, S^{(k+1)}) = \sum_{\mathbf{x} \in S^{(k)} \setminus S^{(k+1)}} f(\mathbf{x}; \mathbf{w}^{(k)}, y) - \sum_{\mathbf{x} \in S^{(k+1)} \setminus S^{(k)}} f(\mathbf{x}; \mathbf{w}^{(k)}, y) \quad (39)$$

$$\geq |S^{(k)} \setminus S^{(k+1)}| \inf \left\{ f(\mathbf{x}; \mathbf{w}^{(k)}, y) : \mathbf{x} \in S^{(k)} \setminus S^{(k+1)} \right\} - |S^{(k+1)} \setminus S^{(k)}| \sup \left\{ f(\mathbf{x}; \mathbf{w}^{(k)}, y) : \mathbf{x} \in S^{(k+1)} \setminus S^{(k)} \right\} \quad (40)$$

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

$$= \eta \left(\inf \left\{ f(\mathbf{x}; \mathbf{w}^{(k)}, y) : \mathbf{x} \in S^{(k)} \setminus S^{(k+1)} \right\} - \sup \left\{ f(\mathbf{x}; \mathbf{w}^{(k)}, y) : \mathbf{x} \in S^{(k+1)} \setminus S^{(k)} \right\} \right) \quad (41)$$

$$= \eta \left(\hat{\nu}_{np+1}^{(k)} - \hat{\nu}_{np}^{(k)} \right) \quad (42)$$

$$\geq 0 \quad (43)$$

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

B.2 Proof of **lemma 7**

Proof. We use the ℓ_2 norm of the gradient to bound L from above.

$$\left\| \nabla_{\mathbf{w}} g(t, \mathbf{w}) \right\|_2 = \left\| \frac{2}{np} \sum_{i=1}^n \mathbb{1}_{t \geq (\mathbf{x}_i^\top \mathbf{w} - y_i)^2} (\mathbf{x}_i (\mathbf{x}_i^\top \mathbf{w} - y_i)) \right\|_2 \quad (44)$$

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

$$= \left\| \frac{2}{np} \sum_{i=1}^m \mathbf{k}_i (\mathbf{k}_i^\top \mathbf{w} - y_i) \right\|_2 \quad (45)$$

$$\stackrel{(a)}{\leq} \frac{2}{np} \left(\left\| \sum_{i=1}^m (\mathbf{k}_i \mathbf{k}_i^\top) \mathbf{w} \right\|_2 + \left\| \sum_{i=1}^m \mathbf{k}_i y_i \right\|_2 \right) \quad (46)$$

$$\stackrel{(b)}{\leq} \frac{2}{np} \left(\left\| \sum_{i=1}^m (\mathbf{k}_i \mathbf{k}_i^\top) \right\|_2 \|\mathbf{w}\|_2 + \left\| \sum_{i=1}^m \mathbf{k}_i \right\|_2 \|y\|_2 \right) \quad (47)$$

$$\stackrel{(c)}{\leq} \frac{2}{np} \left(\sigma_{\max} \left(\sum_{i=1}^n (\mathbf{k}_i \mathbf{k}_i^\top) \right) \|\mathbf{w}\|_2 + \left\| \sum_{i=1}^m \mathbf{k}_i \right\|_2 \|y\|_2 \right) \quad (48)$$

$$\stackrel{(d)}{\leq} \frac{K}{np} \sigma_{\max} \left(\sum_{i=1}^n (\mathbf{k}_i \mathbf{k}_i^\top) \right) + \frac{2}{np} \left\| \sum_{i=1}^m \mathbf{k}_i \right\|_2 \|y\|_2 \quad (49)$$

where (a) follows from Triangle Inequality, (b) follows from Cauchy-Schwarz Inequality, (d) follows from assuming $\|\mathbf{w}\|_2 \leq K$, where K is some positive constant. In (c), we note that \mathbf{X}_m is positive semi-definite.

$$\mathbf{v}^\top \mathbf{X}_m \mathbf{v} = \mathbf{v}^\top \left(\sum_{i=1}^m \mathbf{k}_i \mathbf{k}_i^\top \right) \mathbf{v} = \sum_{i=1}^m \mathbf{v}^\top \mathbf{k}_i \mathbf{k}_i^\top \mathbf{v} = \sum_{i=1}^m (\mathbf{v}^\top \mathbf{k}_i)^2 \leq \sum_{i=1}^n (\mathbf{v}^\top \mathbf{k}_i)^2 \quad (50)$$

This concludes the proof. ■

B.3 Proof of [theorem 25](#)

Proof. This proof follows a similar structure to ([Awasthi u. a., 2022](#)), Lemma A.1

$$\left(\frac{1}{np} \sum_{\mathbf{x} \in S} \mathbf{k}_x (\mathbf{k}_x^\top \mathbf{w} - y) \right)^\top (\mathbf{w} - \mathbf{w}^*) \leq \eta \|\hat{\mathbf{w}} - \mathbf{w}^*\| \quad (51)$$

$$\Leftrightarrow \left(\frac{1}{np} \sum_{\mathbf{x} \in S \cap \hat{P}} \mathbf{k}_x (\mathbf{k}_x^\top \hat{\mathbf{w}} - y) \right)^\top (\hat{\mathbf{w}} - \mathbf{w}^*) \leq \left(-\frac{1}{np} \sum_{\mathbf{x} \in S \cap \hat{Q}} \mathbf{k}_x (\mathbf{k}_x^\top \hat{\mathbf{w}} - y) \right)^\top (\hat{\mathbf{w}} - \mathbf{w}^*) + \eta \|\hat{\mathbf{w}} - \mathbf{w}^*\| \quad (52)$$

Lower Bound on LHS

$$\left(\frac{1}{np} \sum_{\mathbf{x} \in S \cap \hat{P}} \mathbf{k}_x (\mathbf{k}_x^\top \hat{\mathbf{w}} - y) \right)^\top (\hat{\mathbf{w}} - \mathbf{w}^*) \quad (53)$$

$$\stackrel{\text{asm. 14}}{=} \left(\frac{1}{np} \sum_{\mathbf{x} \in S \cap \hat{P}} \mathbf{k}_x (\mathbf{k}_x^\top \mathbf{w} - \mathbf{k}_x^\top \mathbf{w}^* - \xi_i) \right)^\top (\mathbf{w} - \mathbf{w}^*) \quad (54)$$

$$= \frac{1}{np} \sum_{\mathbf{x} \in S \cap \hat{P}} (\mathbf{w} - \mathbf{w}^*)^\top (\mathbf{k}_x \mathbf{k}_x^\top) (\mathbf{w} - \mathbf{w}^*) - \xi_i \mathbf{k}_x^\top (\mathbf{w} - \mathbf{w}^*) \quad (55)$$

$$\geq \frac{1}{np} \sum_{\mathbf{x} \in S \cap \hat{P}} (\mathbf{k}_x^\top (\mathbf{w} - \mathbf{w}^*))^2 - \|\xi_i \mathbf{k}_x\| \|\mathbf{w} - \mathbf{w}^*\| \quad (56)$$

$$\geq \frac{1}{np} \left(K \|\mathbf{w} - \mathbf{w}^*\|^2 - \left\| \sum_{\mathbf{x} \in S \cap \hat{P}} \xi_i \mathbf{k}_x \right\| \|\mathbf{w} - \mathbf{w}^*\|_{\mathcal{H}} \right) \quad (57)$$

where K is a lower bound on $\|\sum_{\mathbf{x} \in S \cap \hat{P}} \mathbf{k}_{\mathbf{x}} \mathbf{k}_{\mathbf{x}}^\top - \mathbb{E}[\mathbf{k}_{\mathbf{x}} \mathbf{k}_{\mathbf{x}}^\top]\|$. Note this term will be significantly greater than the other term in the denominator.

Upper Bound on RHS

$$\left(-\frac{1}{np} \sum_{\mathbf{x} \in S \cap \hat{Q}} \mathbf{k}_{\mathbf{x}} (\mathbf{k}_{\mathbf{x}}^\top \hat{\mathbf{w}} - y) \right)^\top (\hat{\mathbf{w}} - \mathbf{w}^*) \quad (58)$$

$$\stackrel{\text{Cauchy-Schwarz}}{\leq} \left(\frac{1}{np} \sum_{\mathbf{x} \in S \cap \hat{Q}} (\mathbf{k}_{\mathbf{x}}^\top \mathbf{w} - y)^2 \right)^{1/2} \left(\frac{1}{np} \sum_{\mathbf{x} \in S \cap \hat{Q}} (\mathbf{k}_{\mathbf{x}}^\top (\mathbf{w} - \mathbf{w}^*))^2 \right)^{1/2} \quad (59)$$

$$\leq \frac{1}{np} \left(\sum_{\mathbf{x} \in S \cap \hat{Q}} (\mathbf{k}_{\mathbf{x}}^\top \mathbf{w} - y)^2 \right)^{1/2} \left(\left(\sum_{\mathbf{x} \in S \cap \hat{Q}} \|\mathbf{k}_{\mathbf{x}}\|^2 \right) \|\mathbf{w} - \mathbf{w}^*\|^2 \right)^{1/2} \quad (60)$$

$$\stackrel{(a)}{\leq} \frac{1}{np} \left(\sum_{\mathbf{x} \in S \cap \hat{Q}} (\mathbf{k}_{\mathbf{x}}^\top \mathbf{w} - y)^2 \right)^{1/2} \left(\left\| \sum_{\mathbf{x} \in S \cap \hat{Q}} \mathbf{k}_{\mathbf{x}} \right\|^2 \|\mathbf{w} - \mathbf{w}^*\|^2 \right)^{1/2} \quad (61)$$

$$\stackrel{(b)}{\leq} \frac{1}{np} \left(\frac{p}{1-p} \sum_{\mathbf{x} \in \hat{P} \setminus S} (\mathbf{k}_{\mathbf{x}}^\top (\mathbf{w} - \mathbf{w}^*) - \xi_i)^2 \right)^{1/2} \left(\left\| \sum_{\mathbf{x} \in S \cap \hat{Q}} \mathbf{k}_{\mathbf{x}} \right\| \|\mathbf{w} - \mathbf{w}^*\| \right) \quad (62)$$

$$\stackrel{(c)}{\leq} \frac{1}{np} \left(\frac{p}{1-p} \left(\left\| \sum_{\mathbf{x} \in \hat{P} \setminus S} \mathbf{k}_{\mathbf{x}} \right\| \|\mathbf{w} - \mathbf{w}^*\| + \left\| \sum_{\mathbf{x} \in \hat{P} \setminus S} \xi_i \right\| \right)^2 \right)^{1/2} \left(\left\| \sum_{\mathbf{x} \in S \cap \hat{Q}} \mathbf{k}_{\mathbf{x}} \right\| \|\mathbf{w} - \mathbf{w}^*\| \right) \quad (63)$$

$$= \frac{1}{np} \sqrt{\frac{p}{1-p}} \left(\left\| \sum_{\mathbf{x} \in \hat{P} \setminus S} \mathbf{k}_{\mathbf{x}} \right\| \|\mathbf{w} - \mathbf{w}^*\| + \left\| \sum_{\mathbf{x} \in \hat{P} \setminus S} \xi_i \right\| \right) \left(\left\| \sum_{\mathbf{x} \in S \cap \hat{Q}} \mathbf{k}_{\mathbf{x}} \right\| \|\mathbf{w} - \mathbf{w}^*\| \right) \quad (64)$$

(a) holds because all entries in $\mathbf{k}_{\mathbf{x}}$ are strictly non-negative.

(b) holds because the points outside of the subquantile will have greater average error than the points within the subquantile, we use this step as we do not make distributional assumptions on \hat{Q} . We also introduce the $\frac{p}{1-p}$ since this is the fraction of points within the subquantile to outside.

(c) holds because of the inequality, $\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x}\|^2 - 2\mathbf{x}^\top \mathbf{y} + \|\mathbf{y}\|^2 \leq \|\mathbf{x}\|^2 + 2\|\mathbf{x}\| \|\mathbf{y}\| + \|\mathbf{y}\|^2 = (\|\mathbf{x}\| + \|\mathbf{y}\|)^2$. We know analyze the LHS and RHS together to upper bound $\|\mathbf{w} - \mathbf{w}^*\|_{\mathcal{H}}$.

$$K \|\mathbf{w} - \mathbf{w}^*\| - \left\| \sum_{\mathbf{x} \in S \cap \hat{P}} \xi_{\mathbf{x}} \mathbf{k}_{\mathbf{x}} \right\| \leq \sqrt{\frac{p}{1-p}} \left(\left\| \sum_{\mathbf{x} \in \hat{P} \cap S} \mathbf{k}_{\mathbf{x}} \right\| \|\mathbf{w} - \mathbf{w}^*\| + \left\| \sum_{i \in \hat{P} \setminus S} \xi_i \right\| \right) \left\| \sum_{\mathbf{x} \in S \cap \hat{Q}} \mathbf{k}_{\mathbf{x}} \right\| + \eta \quad (65)$$

$$\|\mathbf{w} - \mathbf{w}^*\| \leq \frac{\left\| \sum_{\mathbf{x} \in S \cap \hat{P}} \xi_{\mathbf{x}} \mathbf{k}_{\mathbf{x}} \right\| + \sqrt{\frac{p}{1-p}} \left(\left\| \sum_{i \in \hat{P} \setminus S} \xi_i \right\| \left\| \sum_{\mathbf{x} \in S \cap \hat{P}} \mathbf{k}_{\mathbf{x}} \right\| \right) + \eta}{K - \sqrt{\frac{p}{1-p}} \left\| \sum_{\mathbf{x} \in \hat{P} \setminus S} \mathbf{k}_{\mathbf{x}} \right\|} \quad (66)$$

Can we extend the notion of resilience in (Steinhardt u. a., 2018)(Definition 1, Proposition 2), (Awasthi u. a., 2022)(Proposition D.1, D.3, D.4), and (Jambulapati u. a., 2020)(Lemma 8, Corollary 3, Corollary 4) to the kernel feature space $\phi : \mathcal{X} \rightarrow \mathcal{H}$. (Staib und Jegelka, 2019) has some analysis on gaussian kernels which might be useful.

In (Liao u. a., 2020), equation 1 and equation 2, there is an approximation for the Kernel Gram Matrix for Gaussian Kernel, \mathbf{K} . The original idea is in (Rahimi und Recht, 2007). This could be easier to model.

$$[\mathbf{K}_{\mathbf{X}}]_{ij} = e^{-\frac{1}{2}(\|\mathbf{x}_i\|^2 + \|\mathbf{x}_j\|^2)} (\cosh(\mathbf{x}_i^\top \mathbf{x}_j) + \sinh(\mathbf{x}_i^\top \mathbf{x}_j)) \quad (67)$$

This concludes the proof. ■

C Robust Linear Regression

Theorem 31. (*Approximate Stationary Point is Good*). Let $\hat{\mathbf{w}}$ be a η -stationary point as defined in *definition 23*. It then follows:

$$\|\hat{\mathbf{w}} - \mathbf{w}^*\| \leq \frac{\|\sum_{\mathbf{x} \in S \cap \hat{P}} \xi_{\mathbf{x}} \mathbf{k}_{\mathbf{x}}\| + \sqrt{\frac{p}{1-p}} (\|\sum_{i \in \hat{P} \cap S} \xi_i\| \|\sum_{\mathbf{x} \in S \cap \hat{P}} \mathbf{k}_{\mathbf{x}}\|)}{\|\sum_{\mathbf{x} \in S \cap \hat{P}} \mathbf{k}_{\mathbf{x}}\|^2 - \sqrt{\frac{p}{1-p}} \|\sum_{\mathbf{x} \in \hat{P} \setminus S} \mathbf{k}_{\mathbf{x}}\|} \quad (68)$$

where $S^{(T)}$ represents the np set of points in the subquantile.

Proof. We start from the proof of *theorem 25*.

$$\|\hat{\mathbf{w}} - \mathbf{w}^*\| \leq \frac{\|\sum_{i \in S \cap \hat{P}} \xi_i \mathbf{x}_i\| + \sqrt{\frac{p}{1-p}} (\|\sum_{i \in \hat{P} \cap S} \xi_i\| \|\sum_{i \in S \cap \hat{P}} \mathbf{x}_i\|)}{\|\sum_{i \in S \cap \hat{P}} \mathbf{x}_i\|^2 - \sqrt{\frac{p}{1-p}} \|\sum_{i \in \hat{P} \setminus S} \mathbf{x}_i\|} \quad (69)$$

$$\leq \frac{\sigma \varepsilon \log(1/\varepsilon) + \sqrt{\frac{p}{1-p}} (\varepsilon + \varepsilon \sqrt{\log(1/\varepsilon)})}{(1 - \varepsilon \log(1/\varepsilon))^2 - \sqrt{\frac{p}{1-p}} \varepsilon \sqrt{\log(1/\varepsilon)}} \quad (70)$$

■

D Base Learner Algorithm

Algorithm 5: SUBQ-BASE-LEARNER

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

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

```

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

```

Here we can note the similarity of Algorithm ?? 5 to the algorithm described in (Awasthi u. a., 2022). This is because the Trimmed Maximum Likelihood Estimator is equivalent to minimizing over the subquantile of the likelihood.

Remark 32. Define the function $\Psi(t) \triangleq \min_{\mathbf{w}} g(t, \mathbf{w})$

E Experimental Details

Our datasets are synthetic and are sourced from (Dua und Graff, 2017)

Dataset	Dimension d	Sample Size n	Source
Polynomial	3	1000	Ours
Boston Housing	13	506	(Dua und Graff, 2017)
Concrete Data	8	1030	(Dua und Graff, 2017)
Wine Quality	11	1599	(Dua und Graff, 2017)

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.