
ROBUST LINEAR REGRESSION BY SUB-QUANTILE OPTIMIZATION

Arvind Rathnashyam, Fatih Orhan, Joshua Myers, & Jake Herman *

Department of Computer Science

Rensselaer Polytechnic University

Troy, NY 12180, USA

{rathna, orhanf, myersj5, hermaj2}@rpi.edu

ABSTRACT

Robust Linear Regression is the problem of fitting data to a distribution, P when there exists contaminated samples, Q . We consider the Huber Contamination modeled as $\hat{P} = (1 - \varepsilon)P + \varepsilon Q$ where $\varepsilon \in (0, 0.5)$. Traditional Least Squares Methods fit the empirical risk model to all training data in \hat{P} . In this paper we show theoretical and experimental results of sub-quantile optimization, where we optimize with respect to the p -quantile of the empirical loss. Sub-Quantile Optimization theoretically and empirically works in the case of both oblivious and adversarial outliers.

1 INTRODUCTION

Linear Regression is one of the most widely used statistical estimators throughout science. Robustness Learning in High Dimensions on Huber Contamination Models, Huber & Ronchetti (2009), has gained much attention in the last decade, Diakonikolas & Kane (2019). The key motivating factor in investigating robust linear regression is the sheer vastness of probability distributions that are not drawn from a normal distribution schema. Given that outliers in data sets occur so frequent, the ability for a linear regression model to be robust is necessary to compensate for the various distributions being analyzed.

1.1 MOTIVATIONS

The failure of classical regression techniques being unable to model data highly corrupted by outliers can be conveyed clearly in numerous datasets, including those featuring data in the medical, economic, and meteorological fields. Ultimately, in many real data sets, the samples may not be collected from even or fair distributions; thus, classical analyses such as standard regression or least-squares may not represent the actual distribution of the data well.

The quantile is a statistical measure that is distribution-agnostic, this makes it very suitable for robust estimation in the Huber Contamination Model.

1.2 CONTRIBUTIONS

Our goal is to provide a theoretic analysis and convergence conditions for sub-quantile optimization and offer practitioners a method for robust linear regression. Several popular methods have been utilized due

*Work done as a part of ML and Optimization Spring 2023 Group Project.

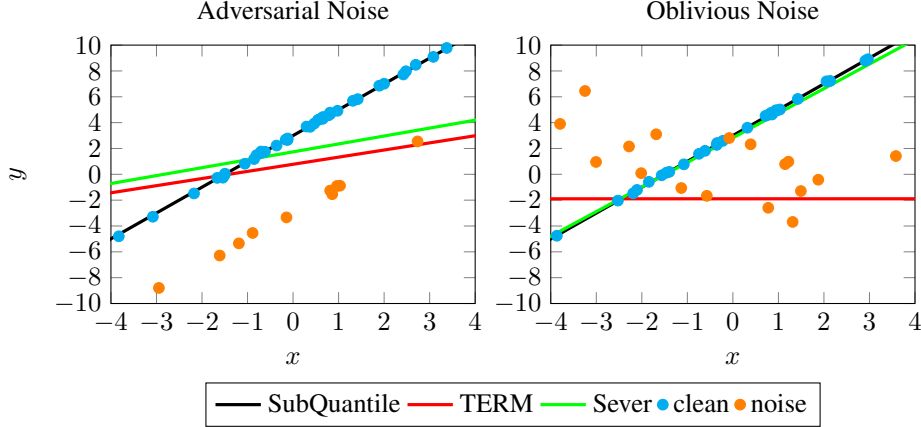


Figure 1: Sub-Quantile Performance on Adaptive Outliers

to their simplicity and high effectiveness including quantile regression Koenker & Hallock (2001), Theil-Sen Estimator Sen (1968), and Huber Regression Huber & Ronchetti (2009). These methods, although rudimentary, serve to show the effectiveness of building resistance against outliers in data. By improving upon existing methods, namely least-squares estimation in these cases, models can be designed to better estimate data sets with considerably corruptive outliers.

Sub-Quantile Optimization aims to address the shortcomings of ERM in applications such as noisy/corrupted data (Khetan et al. (2018), Jiang et al. (2018)), classification with imbalanced classes, (Lin et al. (2017), He & Garcia (2009)), as well as fair learning (Corbett-Davies & Goel (2018)).

As seen in the above comparison, current models fail to estimate data sets corrupted by structured noise, with some models even failing to estimate trends plagued with unstructured noise. Through this, sub-quantile optimization is shown to prevail at overcoming these challenges current models currently face. In Table

Paper	Adversary	Threshold
Sever Diakonikolas et al. (2019)	Adaptive	Gradient of Loss
CRR Bhatia et al. (2017)	Oblivious	
This Paper	Adaptive	Loss

Table 1: A comparison of different iterative thresholding algorithms for Robust Least Squares Regression

2 RELATED WORK

Least Trimmed Squares (LTS) Mount et al. (2014) is an estimator that relies on minimizing the sum of the smallest h residuals given a $(d - 1)$ -dimension hyperplane calculated given n data points in \mathbf{R}^d and an integer trimming parameter h . Given that the outliers comprise less than half the data, this algorithm is more efficient than the more common LMS estimator. However, this algorithm unfortunately suffers from the curse of dimensionality; the computational cost of the algorithm grows exponentially with increasing dimensions of the data. Thus, the necessity to design a more computationally efficient algorithm is expressed.

Tilted Empirical Risk Minimization (TERM) Li et al. (2020) is a framework built to similarly handle the shortcomings of empirical risk minimization (ERM) with respect to robustness. The TERM framework instead minimizes the following quantity, where t is a hyperparameter known as tilt

$$\tilde{R}(t; \theta) := \frac{1}{t} \log \left(\frac{1}{N} \sum_{i \in [N]} e^{t f(\mathbf{x}_i; \theta)} \right) \quad (1)$$

By using the tilt hyperparameter to change the individual impact of each specific loss, the model is more resistant to outliers found in the data.

SMART Awasthi et al. (2022) proposes the *iterative trimmed maximum likelihood estimator* against adversarially corrupted samples in General Linear Models (GLM). The estimator is defined as follows, where $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ represents the training data.

$$\hat{\theta}(S) = \min_{\theta} \min_{\hat{S} \subset S, |\hat{S}|=(1-\epsilon)n} \sum_{(\mathbf{x}_i, y_i) \in \hat{S}} -\log f(y_i | \theta^T \mathbf{x}_i) \quad (2)$$

This estimator is proven to return near-optimal risk on a variety of linear models, including Gaussian regression, Poisson regression, and binomial regression; these achievements can be demonstrated on label and covariate corruptions.

SEVER Diakonikolas et al. (2019) is a gradient filtering algorithm which removes elements whose gradients have the furthest distance from the average gradient of all points

$$\tau_i = \left((\nabla f_i(\mathbf{w}) - \hat{\nabla}) \cdot \mathbf{v} \right)^2 \quad (3)$$

This method is novel in that it is highly scalable, making it robust against high-dimension data with structured outliers. Similarly, SEVER is easily implemented with standard machine learning libraries and can be applied to many typical learning problems, including classification and regression. Despite this, the algorithm still falls short when features have high covariance or when features have low predictive power of the target. Moreover, SEVER requires approximate learners to be run after every iteration, making SEVER unfeasible for large-scale machine learning tasks.

Quantile Regression Yu et al. (2003) relies on splitting data into quantiles to better represent data that is not evenly distributed. The paper introduces various estimation methods for quantile regression and apply them to a multitude of datasets. In doing so, they prove quantile regression is suitable at estimating both linear and nonlinear response models.

Super-Quantile Optimization Rockafellar et al. (2014) aims to solve error minimization problems by building upon the aforementioned quantile regression by centering around a conditional value-at-risk, or a superquantile. For $\alpha \in [0, 1)$, the α -superquantile for a random variable Y is defined as

$$\bar{q}_\alpha(Y) := \frac{1}{1-\alpha} \int_\alpha^1 q_\beta(Y) d\beta \quad (4)$$

In doing so, more conservatively fitted curves are produced. As with quantile regression, such curves do require the solution of a linear program. This concept of superquantile error provides insight into tail behavior for quantities of error and an overall unique approach to linear regression.

Robust Risk Minimization Osama et al. (2020) is a method in which given an upper bound on the corrupted data fraction ϵ , the risk function can be minimized as follows:

$$\hat{\theta}_{RRM} = \operatorname{argmin}_{\theta \in \Theta} \min_{\pi \in \Pi: \mathbf{H}(\pi) \geq \lfloor n(1-\epsilon) \rfloor} R(\theta, \pi) \quad (5)$$

This method is popular as it does not require the removal of corrupted data points and does not rely on a specified corruption fraction.

3 SUB-QUANTILE OPTIMIZATION

Definition 1. Let F_X represent the Cumulative Distribution Function (CDF) of the random variable X . The **p-Quantile** of a Random Variable X is defined as follows

$$Q_p(p) = \inf\{x \in \mathbb{R} : p \leq F(x)\} \quad (6)$$

Definition 2. Let ℓ be the loss function. **Risk** is defined as follows

$$U = \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}} [\ell(f(\mathbf{x}; \boldsymbol{\theta}, y))] \quad (7)$$

The **p-Quantile** of the Empirical Risk is given

$$\mathbb{L}_p(U) = \frac{1}{p} \int_0^p Q_q(U) dq = \mathbb{E}[U|U \leq Q_p(U)] = \max_{t \in \mathbb{R}} \left\{ t - \frac{1}{p} \mathbb{E}[(t - U)^+] \right\} \quad (8)$$

In equation 8, t represents the p -quantile of U . We also show that we can calculate t by a maximizing optimization function. The Sub-Quantile Optimization problem is posed as follows

$$\boldsymbol{\theta}_{SM} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \max_{t \in \mathbb{R}} \left\{ t - \frac{1}{p} \mathbb{E}(t - \ell(f(\mathbf{x}; \boldsymbol{\theta}), y))^+ \right\} \quad (9)$$

For the linear regression case, this equation becomes

$$\boldsymbol{\theta}_{SM} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \max_{t \in \mathbb{R}} \left\{ t - \frac{1}{np} \sum_{i=1}^n (t - (\boldsymbol{\theta}^T \mathbf{x}_i - y_i)^2)^+ \right\} \quad (10)$$

The two-step optimization for Sub-Quantile optimization is given as follows

$$t_{k+1} = \arg \max_t g(t, \boldsymbol{\theta}_k) \quad (11)$$

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \alpha \nabla_{\boldsymbol{\theta}_k} g(t, \boldsymbol{\theta}_k) \quad (12)$$

This algorithm is adopted from Razaviyayn et al. (2020). Theoretically, it has been proven to converge in research by Jin et al. (2019).

Algorithm 1: Sub-Quantile Minimization Gradient Descent

Input: Training iterations T , Quantile p ,

SubQuantile Update: j

Output: Trained Parameters, $\boldsymbol{\theta}_{(T)}$

```

1:  $\boldsymbol{\theta}_{(0)} \leftarrow (X^T X)^{-1} X^T y$ 
2: for  $k \in 1, 2, \dots, T$  do
3:    $\boldsymbol{\nu} \leftarrow (X\boldsymbol{\theta}_{(k)} - y)^2$ 
4:   if  $k \% j = 0$  then
5:      $\hat{\boldsymbol{\nu}} \leftarrow \text{sorted}(\boldsymbol{\nu})$ 
6:      $t_{(k+1)} \leftarrow \hat{\nu}_{np}$ 
7:   end
8:    $S_{(k)} \leftarrow \|\mathbf{x}_i \text{ if } (\boldsymbol{\theta}_{(k)}^T \mathbf{x}_i - y_i)^2 \leq t_{(k+1)}\|$ 
9:    $L_{(k)} \leftarrow \frac{1}{np} \|S^T S\|_2$ 
10:   $\alpha_{(k)} \leftarrow \frac{1}{2L_{(k)}}$ 
11:   $\boldsymbol{\theta}_{(k+1)} \leftarrow \boldsymbol{\theta}_k - \alpha_{(k)} \nabla_{\boldsymbol{\theta}} g(t_{(k+1)}, \boldsymbol{\theta}_{(k)})$ 
12: end
13: return  $\boldsymbol{\theta}_T$ 

```

Algorithm 2: Sub-Quantile Minimization for Ridge Regression

Input: Training Iterations T , Quantile p

Output: Trained Parameters, $\boldsymbol{\theta}_{(T)}$

```

1:  $\boldsymbol{\theta}_{(0)} \leftarrow (X^T X + \lambda I)^{-1} X^T y$ 
2: for  $k \in \{1, 2, \dots, T\}$  do
3:    $\boldsymbol{\nu} \leftarrow (X\boldsymbol{\theta}_{(k)} - y)^2$ 
4:    $t_{(k+1)} \leftarrow \hat{\nu}_{np}$ 
5:    $S_{(k)} \leftarrow \|\mathbf{x}_i \text{ if } (\boldsymbol{\theta}_{(k)}^T \mathbf{x}_i - y_i)^2 \leq t_{(k+1)}\|$ 
6:    $\boldsymbol{\theta}_{(k+1)} \leftarrow (S^T S + \lambda I)^{-1} S^T y_S$ 
7: end
8: return  $\boldsymbol{\theta}_T$ 

```

3.1 MOTIVATION

Assumption 1. To provide theoretical bounds on the effectiveness of Sub-Quantile Minimization, we make the General Linear Model Assumption that

$$\mathbf{y}_P = P\beta_P + \epsilon_P \quad (13)$$

and similarly

$$\mathbf{y}_Q = Q\beta_Q + \epsilon_Q \quad (14)$$

where β_P and β_Q the oracle regressors for \mathbb{P} and \mathbb{Q} and ϵ_P and ϵ_Q are both Normally Distributed with mean 0.

Since we are interested in learning the optimal model for distributions, our goal is to learn the parameters β_P from the distribution \hat{P} . We want to clarify the corruption is not adversarially chosen.

In this section we quantify the effect of corruption on the desired model. To introduce notation, let P represent the data from distribution \mathbb{P} and let Q represent the training data for \mathbb{Q} . Let \mathbf{y}_P represent the target data for \mathbb{P} and let \mathbf{y}_Q represent the target data for \mathbb{Q} .

Assumption 2. We assume the rows of P and Q are sampled from the same multivariate normal distribution.

$$P_i, Q_j \sim \mathcal{N}_p(\mathbf{0}, \Sigma) \quad (15)$$

We will use our assumptions to quantify the effect of the corrupted data on an optimal least squares regression model. We are interested in $(X^T X)^{-1} X^T y - (P^T P)^{-1} P^T y$. It is known the least squares optimal solution for X is equal to $(X^T X)^{-1} X^T y$

Note $X = \begin{pmatrix} P \\ Q \end{pmatrix}$ and $y = \begin{pmatrix} \mathbf{y}_P \\ \mathbf{y}_Q \end{pmatrix}$ so $X^T = (P^T \quad Q^T)$

Theorem 1. The expected optimal parameters of the corrupted model $\hat{\mathbb{P}}$

$$\mathbb{E}[X^\dagger y] = \beta_P + \epsilon(\beta_Q - \beta_P) \quad (16)$$

The proof is reliant on assumption 2, this allows us to utilize the Wishart Distribution, \mathcal{W} , and the inverse Wishart Distribution, \mathcal{W}^{-1} . Please refer to Appendix B.1. By Theorem 1 we can see the level of corruption is dependent upon ϵ , which represents the percentage of corrupted samples, and the distance between the optimal parameters for \mathbb{P} , which is β_P and the optimal parameters for \mathbb{Q} , which is β_Q .

Here we utilize the idea of *influence* from McWilliams et al. (2014).

Theorem 1 finds the optimal model when the corrupted distribution is sampled from the same distribution as the target distribution but has different optimal parameters. We will now look at the case of feature corruption. This is where the optimal parameters of the two distributions are the same but the data from \mathbb{P} and \mathbb{Q} are sampled differently.

In equation 16, note as $\epsilon \rightarrow 0$ we are returned β_P . This is the intuition behind SubQuantile Minimization. By minimizing over the SubQuantile, we seek to reduce ϵ , and thus our model will return a model which is by expectation closer to β_P .

4 THEORY

4.1 ANALYSIS OF $g(t, \theta)$

In this section, we will explore the fundamental aspects of $g(t, \theta)$. This will motivate the convergence analysis in the next section.

Lemma 1.1. $g(t_{k+1}, \theta_k)$ is concave with respect to t .

Proof. We provide a simple argument for concavity. Note t is a concave and convex function. Also $(\cdot)^+$ is a convex strictly non-negative function. Therefore we have a concave function minus the non-negative multiple of a summation of an affine function composed with a convex function. Therefore this is a concave function with respect to t . \square

Lemma 1.2. The maximizing value of t in $g(t, \theta)$ in t -update step of optimization as described by Equation 11 is maximized when $t = Q_p(U)$

Proof. Since $g(t, \theta)$ with respect to t is a concave function. Maximizing $g(t, \theta)$ is equivalent to minimizing $-g(t, \theta)$. We will find fermat's optimality condition for the function $-g(t, \theta)$, which is convex. Let $\hat{\nu} = \text{sorted}((\theta^T X - y)^2)$ and note $0 < p < 1$

$$\partial(-g(t, \theta)) = -1 + \frac{1}{np} \sum_{i=1}^n \begin{cases} 1, & \text{if } t > \hat{\nu}_i \\ 0, & \text{if } t < \hat{\nu}_i \\ [0, 1], & \text{if } t = \hat{\nu}_i \end{cases} \quad (17)$$

$$= 0 \text{ when } t = \hat{\nu}_{np} \quad (18)$$

This is the p -quantile of U . A full proof is provided in Appendix D.1. \square

Lemma 1.3. Let $t = \hat{\nu}_{np}$. The θ -update step described in Equation 10 is equivalent to minimizing the least squares loss of the np elements with the lowest squared loss.

$$\nabla_{\theta} g(t_{k+1}, \theta_k) = \frac{1}{np} \sum_{i=1}^{np} 2\mathbf{x}_i(\theta_k^T \mathbf{x}_i - y_i) \quad (19)$$

We provide a proof in Appendix D.2. However, this result is quite intuitive as it shows we are optimizing over the p Sub-Quantile of the Risk.

Interpretation 1. Sub-Quantile Minimization continuously minimizes the risk over the p -quantile of the error. In each iteration, this means we reduce the error of the points within the lowest np errors.

Lemma 1.4. $g(t_{k+1}, \theta_k)$ is convex with respect to θ_k .

Proof. We see by lemma 1.2 and interpretation 1, we are optimizing by the np points with the lowest squared error. Mathematically,

$$\begin{aligned} g(t_{k+1}, \theta_k) &= t_{k+1} - \frac{1}{np} \sum_{i=1}^n (t_{k+1} - (\theta^T \mathbf{x}_i - y_i)^2)^+ \\ &= t - t + \frac{1}{np} \sum_{i=1}^{np} (\theta^T \mathbf{x}_i - y_i)^2 \\ &= \frac{1}{np} \sum_{i=1}^{np} (\theta^T \mathbf{x}_i - y_i)^2 \end{aligned}$$

Now we can make a simple argument for convexity. We have a non-negative multiple of the sum of the composition of an affine function with a convex function. Thus $g(t, \theta)$ is convex with respect to θ . \square

Lemma 1.5. $g(t, \theta)$ is L -smooth with respect to θ with $L = \left\| \frac{2}{np} \sum_{i=1}^{np} \|x_i\|^2 \right\|$

Now we will state two properties regarding the effect of the t -update step and the θ -update step as described in Equations 11 and 12, respectively.

Lemma 1.6. If $t_{k+1} \leq t_k$ then $g(t_{k+1}, \theta_k) = g(t_k) + \frac{1}{np} \sum_{i=np}^n (t_k - \nu_i)^+$. If $t_{k+1} > t_k$, then $g(t_{k+1}, \theta_k) = g(t_k) + \frac{1}{np} \sum_{i=n(p-\delta)}^{np} (t - \nu_i)^+ - \delta t$. For a small δ .

Proof Sketch. When $t_{k+1} \leq t_k$ this result is quite intuitive, as we are simply removing the error of the elements outside elements within the lowest np squared losses. We delegate the rest of the proof to Appendix D.4 \square

4.2 OPTIMIZATION

We are solving a min-max convex-concave problem, thus we are looking for a Nash Equilibrium Point.

Definition 3. (t^*, θ^*) is a **Nash Equilibrium** of g if for any $(t, \theta) \in \mathbb{R} \times \mathbb{R}^d$

$$g(t^*, \theta) \leq g(t^*, \theta^*) \leq g(t, \theta^*) \quad (20)$$

Definition 4. (t^*, θ^*) is a **Local Nash Equilibrium** of g if there exists $\delta > 0$ such that for any t, θ (t, θ) satisfying $\|t - t^*\| \leq \delta$ and $\|\theta - \theta^*\| \leq \delta$ then:

$$g(t^*, \theta) \leq g(t^*, \theta^*) \leq g(t, \theta^*) \quad (21)$$

Proposition 1. As g is first-order differentiable, any local Nash Equilibrium satisfies $\nabla_{\theta} g(t, \theta) = \mathbf{0}$ and $\nabla_t g(t, \theta) = 0$

We are now interested in what it means to be at a Local Nash Equilibrium. By Proposition 1, this means both first-order partial derivatives are equal to 0. By lemma 1.2, we have shown $\nabla_t g(t, \theta) = 0$ when $\nu_{np} \leq t < \nu_{np+1}$. Furthermore, by lemma 1.3, we have shown $\nabla_{\theta} g(t, \theta) = 0$ when the least squares error is minimized for the np points with lowest squared error. In other words:

$$\mathbb{E} [\nabla_{\theta} g(t_{k+1}, \theta_k)] = 0$$

$$2(\mu\mu^T + \Sigma)(\theta_k - (1 - \varepsilon)\beta_P - \varepsilon\beta_Q) = 0$$

Since the first term is non-zero, the equality is satisfied when:

$$\begin{aligned} (\theta_k - (1 - \varepsilon)\beta_P - \varepsilon\beta_Q) &= 0 \\ \theta_k &= (1 - \varepsilon)\beta_P + \varepsilon\beta_P \end{aligned}$$

Note this aligns with the results of Theorem 1. This means that for a subset of np points from X , the least squares error is minimized. What we are interested in is how many points within those np points come from \mathbb{P} and how many of those points from \mathbb{Q} . Our goal is to minimize the number of points within the np lowest squared losses from Q , as they will introduce error to our predictions on points from P .

Theorem 2. Let $g(t, \theta)$ be differentiable and $g(t, \theta)$ be L -smooth in θ . Let $t_{(k)}$ and $\theta_{(k)}$ be iterates from algorithm 2. Then, $\lim_{k \rightarrow \infty} \mathbb{E} [\|\nabla_{\theta} g(t_{(k+1)}, \theta_{(k)})\|] = 0$.

Proof Sketch. From the intuitions we have gained on Subquantile Minimization. We can state the following:

$$\arg \min_{\theta \in \mathbb{R}^d} \min_{S \in \Pi} \|\theta^T S - y_S\|_2^2 \iff \arg \min_{\theta \in \mathbb{R}^d} \max_{t \in \mathbb{R}} \|\theta^T S - y_S\|_2^2 \quad (22)$$

where Π represents the entire distribution of Subquantile matrices. \square

4.3 CONVERGING TO β_P

We will start by defining the two types of noise we are interesting in.

Definition 5. *Oblivious Noise* is noise that is not dependent on the input data, i.e., $\mathbb{P}[y|X] = \mathbb{P}[y]$

Definition 6. *Adaptive Noise* is noise that is made from a linear combination of the input data, i.e. $y = \beta_Q X + \epsilon$

Also note we often consider Gaussian Noise as Oblivious Noise, but it can be modeled as Adaptive Noise where $\beta_Q = \mathbf{0}$.

Lemma 2.1. *The expected value of error on points in \mathbb{P} will be lower than the expected value of error on points in \mathbb{Q} if $\text{proj}_{\beta_P}(\theta) - \beta_P < \text{proj}_{\beta_Q}(\theta) - \beta_Q$*

Lemma 2.1 gives us an intuitive result, the proof is in appendix ?? . If in each optimization step, our projection on β_P is closer than our projection on β_Q , we know the number of steps from \mathbb{Q} will increase from the previous iteration.

Theorem 3. *After an optimization step,*

$$|\text{proj}_{\beta_P} \theta_{(t+1)} - 1| - |\text{proj}_{\beta_P} \theta_{(t)} - 1| < |\text{proj}_{\beta_Q} \theta_{(t+1)} - 1| - |\text{proj}_{\beta_Q} \theta_{(t)} - 1|$$

if the following holds

$$\left\| \left((1 - \varepsilon^{(t)}) - \alpha_1 \right) \beta_P \right\| > \left\| \left(\varepsilon^{(t)} - \alpha_2 \right) \beta_Q \right\| \quad (23)$$

where α_1 and α_2 represents the coefficients for the linear combination of θ in the basis defined as $\mathbf{B} = [\beta_P \ \beta_Q \ \mathbf{R}]$

We are also interested in Algorithm 2, so we provide theory for it.

Corollary 3.1. *Theorem 2 by expectation converges to a good solution.*

Proof Sketch. Let us first note $\varepsilon^{(0)} < 0.5$, and by Theorem 1, $\theta_{(0)} = (1 - \varepsilon^{(0)})\beta_P + \varepsilon^{(0)}\beta_Q$. It thus follows:

$$\|\theta_{(0)} - \beta_P\| = \|\varepsilon^{(0)}\beta_P + \varepsilon^{(0)}\beta_Q\| \quad (24)$$

and similarly

$$\|\theta_{(0)} - \beta_Q\| = \|(1 - \varepsilon^{(0)})\beta_P + (1 - \varepsilon^{(0)})\beta_Q\| \quad (25)$$

Therefore in the linear case

$$\|\varepsilon^{(0)}\beta_P + \varepsilon^{(0)}\beta_Q\| < \|(1 - \varepsilon^{(0)})\beta_P + (1 - \varepsilon^{(0)})\beta_Q\| \quad (26)$$

Then if $\text{Var}(\epsilon_P) \leq \text{Var}(\epsilon_Q)$, it follows

$$\mathbb{E} \left[\left(\theta_{(0)}^T \mathbf{p} - y_i \right)^2 \right] \leq \mathbb{E} \left[\left(\theta_{(0)}^T \mathbf{q} - y_i \right)^2 \right]$$

This concludes the proof. \square

Note theorem 3.1 is enough to show the final subquantile matrix will have more elements from P . It thus follows, in the case of no feature noise, by expectation, θ will move towards β_P faster than it moves to β_Q . We instead calculate the probability in theorem ??.

We are now interested in theoretical guarantees with no distributional assumptions. First we will consider some intuition on why this problem is not as hard as compared to when the corruption is linearly structured. Let us say the noise is of non-linear regression, in other words $\mathbf{y}_Q \sim f(X, \beta_Q)$, where f is a non-linear combination of features of X . In this case, it is not possible to model the non-linear regression by a linear combination of the features, thus, if we have elements from \mathbb{Q} within the lowest np losses, then training on these points will not generalize well to points from \mathbb{Q} , so their error will not decrease.

5 EMPIRICAL RESULTS

We also present a batch algorithm which improves training speed significantly. In accordance with Minibatch theory, if the subset I of all data is representative of all the data, then this will have similar results to Algorithm 1.

5.1 SYNTHETIC DATA

We now demonstrate SubQuantile Regression in the presence of Gaussian Random Noise.

In our first synthetic experiment, we run Algorithm 1 on synthetically generated structured linear regression data, the noise is sampled from a linear distribution that is dependent on the vector of X . Our results show the near optimal performance of Sub-Quantile Minimization. The results and comparison with other methods can be seen in Table 3. We see in Table 3, Sub-Quantile Minimization produces State of the Art Results in the Quadratic Regression Case. Furthermore, it performs significantly better than baseline methods in the high-noise regimes ($\epsilon = 0.4$), this is confirmed in both the small data and large data datasets. Please refer to Appendix I for more details on the Structured Linear Regression Dataset.

5.2 REAL DATA

We provide results on the Drug Discovery Dataset in Diakonikolas et al. (2019) utilizing the noise procedure described in Li et al. (2020). For each algorithm, if possible we use Ridge Regression, else we use typical least squares. SubQuantile Minimization, ERM, RANSAC, SEVER, TERM, and SMART are all capable of Ridge Regression.

As we can see in Table 2, we obtain state of the art results throughout all noise regimes. This makes our model the strongest among the tested, due to our strength throughout the whole range of noises. This dataset is also

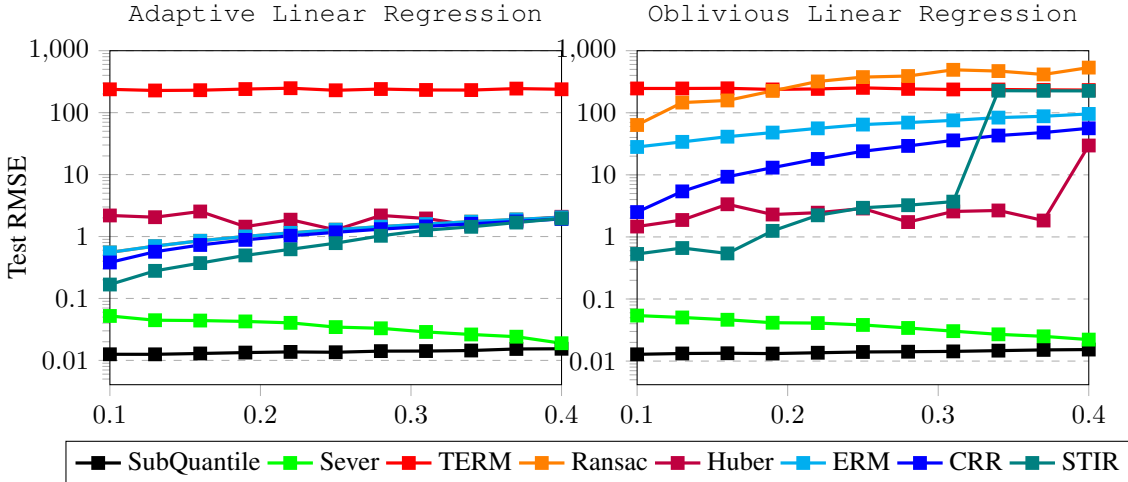


Figure 2: Structured Linear Regression & Noisy Linear Regression Datasets

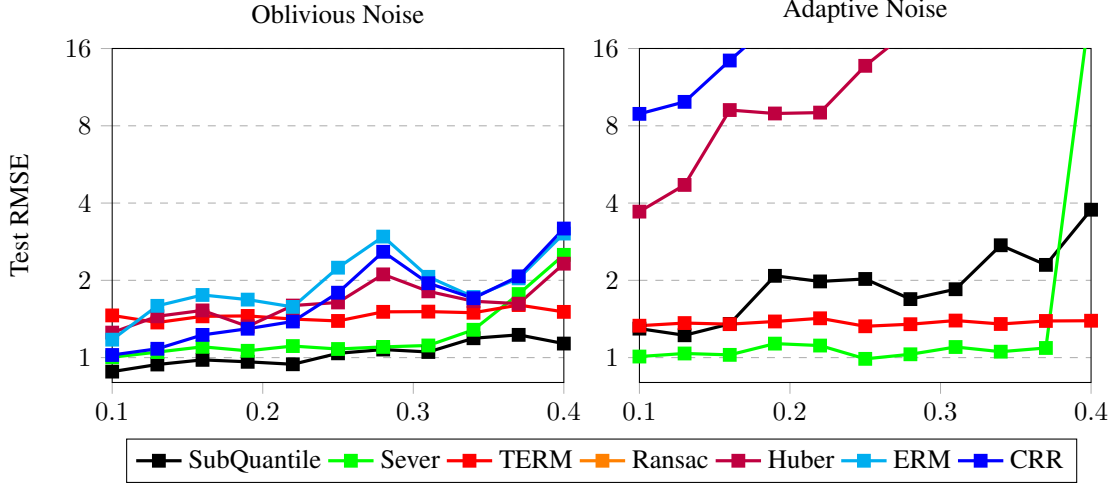


Figure 3: Drug Discovery Dataset with Normal Noise and Structured Noise

Objectives	Test RMSE (Drug Discovery)			
	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon = 0.4$
ERM	1.303 _(0.0665)	1.790 _(0.0849)	2.198 _(0.0645)	2.623 _(0.1010)
CRR Bhatia et al. (2017)	1.079 _(0.0899)	1.125 _(0.0832)	1.385 _(0.1372)	1.725 _(0.1136)
STIR Mukhoty et al. (2019)	1.087 _(0.1256)	1.167 _(0.0750)	1.403 _(0.0987)	1.668 _(0.1142)
Robust Risk Osama et al. (2020)	1.176 _(0.1110)	1.336 _(0.1882)	1.437 _(0.1723)	1.800 _(0.0820)
SMART Awasthi et al. (2022)	1.094 _(0.1065)	1.323 _(0.0758)	1.578 _(0.0799)	1.984 _(0.2020)
TERM Li et al. (2020)	1.029 _(0.0707)	1.126 _(0.0776)	1.191 _(0.1091)	1.201 _(0.1409)
SEVER Diakonikolas et al. (2019)	1.011 _(0.0838)	1.067 _(0.0457)	1.071 _(0.0807)	1.138 _(0.1162)
Huber Huber & Ronchetti (2009)	1.412 _(0.0474)	1.501 _(0.2918)	2.231 _(0.9054)	2.247 _(1.0399)
RANSAC Fischler & Bolles (1981)	1.238 _(0.0529)	1.643 _(0.1331)	2.092 _(0.1935)	2.679 _(0.1365)
SubQuantile($p = 1 - \epsilon$)	0.966 _(0.1119)	1.002 _(0.1025)	1.010 _(0.0630)	1.089 _(0.1129)
Genie ERM	0.960 _(0.0845)	0.982 _(0.0842)	1.006 _(0.0879)	1.030 _(0.0578)

Table 2: Drug Discovery Dataset. Empirical Risk over P with oblivious noise

6 CONCLUSION

In this work we provide a theoretical analysis for robust linear regression by minimizing the *Sub-Quantile* of the Empirical Risk. Furthermore, we run various numerical experiments and compare against the current State of the Art in Robust Linear Regression. Since minimizing over the subquantile is a general machine learning framework, it is scalable to larger scale machine learning problems. In future work, more real world applications can be explored and the theory can be expanded beyond linear regression.

REFERENCES

- Pranjal Awasthi, Abhimanyu Das, Weihao Kong, and Rajat Sen. Trimmed maximum likelihood estimation for robust learning in generalized linear models, 2022. URL <https://arxiv.org/abs/2206.04777>.
- 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 (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/e702e51da2c0f5be4dd354bb3e295d37-Paper.pdf.
- Sam Corbett-Davies and Sharad Goel. The measure and mismeasure of fairness: A critical review of fair machine learning. *CoRR*, abs/1808.00023, 2018. URL <http://arxiv.org/abs/1808.00023>.
- Ilias Diakonikolas and Daniel M. Kane. Recent advances in algorithmic high-dimensional robust statistics. *ArXiv*, abs/1911.05911, 2019.
- 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*, pp. 1596–1606. JMLR, Inc., 2019.
- Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- 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. ISSN 0001-0782. doi: 10.1145/358669.358692. URL <https://doi.org/10.1145/358669.358692>.
- Haibo He and Edwardo A. Garcia. Learning from imbalanced data. *IEEE Transactions on Knowledge and Data Engineering*, 21(9):1263–1284, 2009. doi: 10.1109/TKDE.2008.239.
- Peter J. Huber and Elvezio Ronchetti. *Robust statistics*. Wiley series in probability and statistics. Wiley, Hoboken, N.J., 2nd ed. edition, 2009. URL <http://catdir.loc.gov/catdir/toc/ecip0824/2008033283.html>.
- 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.
- Chi Jin, Praneeth Netrapalli, and Michael I. Jordan. What is local optimality in nonconvex-nonconcave minimax optimization?, 2019. URL <https://arxiv.org/abs/1902.00618>.
- Ashish Khetan, Zachary C. Lipton, and Anima Anandkumar. Learning from noisy singly-labeled data. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=H1sUHgb0Z>.
- Roger Koenker and Kevin F. Hallock. Quantile regression. *Journal of Economic Perspectives*, 15(4):143–156, December 2001. doi: 10.1257/jep.15.4.143. URL <https://www.aeaweb.org/articles?id=10.1257/jep.15.4.143>.
- Tian Li, Ahmad Beirami, Maziar Sanjabi, and Virginia Smith. Tilted empirical risk minimization. *arXiv preprint arXiv:2007.01162*, 2020.

-
- T. Lin, P. Goyal, R. Girshick, K. He, and P. Dollar. Focal loss for dense object detection. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 2999–3007, Los Alamitos, CA, USA, oct 2017. IEEE Computer Society. doi: 10.1109/ICCV.2017.324. URL <https://doi.ieeecomputersociety.org/10.1109/ICCV.2017.324>.
- Brian McWilliams, Gabriel Krummenacher, Mario Lucic, and Joachim M. Buhmann. Fast and robust least squares estimation in corrupted linear models. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 1, NIPS’14*, pp. 415–423, Cambridge, MA, USA, 2014. MIT Press.
- David M. Mount, Nathan S. Netanyahu, Christine D. Piatko, Ruth Silverman, and Angela Y. Wu. On the least trimmed squares estimator. *Algorithmica*, 69(1):148–183, 2014. doi: 10.1007/s00453-012-9721-8. URL <https://doi.org/10.1007/s00453-012-9721-8>.
- 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 (eds.), *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pp. 313–322. PMLR, 16–18 Apr 2019. URL <https://proceedings.mlr.press/v89/mukhoty19a.html>.
- 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.
- Kelley Pace and Ronald Barry. Sparse spatial autoregressions. *Statistics ‘I& Probability Letters*, 33(3):291–297, 1997. URL <https://EconPapers.repec.org/RePEc:eee:stapro:v:33:y:1997:i:3:p:291-297>.
- Meisam Razaviyayn, Tianjian Huang, Songtao Lu, Maher Nouiehed, Maziar Sanjabi, and Mingyi Hong. Non-convex min-max optimization: Applications, challenges, and recent theoretical advances, 06 2020.
- 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. ISSN 0377-2217. doi: <https://doi.org/10.1016/j.ejor.2013.10.046>. URL <https://www.sciencedirect.com/science/article/pii/S0377221713008692>.
- Pranab Kumar Sen. Estimates of the regression coefficient based on kendall’s tau. *Journal of the American Statistical Association*, 63(324):1379–1389, 1968. doi: 10.1080/01621459.1968.10480934. URL <https://www.tandfonline.com/doi/abs/10.1080/01621459.1968.10480934>.
- D.D Wackerly, W. Mendenhall, and R.L. Scheaffer. *Mathematical Statistics with Applications, 7th Edition*. Thompson Learning, Inc., USA, 2008.
- Keming Yu, Zudi Lu, and Julian Stander. Quantile regression: applications and current research areas. *Journal of the Royal Statistical Society: Series D (The Statistician)*, 52(3):331–350, 2003. doi: <https://doi.org/10.1111/1467-9884.00363>. URL <https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/1467-9884.00363>.

A	Linear Algebra and Probability Theory Preliminaries	14
B	Theory for Ridge Regression Algorithm 2	15
B.1	Proof of Theorem 1	15
C	Proofs for $\varepsilon^{(t)}$	17
D	Theory for Subquantile Minimization	19
D.1	Derivation of Lemma 1.2	19
D.2	Derivation of Lemma 1.3	19
D.3	Derivation of Lemma 1.5	19
D.4	Proof of Lemma 1.6	20
E	Theory for Adaptive Linear Corruption	22
E.1	Proof of Theorem 3	22
E.2	Proof of Theorem ??	23
F	Proofs for Convergence	25
F.1	Proof of Theorem 2	25
F.2	Expectation of Improvement	25
G	Stochastic Sub-Quantile Optimization	27
H	Additional Experiments	28
H.1	Quadratic Regression	28
H.2	Abalone	28
H.3	Cal-Housing	28
I	Experimental Details	29
I.1	Adaptive Linear Regression Dataset	29
I.2	Oblivious Linear Regression Dataset	29
I.3	Quadratic Regression Dataset	29
I.4	Drug Discovery Dataset	30
I.5	Feature Noise	30

A LINEAR ALGEBRA AND PROBABILITY THEORY PRELIMINARIES

Fact 1. *The spectral norm of a matrix, A , an $(m \times n)$ matrix, is defined as follows*

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)} = \sigma_{\max}(A) \quad (27)$$

It similarly follows:

$$\|A^T A\|_2 = \|A\|_2^2 \quad (28)$$

Fact 2. Weyl's Inequality states the following:

If M , N , and R are $n \times n$ Hermitian Matrices with the following eigenvalues where $M = N + R$:

$$\begin{aligned} M : \mu_1 \geq \dots \geq \mu_n \\ N : \nu_1 \geq \dots \geq \nu_n \\ R : \rho_1 \geq \dots \geq \rho_n \end{aligned}$$

Then the following equalities hold:

$$\nu_i + \rho_n \leq \mu_i \leq \nu_i + \rho_1 \text{ for } i = 1, \dots, n$$

Fact 3. *Let A be a $n \times m$ matrix with $n \gg m$. It then follows:*

$$A^T A = (U \Sigma V^T)^T (U \Sigma V^T) = (V \Sigma^T U^T) (U \Sigma V^T) = V \Sigma^T \Sigma V^T = V D V^T \quad (29)$$

where $D = \Sigma^T \Sigma = \begin{pmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & \\ & & \ddots & \\ & & & \sigma_m^2 \end{pmatrix}$

Fact 4. *This is a restatement from Wackerly et al. (2008).*

Let X_1, \dots, X_n be i.i.d continuous random variables with common distribution function $F(y)$ and common probability density function $f(x)$. If $X_{(k)}$ denotes the k th-order statistic, then the density function of $X_{(k)}$ is given by:

$$g_{(k)}(x_k) = \frac{n!}{(k-1)!(n-k)!} [F(x_k)]^{k-1} [1 - F(x_k)]^{n-k} f(x_k) \quad (30)$$

Fact 5. *The cdf of the k th order statistic from a sample of n is:*

$$F_{(k,n)} = \mathbb{P}[X_{(k)} \leq x] = \sum_{j=k}^n \binom{n}{j} (1 - F(x))^{n-j} F(x)^j \quad (31)$$

B THEORY FOR RIDGE REGRESSION ALGORITHM 2

B.1 PROOF OF THEOREM 1

Proof.

Assumption 3. The rows of P and Q are sampled from $\mathbf{0}$ centered Normal Distributions.

$$\begin{aligned} P_i &\sim \mathcal{N}(\mathbf{0}, \Sigma_P) \\ Q_i &\sim \mathcal{N}(\mathbf{0}, \Sigma_Q) \end{aligned} \quad (32)$$

Assumption 4. By assumption 3, it thus follows that the matrices $P^T P$ and $Q^T Q$ are sampled from Wishart Distributions.

$$P^T P \sim \mathcal{W}(n, \Sigma_P) \quad (33)$$

$$Q^T Q \sim \mathcal{W}(n, \Sigma_Q) \quad (34)$$

Assumption 5. Similar to the assumption made in Bhatia et al. (2017), to give theoretical bounds on the our algorithm, we assume the following:

$$\Sigma_P = \xi_P I \quad (35)$$

$$\Sigma_Q = \xi_Q I \quad (36)$$

where $\xi_P, \xi_Q \geq 0$

The closed form solution for Ridge Regression with regularization parameter λ is equal to the following:

$$\hat{\beta} = (X^T X + \lambda I)^{-1} X^T y \quad (37)$$

We will use this to bound the difference of the β_P and $\theta_{(k)}$.

$$\|\beta_P - \theta_{(k)}\|_2 = \|\beta_P - (S^T S + \lambda I)^{-1} S^T y\|_2$$

Note the subquantile matrix S , consists of data points from P and Q , we will reorganize S into the following:

$$S = \begin{pmatrix} P \\ Q \end{pmatrix} = \begin{pmatrix} \leftarrow & \mathbf{p}_1 & \rightarrow \\ \vdots & \vdots & \vdots \\ \leftarrow & \mathbf{p}_{n(1-\varepsilon^{(t)})} & \rightarrow \\ \leftarrow & \mathbf{q}_1 & \rightarrow \\ \vdots & \vdots & \vdots \\ \leftarrow & \mathbf{q}_{n\varepsilon^{(t)}} & \rightarrow \end{pmatrix} \text{ and } \mathbf{y} = \begin{pmatrix} \mathbf{y}_P \\ \mathbf{y}_Q \end{pmatrix} = \begin{pmatrix} \beta_P \mathbf{p}_1 + \epsilon_P \\ \vdots \\ \beta_P \mathbf{p}_{n(1-\varepsilon^{(t)})} + \epsilon_P \\ \beta_Q \mathbf{q}_1 + \epsilon_Q \\ \vdots \\ \beta_Q \mathbf{q}_{n\varepsilon^{(t)}} + \epsilon_Q \end{pmatrix}$$

Let us also assume $\text{Var}(\epsilon_P) = \eta_P$ and $\text{Var}(\epsilon_Q) = \eta_Q$. Then we can make the following manipulations:

$$\begin{aligned} &= \left\| \beta_P - (P^T P + Q^T Q + \lambda I)^{-1} (P^T Q^T) \begin{pmatrix} \mathbf{y}_P \\ \mathbf{y}_Q \end{pmatrix} \right\|_2 \\ &= \left\| (P^T P + \lambda I)^{-1} P^T \mathbf{y}_P - (P^T P + Q^T Q + \lambda I)^{-1} P^T \mathbf{y}_P - (P^T P + Q^T Q + \lambda I)^{-1} Q^T \mathbf{y}_Q \right\|_2 \\ &\leq \left\| (P^T P + \lambda I)^{-1} P^T \mathbf{y}_P - (P^T P + Q^T Q + \lambda I)^{-1} P^T \mathbf{y}_P \right\|_2 + \left\| (P^T P + Q^T Q + \lambda I)^{-1} Q^T \mathbf{y}_Q \right\|_2 \\ &\leq \left\| (P^T P + \lambda I)^{-1} \right\|_2 \|P^T \mathbf{y}_P\|_2 + \left\| (P^T P + Q^T Q + \lambda I)^{-1} \right\|_2 \|P^T \mathbf{y}_P\|_2 + \left\| (P^T P + Q^T Q + \lambda I)^{-1} \right\|_2 \|Q^T \mathbf{y}_Q\|_2 \\ &\leq \sqrt{\lambda_{\max}(P^T P)} \left(\left\| (P^T P + \lambda I)^{-1} \right\|_2 + \left\| (P^T P + Q^T Q + \lambda I)^{-1} \right\|_2 \right) \|\mathbf{y}_P\|_2 + \sqrt{\lambda_{\max}(Q^T Q)} \left\| (P^T P + Q^T Q + \lambda I)^{-1} \right\|_2 \|\mathbf{y}_Q\|_2 \\ &= \frac{\sigma_{\max}(P) \|\mathbf{y}_P\|_2}{\sqrt{\lambda_{\max}(P^T P + \lambda I)}} + \frac{\sigma_{\max}(P) \|\mathbf{y}_P\|_2}{\sqrt{\lambda_{\max}(P^T P + Q^T Q + \lambda I)}} + \frac{\sigma_{\max}(Q) \|\mathbf{y}_Q\|_2}{\sqrt{\lambda_{\max}(P^T P + Q^T Q + \lambda I)}} \\ &\stackrel{(a)}{\leq} \frac{2\sigma_{\max}(P) \|P\beta_P + \epsilon_P\|_2 + \sigma_{\max}(Q) \|Q\beta_Q + \epsilon_Q\|_2}{\sqrt{\lambda_{\max}(P^T P + \lambda I)}} \\ &\stackrel{(b)}{\leq} \frac{2\sigma_{\max}^2(P) \|\beta_P\|_2 + 4\sigma_{\max}(P)n(1-\varepsilon^{(t)})\eta_P + 2\sigma_{\max}^2(Q) \|\beta_Q\|_2 + 4\sigma_{\max}(Q)n\varepsilon^{(t)}\eta_Q}{\sqrt{\lambda_{\max}(P^T P) + \lambda_{\min}(\lambda I)}} \\ &\leq \frac{2\sigma_{\max}^2(P) \|\beta_P\|_2 + 4\sigma_{\max}(P)n(1-\varepsilon^{(t)})\eta_P + 2\sigma_{\max}^2(Q) \|\beta_Q\|_2 + 4\sigma_{\max}(Q)n\varepsilon^{(t)}\eta_Q}{\sqrt{\sigma_{\max}^2(P) + \lambda}} \end{aligned}$$

(a) is due to $Q^T Q$ being a positive semi definite symmetric matrix.

(b) holds with probability $(1 - (\Phi(2) - \Phi(-2)))^{n(1-\varepsilon^{(t)})}$

Thus we have shown that $\|\beta_P - \theta_{(t)}\|_2$ is bounded above at any time-step (t) in terms of the maximal singular values of the data matrices and the variance of the white noise.

□

C PROOFS FOR $\varepsilon^{(t)}$

Let us note the Subquantile Matrix $S = \begin{pmatrix} P \\ Q \end{pmatrix} = \begin{pmatrix} \leftarrow & \mathbf{p}_1 & \rightarrow \\ \vdots & \vdots & \vdots \\ \leftarrow & \mathbf{p}_{\eta_P} & \rightarrow \\ \leftarrow & \mathbf{q}_1 & \rightarrow \\ \vdots & \vdots & \vdots \\ \leftarrow & \mathbf{q}_{\eta_Q} & \rightarrow \end{pmatrix}$. Thus we can define: $\varepsilon^{(t)} = \frac{\eta_Q}{\eta_P + \eta_Q} = \frac{\eta_Q}{np}$. Where

n is the number of training examples and p is the Subquantile we are minimizing over. In this section, we want to provide a theoretical upper bound on η_Q , from where we can upper bound $\varepsilon^{(t)}$ which is stronger than the trivial upper bound of $\frac{\epsilon}{p}$. We will approach this problem with Order Statistics.

Assumption 6. *The optimal regressors are linearly independent, i.e. $\beta_P \neq \gamma\beta_Q \forall \gamma \in \mathbb{R}$*

Let $\theta_{(t)} = \alpha_1\beta_P + \alpha_2\beta_Q + \sum_{i=3}^d \mathbf{r}_i$.

Let us denote $P_1 < P_2 < \dots < P_{n(1-\epsilon)}$ as the order statistics of the random variable $P \sim (\mathbf{p}_i\theta - (\mathbf{p}_i\beta_P + \epsilon_P))^2$ where $\mathbf{p}_i \sim \mathcal{N}_d(\mathbf{0}, \xi_P I)$. Let us also denote $Q_1 < Q_2 < \dots < Q_{n(1-\epsilon)}$ as the order statistics of the random variable $Q \sim (\mathbf{q}_i\theta - (\mathbf{q}_i\beta_Q + \epsilon_Q))^2$ where $\mathbf{q}_i \sim \mathcal{N}_d(\mathbf{0}, \xi_Q I)$.

First we will formalize the CDF of P and Q . Note \mathbf{p}_i and \mathbf{q}_i represent the normally sampled gaussian data.

$$P_i = (\mathbf{p}_i\theta_{(t)} - (\mathbf{p}_i\beta_P + \epsilon_P))^2 \quad (38)$$

$$= (\mathbf{p}_i(\theta_{(t)} - \beta_P) - \epsilon_P)^2 \quad (39)$$

$$= (\mathbf{p}_i(\theta_{(t)} - \beta_P))^2 - 2\epsilon_P (\mathbf{p}_i(\theta_{(t)} - \beta_P)) + \epsilon_P^2 \quad (40)$$

As $\mathbb{E}[\epsilon_P] = 0$, we will only consider the case $\epsilon_P = 0$. This simplifies P_i and Q_i :

$$P_i = (\mathbf{P}_i(\theta_{(t)} - \beta_P))^2 \quad (41)$$

$$Q_i = (\mathbf{Q}_i(\theta_{(t)} - \beta_Q))^2 \quad (42)$$

Let us note all the entries of \mathbf{P}_i are sampled from $\mathcal{N}(0, \xi_P)$ and all entries of \mathbf{Q}_i are sampled from $\mathcal{N}(0, \xi_Q)$. It thus follows:

$$\begin{aligned} (\mathbf{P}_i(\theta_{(t)} - \beta_P))^2 &= \left(\sum_{j=1}^d \mathcal{N}(0, \xi_P) (\theta_j^{(t)} - \beta_{Pj}) \right)^2 \\ &= \left(\sum_{j=1}^d \mathcal{N}\left(0, \xi_P (\theta_j^{(t)} - \beta_{Pj})^2\right) \right)^2 \\ &= \left(\mathcal{N}\left(0, \sum_{j=1}^d \xi_P (\theta_j^{(t)} - \beta_{Pj})^2\right) \right)^2 \end{aligned}$$

It thus similarly follows for Q_i :

$$(\mathbf{Q}_i(\theta_{(t)} - \beta_Q))^2 = \left(\mathcal{N}\left(0, \sum_{j=1}^d \xi_Q (\theta_j^{(t)} - \beta_{Qj})^2\right) \right)^2$$

Now we can define the cumulative distribution functions.

$$F_P(z) = \frac{1}{\sqrt{2\pi}} \int_{-\sqrt{z}}^{\sqrt{z}} \exp\left(-\frac{u^2}{2 \sum_{j=1}^d \xi_P (\theta_j^{(t)} - \beta_{Pj})^2}\right) du \quad (43)$$

$$= \Phi\left(\frac{\sqrt{z}}{\sqrt{\sum_{j=1}^d \xi_P (\theta_j^{(t)} - \beta_{Pj})^2}}\right) - \Phi\left(\frac{-\sqrt{z}}{\sqrt{\sum_{j=1}^d \xi_Q (\theta_j^{(t)} - \beta_{Pj})^2}}\right) \quad (44)$$

Let us define $\phi = \sum_{j=1}^d \xi_P \left(\theta_j^{(t)} - \beta_{Pj} \right)^2$ and $\psi = \sum_{j=1}^d \xi_Q \left(\theta_j^{(t)} - \beta_{Qj} \right)^2$. Therefore it follows:

$$F_P(z) = \Phi \left(\frac{\sqrt{z}}{\sqrt{\phi}} \right) - \Phi \left(\frac{-\sqrt{z}}{\sqrt{\phi}} \right) \quad (45)$$

Similarly for $F_Q(z)$ it follows:

$$F_Q(z) = \Phi \left(\frac{\sqrt{z}}{\sqrt{\psi}} \right) - \Phi \left(\frac{-\sqrt{z}}{\sqrt{\psi}} \right) \quad (46)$$

Here we can note if $\phi < \psi$, then for all $z > 0$, it follows $F_P(z) > F_Q(z)$.

We now want to find the max integer, η such that with probability greater than 0.9, $Q_\eta < P_{np-\eta}$ and for all $m \in \mathbb{N}$ such that $m > \eta$, $Q_m > P_{np-\eta}$. This is equivalent to the probability there will be η elements from Q within the Subquantile matrix.

$$\begin{aligned} \mathbb{P} \left[Q_\eta < P_{np-\eta} \bigcap_{i=1}^{n\epsilon-\eta} Q_{\eta+i} > P_{np-\eta} \right] &= \mathbb{P} [Q_\eta < P_{np-\eta}] \prod_{i=1}^{n\epsilon-\eta} \mathbb{P} [Q_{n+i} > P_{np-\eta}] \\ &= \mathbb{P} [Q_n < P_{np-\eta}] \prod_{i=1}^{n\epsilon-\eta} (1 - \mathbb{P} [Q_{n+i} < P_{np-\eta}]) \\ &= \mathbb{P} [Q_n - P_{np-\eta} < 0] \prod_{i=1}^{n\epsilon-\eta} (1 - \mathbb{P} [Q_{n+i} - P_{np-\eta} < 0]) \end{aligned}$$

It follows η will be lower depending on the difference of ϕ and ψ . In a further work we will give a tight upper bound on the probability.

D THEORY FOR SUBQUANTILE MINIMIZATION

D.1 DERIVATION OF LEMMA 1.2

Since $g(t, \theta)$ is a concave function. Maximizing $g(t, \theta)$ is equivalent to minimizing $-g(t, \theta)$. We will find fermat's optimality condition for the function $-g(t, \theta)$, which is convex. Let $\hat{\nu} = \text{sorted}((\theta^T X - y)^2)$ and note $0 < p < 1$

$$\partial(-g(t, \theta)) = \partial\left(-t + \frac{1}{np} \sum_{i=1}^n (t - \hat{\nu}_i)^+\right) \quad (47)$$

$$= \partial(-t) + \partial\left(\frac{1}{np} \sum_{i=1}^n (t - \hat{\nu}_i)^+\right) \quad (48)$$

$$= -1 + \frac{1}{np} \sum_{i=1}^n \partial(t - \hat{\nu}_i)^+ \quad (49)$$

$$= -1 + \frac{1}{np} \sum_{i=1}^n \begin{cases} 1, & \text{if } t > \hat{\nu}_i \\ 0, & \text{if } t < \hat{\nu}_i \\ [0, 1], & \text{if } t = \hat{\nu}_i \end{cases} \quad (50)$$

$$= 0 \text{ when } t = \hat{\nu}_{np} \quad (51)$$

This is the p -quantile of ν . Assuming no two points are equal in the dataset, this means the minimizing value for t has a range of values, $\hat{\nu}_{np} \leq t < \hat{\nu}_{np+1}$. This means $g(t, \theta)$ is not strongly convex with respect to t .

D.2 DERIVATION OF LEMMA 1.3

Note that $t_k = \nu_{np}$ which is equivalent to $(\theta_k^T \mathbf{x}_{np} - y_{np})^2$

$$\begin{aligned} \nabla_{\theta_k} g(t_{k+1}, \theta_k) &= \nabla_{\theta_k} \left(\nu_{np} - \frac{1}{np} \sum_{i=1}^n (\nu_{np} - (\theta_k^T \mathbf{x}_i - y_i)^2)^+ \right) \\ &= \nabla_{\theta_k} \left((\theta_k^T \mathbf{x}_{np} - y_{np})^2 - \frac{1}{np} \sum_{i=1}^n ((\theta_k^T \mathbf{x}_{np} - y_{np})^2 - (\theta_k^T \mathbf{x}_i - y_i)^2)^+ \right) \\ &= \nabla_{\theta_k} (\theta_k^T \mathbf{x}_{np} - y_{np})^2 - \frac{1}{np} \sum_{i=1}^n \nabla_{\theta_k} ((\theta_k^T \mathbf{x}_{np} - y_{np})^2 - (\theta_k^T \mathbf{x}_i - y_i)^2)^+ \\ &= 2\mathbf{x}_{np}(\theta_k^T \mathbf{x}_{np} - y_{np}) - \frac{1}{np} \sum_{i=1}^n 2\mathbf{x}_{np}(\theta_k^T \mathbf{x}_{np} - y_{np}) \\ &\quad - 2\mathbf{x}_i(\theta_k^T \mathbf{x}_i - y_i) \begin{cases} 1, & \text{if } t > v_i \\ 0, & \text{if } t < v_i \\ [0, 1], & \text{if } t = v_i \end{cases} \\ &= 2\mathbf{x}_{np}(\theta_k^T \mathbf{x}_{np} - y_{np}) - \frac{1}{np} \sum_{i=1}^{np} 2\mathbf{x}_{np}(\theta_k^T \mathbf{x}_{np} - y_{np}) - 2\mathbf{x}_i(\theta_k^T \mathbf{x}_i - y_i) \\ &= 2\mathbf{x}_{np}(\theta_k^T \mathbf{x}_{np} - y_{np}) - 2\mathbf{x}_{np}(\theta_k^T \mathbf{x}_{np} - y_{np}) + \frac{1}{np} \sum_{i=1}^{np} 2\mathbf{x}_i(\theta_k^T \mathbf{x}_i - y_i) \\ &= \frac{1}{np} \sum_{i=1}^{np} 2\mathbf{x}_i(\theta_k^T \mathbf{x}_i - y_i) \end{aligned}$$

This is the derivative of the np samples with lowest error with respect to θ .

D.3 DERIVATION OF LEMMA 1.5

The objective function $g(\theta, t)$ is L -smooth w.r.t θ iff

$$\|\nabla_{\theta} g(\theta', t) - \nabla_{\theta} g(\theta, t)\| \leq L \|\theta' - \theta\| \quad (52)$$

$$\left\| \nabla_{\theta} g(\theta', t) - \nabla_{\theta} g(\theta, t) \right\| = \left\| \frac{1}{np} \sum_{i=1}^{np} 2\mathbf{x}_i(\theta_k'^{\top} \mathbf{x}_i - y_i) - \frac{1}{np} \sum_{i=1}^{np} 2\mathbf{x}_i(\theta_k^T \mathbf{x}_i - y_i) \right\| \quad (53)$$

$$= \left\| \frac{1}{np} \sum_{i=1}^{np} 2\mathbf{x}_i(\theta_k'^{\top} \mathbf{x}_i - \theta_k^T \mathbf{x}_i) \right\| \quad (54)$$

$$= \left\| \frac{1}{np} \sum_{i=1}^{np} 2\mathbf{x}_i \mathbf{x}_i^T (\theta_k'^{\top} - \theta_k^T) \right\| \quad (55)$$

$$\stackrel{\text{Cauchy-Schwarz}}{\leq} \left\| \frac{2}{np} \sum_{i=1}^{np} \mathbf{x}_i \mathbf{x}_i^T \right\| \left\| \theta_k'^{\top} - \theta_k^T \right\| \quad (56)$$

$$= L \left\| \theta_k'^{\top} - \theta_k^T \right\| \quad (57)$$

where $L = \left\| \frac{2}{np} X^T X \right\|$

This concludes the derivation.

D.4 PROOF OF LEMMA 1.6

Proof. We will investigate the two cases $t_{k+1} \leq t$ and $t_{k+1} > t_k$.

Case (i) $t_{k+1} \leq t_k$

Let us first expand out $g(t_k, \theta_k)$ with the knowledge that $t_k \geq \nu_k$

$$g(t_k, \theta_k) = t_k - \frac{1}{np} \sum_{i=1}^n (t_k - \nu_i)^+ \quad (58)$$

$$= t_k - \frac{1}{np} (np)t_k + \frac{1}{np} \sum_{i=1}^{np} \nu_i + \frac{1}{np} \sum_{i=np}^n (t_k - \nu_i)^+ \quad (59)$$

$$= \frac{1}{np} \sum_{i=1}^{np} \nu_i + \frac{1}{np} \sum_{i=np}^n (t_k - \nu_i)^+ \quad (60)$$

$$g(t_{k+1}, \theta_k) - g(t_k, \theta_k) = \frac{1}{np} \sum_{i=1}^{np} \nu_i - \left(\frac{1}{np} \sum_{i=1}^{np} \nu_i + \frac{1}{np} \sum_{i=np}^n (t_k - \nu_i)^+ \right) \quad (61)$$

$$= -\frac{1}{np} \sum_{i=np}^n (t_k - \nu_i)^+ \quad (62)$$

Case (ii) $t_{k+1} > t_k$

Since we know t_k is less than ν_{np} , WLOG we will say t_k is greater than the lowest $n(p - \delta)$ elements, where $\delta \in (0, p)$.

$$g(t_k, \theta_k) = t_k - \frac{1}{np} \sum_{i=1}^n (t_k - \nu_i)^+ \quad (63)$$

$$= t_k - \frac{1}{np} \sum_{i=1}^{n(p-\delta)} (t_k - \nu_i)^+ \quad (64)$$

$$= t_k - \frac{1}{np} (n(p - \delta))t_k + \frac{1}{np} \sum_{i=1}^{n(p-\delta)} \nu_i \quad (65)$$

$$g(t_k, \theta_{k+1}) - g(t_k, \theta_k) = \frac{1}{np} \sum_{i=1}^{np} \nu_i - \left(\delta t_k + \frac{1}{np} \sum_{i=1}^{n(p-\delta)} \nu_i \right) \quad (66)$$

$$= \left(\frac{1}{np} \sum_{i=n(p-\delta)}^n \nu_i \right) - \delta t_k \quad (67)$$

This concludes the proof.

□

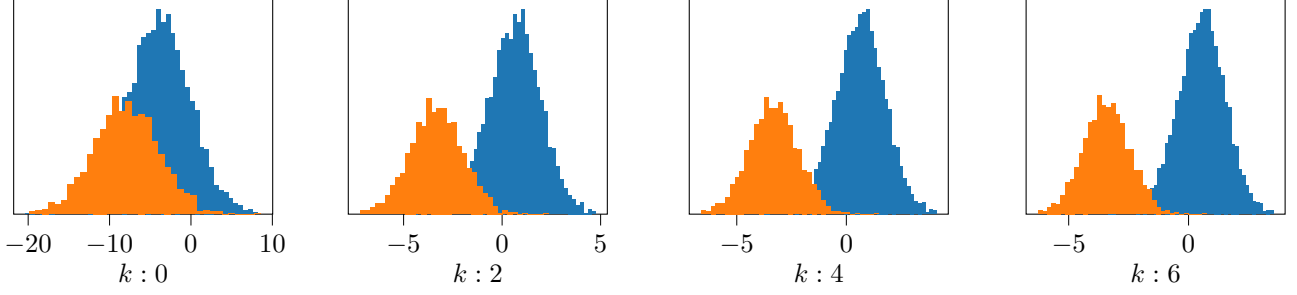


Figure 4: Residuals with respect to \mathbb{P} and \mathbb{Q} , k represents optimization step.

E THEORY FOR ADAPTIVE LINEAR CORRUPTION

In this section, we provide rigorous theory for why Sub-Quantile Minimization works so well in the case of corruption of the form $\beta_Q^T \mu = y_P + \epsilon_Q$.

Assumption 7. *The residuals of θ_k are normally distributed with respect to \mathbb{P} and \mathbb{Q} . In other words, $\theta_k^T p - y_P$ and $\theta_k^T q - y_Q$ are normally distributed.*

Assumption 7 can be visually verified in figure 4. Even after multiple iteration steps the residuals with respect to \mathbb{P} and \mathbb{Q} are still normal. Thus it follows by decreasing $\|\theta - \beta_P\|_1$ more relative to $\|\theta - \beta_Q\|_1$ then the SubQuantile will contain more points from P by expectation.

E.1 PROOF OF THEOREM 3

Proof. To show the change in ε we will first calculate the expected change in θ by the θ -update described in Equation 12. We will also introduce some notation, \mathcal{S} represents all $x \in X$ that are within the lowest np losses, i.e. within the subquantile, $p \in \mathcal{S}$ represent all data vectors from \mathbb{P} that are within the SubQuantile, similarly $q \in \mathbb{Q}$ represent all data vectors from \mathbb{Q} that are within the SubQuantile. Furthermore, $|\mathcal{S}| = np$, there are εn points from \mathbb{Q} in \mathcal{S} and $(1 - \varepsilon)n$ points from \mathbb{P} within \mathcal{S} .

$$\begin{aligned} \mathbb{E}[\theta_{k+1}] &= \theta_k - \mathbb{E}[\alpha \nabla g(\theta_k, t_{k+1})] \\ &= \theta_k - \alpha \mathbb{E} \left[\sum_{x \in \mathcal{S}} x(\theta_k^T x - y) \right] \\ &= \theta_k - \alpha \mathbb{E} \left[\sum_{x \in \mathcal{S}} x x^T \theta_k - x y \right] \\ &= \theta_k - \alpha \mathbb{E} \left[\sum_{p \in \mathcal{S}} p p^T \theta_k - p y_p + \sum_{q \in \mathcal{S}} q q^T \theta_k - q y_q \right] \end{aligned}$$

We will use Assumption 1 to rewrite y_p and y_q

$$\begin{aligned} &= \theta_k - \alpha \mathbb{E} \left[\sum_{p \in \mathcal{S}} p p^T \theta_k - p(\beta_P^T p + \epsilon_P) + \sum_{q \in \mathcal{S}} q q^T \theta_k - q(\beta_Q^T p + \epsilon_Q) \right] \\ &= \theta_k - \alpha \left(\sum_{p \in \mathcal{S}} (\mu \mu^T + \Sigma) \theta_k - (\mu \mu^T + \Sigma) \beta_P - \sum_{q \in \mathcal{S}} (\mu \mu^T + \Sigma) \theta_k + (\mu \mu^T + \Sigma) \beta_Q \right) \end{aligned}$$

$$\mathbb{E}[\theta_{(t+1)}] = \theta_{(t)} - \alpha np C (\theta_{(t)} - (1 - \varepsilon) \beta_P - \varepsilon \beta_Q) \quad (68)$$

Now that we have the expected update for θ in terms of the linear regression coefficients, we now want to utilize Lemma 2.1.

Let $\theta_k = \alpha_1 \beta_P + \alpha_2 \beta_Q + \sum_{i=3}^d \alpha_i r_i$ in the same basis \mathcal{B} defined in Lemma 2.1. Note in the case of data that is normally

distributed about 0 with no covariance amongst the predictor variables, then \mathbf{C} is a multiple of the identity. If we preprocess the data using a standard normal scalar such that all features follow a $\mathcal{N}(0, 1)$ distribution, then we can assume it follows $\mathbf{C} = (\mathbf{00}^T + I) = I$. Let $\gamma := \alpha n p$. Then the following manipulations hold:

$$\begin{aligned}\mathbb{E}[\boldsymbol{\theta}_{(t+1)}] &= \boldsymbol{\theta}_{(t)} - \gamma \left(\boldsymbol{\theta}_{(t)} - (1 - \varepsilon^{(t)})\boldsymbol{\beta}_P - \varepsilon^{(t)}\boldsymbol{\beta}_Q \right) \\ &= \boldsymbol{\theta}_{(t)}(1 - \gamma) + \gamma(1 - \varepsilon^{(t)})\boldsymbol{\beta}_P + (1 - \gamma)\varepsilon^{(t)}\boldsymbol{\beta}_Q \\ &= (1 - \gamma) \left(\alpha_1\boldsymbol{\beta}_P + \alpha_2^{(t)}\boldsymbol{\beta}_Q + \sum_{i=3}^d \alpha_i^{(t)}\mathbf{r}_i \right) + \gamma(1 - \varepsilon^{(t)})\boldsymbol{\beta}_P + \gamma\varepsilon^{(t)}\boldsymbol{\beta}_Q \\ &= \left(\alpha_1^{(t)}(1 - \gamma) + \gamma(1 - \varepsilon^{(t)}) \right) \boldsymbol{\beta}_P + \left(\alpha_2^{(t)}(1 - \gamma) + \gamma\varepsilon^{(t)} \right) \boldsymbol{\beta}_Q + (1 - \gamma) \sum_{i=3}^d \alpha_i^{(t)}\mathbf{r}_i\end{aligned}$$

We will now calculate the difference.

$$\begin{aligned}\mathbb{E}[\boldsymbol{\theta}_{(t+1)} - \boldsymbol{\theta}_{(t)}] &= \left(-\gamma\alpha_1^{(t)} + \gamma(1 - \varepsilon^{(t)}) \right) \boldsymbol{\beta}_P + \left(-\gamma\alpha_2^{(t)} + \gamma\varepsilon^{(t)} \right) \boldsymbol{\beta}_Q + (1 - \gamma) \sum_{i=3}^d \alpha_i^{(t)}\mathbf{r}_i \\ &= \gamma \left((1 - \varepsilon^{(t)}) - \alpha_1^{(t)} \right) \boldsymbol{\beta}_P + \gamma \left(\varepsilon^{(t)} - \alpha_2^{(t)} \right) \boldsymbol{\beta}_Q + (1 - \gamma) \sum_{i=3}^d \alpha_i^{(t)}\mathbf{r}_i\end{aligned}$$

Thus the conditions for ε to decrease by expectation are:

$$\left\| \left((1 - \varepsilon^{(t)}) - \alpha_1^{(t)} \right) \boldsymbol{\beta}_P \right\| > \left\| \left(\varepsilon^{(t)} - \alpha_2^{(t)} \right) \boldsymbol{\beta}_Q \right\| \quad (69)$$

This concludes the proof. \square

E.2 PROOF OF THEOREM ??

Proof. To calculate the probability $\varepsilon^{(t+1)} < \varepsilon^{(t)}$, we will first calculate $\mathbb{E}[\varepsilon^{(t+1)}]$. We will start by calculating the expectation of the loss.

$$\begin{aligned}\mathbb{E} \left[\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{p}_i - y_i \right)^2 \right] &= \mathbb{E} \left[\boldsymbol{\theta}_{(t+1)}^T \mathbf{p}_i - y_i \right]^2 + \text{Var} \left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{p}_i - y_i \right) \\ &= \mathbb{E} \left[\boldsymbol{\theta}_{(t+1)}^T \mathbf{p}_i - \boldsymbol{\beta}_P^T \mathbf{p}_i - \epsilon_P \right]^2 + \text{Var} \left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{p}_i - \boldsymbol{\beta}_P^T \mathbf{p}_i - \epsilon_P \right) \\ &= \left(\left(\boldsymbol{\theta}_{(t+1)}^T - \boldsymbol{\beta}_P^T \right) \mathbb{E}[\mathbf{p}_i] \right)^2 + \text{Var} \left(\left(\boldsymbol{\theta}_{(t+1)}^T - \boldsymbol{\beta}_P^T \right) \mathbf{p}_i \right) + \text{Var}(\epsilon_P) \\ &= \left(\left(\boldsymbol{\theta}_{(t+1)}^T - \boldsymbol{\beta}_P^T \right) \boldsymbol{\mu} \right)^2 + \left(\boldsymbol{\theta}_{(t+1)}^T - \boldsymbol{\beta}_P^T \right) \Sigma \left(\boldsymbol{\theta}_{(t+1)} - \boldsymbol{\beta}_P \right) + \text{Var}(\epsilon_P)\end{aligned}$$

It thus follows similarly:

$$\mathbb{E} \left[\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{q}_i - y_i \right)^2 \right] = \left(\left(\boldsymbol{\theta}_{(t+1)}^T - \boldsymbol{\beta}_Q^T \right) \boldsymbol{\mu} \right)^2 + \left(\boldsymbol{\theta}_{(t+1)}^T - \boldsymbol{\beta}_Q^T \right) \Sigma \left(\boldsymbol{\theta}_{(t+1)} - \boldsymbol{\beta}_Q \right) + \text{Var}(\epsilon_Q)$$

To simplify our notation, let $\zeta_P := \mathbb{E} \left[\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{p}_i - y_i \right)^2 \right]$ and $\zeta_Q := \mathbb{E} \left[\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{q}_i - y_i \right)^2 \right]$ and let $\eta_P := \mathbb{E} \left[\boldsymbol{\theta}_{(t+1)}^T \mathbf{p}_i - y_i \right]$

and $\eta_Q := \mathbb{E} \left[\boldsymbol{\theta}_{(t+1)}^T \mathbf{q}_i - y_i \right]$. We are now interested in calculating the variance of the loss.

Let $\sigma_P^2 := \text{Var} \left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{p}_i - y_i \right) = \left(\boldsymbol{\theta}_{(t+1)}^T - \boldsymbol{\beta}_P^T \right) \Sigma \left(\boldsymbol{\theta}_{(t+1)} - \boldsymbol{\beta}_P \right) + \text{Var}(\epsilon_P)$.

$$\begin{aligned}\text{Var} \left(\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{p}_i - y_i \right)^2 \right) &= \mathbb{E} \left[\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{p}_i - y_i \right)^4 \right] - \mathbb{E} \left[\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{p}_i - y_i \right)^2 \right]^2 \\ &= \mathbb{E} \left[\boldsymbol{\theta}_{(t+1)}^T \mathbf{p}_i - y_i \right]^4 + 6 \mathbb{E} \left[\boldsymbol{\theta}_{(t+1)}^T \mathbf{p}_i - y_i \right]^2 \sigma_P^2 + 3 \sigma_P^4 - \zeta_P^2 \\ &= \eta_P^4 + 6 \eta_P^2 \sigma_P^2 + 3 \sigma_P^4 - \zeta_P^2\end{aligned}$$

It similarly follows:

$$\text{Var} \left(\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{q}_i - y_i \right)^2 \right) = \eta_Q^4 + 6 \eta_Q^2 \sigma_Q^2 + 3 \sigma_Q^4 - \zeta_Q^2$$

Recall that the expected value of $\varepsilon^{(t+1)}$ is equal to the expected number of elements from Q within the p -Quantile, \mathcal{Q}_p . So we will now calculate \mathcal{Q}_p . To this we will utilize the loss of the distribution.

$$\begin{aligned}\mathbb{E} \left[\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{x}_i - y_i \right)^2 \right] &= (1 - \varepsilon^{(t)}) \mathbb{E} \left[\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{p}_i - y_i \right)^2 \right] + \varepsilon^{(t)} \mathbb{E} \left[\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{q}_i - y_i \right)^2 \right] \\ &= (1 - \varepsilon^{(t)}) \zeta_P + \varepsilon^{(t)} \zeta_Q \\ \text{Var} \left(\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{x}_i - y_i \right)^2 \right) &= (1 - \varepsilon^{(t)})^2 \text{Var} \left(\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{p}_i - y_i \right)^2 \right) + \varepsilon^{(t)2} \text{Var} \left(\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{q}_i - y_i \right)^2 \right) \\ &= (1 - \varepsilon^{(t)})^2 (\eta_P^4 + 6\eta_P^2 \sigma_P^2 + 3\sigma_P^4 - \zeta_P^2) + (\varepsilon^{(t)})^2 (\eta_Q^4 + 6\eta_Q^2 \sigma_Q^2 + 3\sigma_Q^4 - \zeta_Q^2)\end{aligned}$$

It thus follows $\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{q}_i - y_i \right)$ follows a χ^2 distribution with 1 degree of freedom and

$$\lambda = \frac{\zeta_Q}{\eta_Q^4 + 6\eta_Q^2 \sigma_Q^2 + 3\sigma_Q^4 - \zeta_Q^2}$$

Let $\mathcal{Q}_p^{(t+1)}$ be given as the np greatest loss of the data, i.e. $\mathcal{Q}_p^{(t+1)} = \nu_{np}^{(t+1)}$. It then follows:

$$\mathbb{P} \left[\left(\boldsymbol{\theta}_{(t+1)}^T \mathbf{q}_i - y_i \right)^2 \leq \mathcal{Q}_p \right] = 1 - Q_1 \left(\sqrt{\lambda}, \sqrt{\mathcal{Q}_p} \right)$$

It thus follows:

$$\mathbb{E} \left[\varepsilon^{(t+1)} \right] = 1 - Q_1 \left(\sqrt{\lambda}, \sqrt{\mathcal{Q}_p} \right)$$

Now we can calculate the probability of improvement per iteration by invoking Markov's Inequality:

$$\begin{aligned}\mathbb{P} \left[\varepsilon^{(t+1)} \geq \varepsilon^{(t)} \right] &\leq \frac{\mathbb{E} \left[\varepsilon^{(t+1)} \right]}{\varepsilon^{(t)}} \\ &= \frac{1 - Q_1 \left(\sqrt{\lambda}, \sqrt{\mathcal{Q}_p} \right)}{\varepsilon^{(t)}}\end{aligned}$$

It thus follows:

$$\mathbb{P} \left[\varepsilon^{(t+1)} \leq \varepsilon^{(t)} \right] \geq 1 - \frac{1 - Q_1 \left(\sqrt{\lambda}, \sqrt{\mathcal{Q}_p} \right)}{\varepsilon^{(t)}}$$

□

This concludes the proof.

F PROOFS FOR CONVERGENCE

F.1 PROOF OF THEOREM 2

Proof. We will first start by introducing new notation. Let S represent a matrix with np data points from X , in other words it is a possible SubQuantile Matrix. Let Π represent the set of all such possible matrices S of X . Note $|\Pi| = \binom{n}{np}$. We can now redefine the min-max optimization problem of g to a min-min optimization problem. Let us define the function $f(\theta, S) = \|\theta^T S - y_S\|_2^2$

$$\theta^*, S^* = \arg \min_{\theta \in \mathbb{R}^d} \arg \min_{S \in \Pi} \|\theta^T S - y_S\|_2^2 \quad (70)$$

Note we have a $\mathcal{O}(n)$ oracle for the $\arg \min_{S \in \Pi} f(\theta_T, S_T)$.

Lemma 3.1. *The resultant $\tilde{S} = \arg \min_{S \in \Pi} f(\theta, S)$ is a unique minimizer iff all points in X are different.*

We will now show f is a monotonically decreasing function.

First let us define $\phi(\cdot) = \min_{S \in \Pi} f(\cdot, S)$. Let us also note f is ℓ smooth with respect to θ . This is following notation from Jin et al. (2019). It thus follows:

$$\begin{aligned} f(\theta_{k+1}, S_k) &\leq f(\theta_k, S_k) + \langle \nabla_{\theta} f(\theta_k, S_k), \theta_{k+1} - \theta_k \rangle + \frac{\ell}{2} \|\theta_{k+1} - \theta_k\|_2^2 \\ &= \phi(\theta_k) + \langle \nabla_{\theta} f(\theta_k, S_k), -\frac{1}{\ell} \nabla_{\theta} f(\theta_k, S_k) \rangle + \frac{\ell}{2} \left\| \frac{1}{\ell} \nabla_{\theta} f(\theta_k, S_k) \right\|_2^2 \\ &= \phi(\theta_k) - \frac{1}{\ell} (\nabla_{\theta} f(\theta_k, S_k))^2 + \frac{1}{2\ell} (\nabla_{\theta} f(\theta_k, S_k))^2 \\ &= \phi(\theta_k) - \frac{1}{2\ell} (\nabla_{\theta} f(\theta_k, S_k))^2 \end{aligned}$$

Thus we have proved the inner optimization problem is monotonically decreasing. Since the outer minimization is strictly less than or equal to the result from the inner optimization, it follows after each two step optimization:

$$\phi(\theta_{k+1}) \leq \phi(\theta_k)$$

Since f is lower bounded by 0. We can invoke Monotonicity Convergence Theorem, since f is a monotonically decreasing function and is lower bounded, it therefore converges to either a local or global minimum. \square

F.2 EXPECTATION OF IMPROVEMENT

We are also interested in expected improvement after the Subquantile matrix update.

$$\begin{aligned} \mathbb{E} \left[f(\theta_{(k)}, S^{(k+1)}) - f(\theta_{(k)}, S^{(k)}) \right] &= \mathbb{E} [\varepsilon^{(t)} - \varepsilon^{(t-1)}] \mathbb{E} \left[\sum_{i=1}^{n(\varepsilon^{(t)} - \varepsilon^{(t-1)})} \left\| \theta_{(k)}^T p_i - y_i \right\|_2^2 - \left\| \theta_{(k)}^T q_i - y_i \right\|_2^2 \right] \\ &= \mathbb{E} [\varepsilon^{(t)} - \varepsilon^{(t-1)}] \mathbb{E} \left[\sum_{i=1}^{n(\varepsilon^{(t)} - \varepsilon^{(t-1)})} \left\| \theta_{(k)}^T p_i - y_i \right\|_2^2 - \left\| \theta_{(k)}^T q_i - y_i \right\|_2^2 \right] \end{aligned}$$

We can now use reverse triangle inequality.

$$\begin{aligned} &\leq \mathbb{E} [\varepsilon^{(t)} - \varepsilon^{(t-1)}] \mathbb{E} \left[\sum_{i=1}^{n(\varepsilon^{(t)} - \varepsilon^{(t-1)})} \left\| \theta_{(k)}^T p_i - \beta_P p_i - \epsilon_P - \theta_{(k)}^T q_i + \beta_Q q_i + \epsilon_Q \right\|_2^2 \right] \\ &= \mathbb{E} [\varepsilon^{(t)} - \varepsilon^{(t-1)}] \mathbb{E} \left[\sum_{i=1}^{n(\varepsilon^{(t)} - \varepsilon^{(t-1)})} \left\| (\theta_{(k)}^T - \beta_P) p_i + (\beta_Q - \theta_{(k)}^T) q_i - \epsilon_P + \epsilon_Q \right\|_2^2 \right] \end{aligned}$$

Now we can use theorem 1. To simplify notation, let $\Xi \triangleq \mathbb{E} [\varepsilon^{(t)} - \varepsilon^{(t-1)}]$

$$\begin{aligned}
&= \Xi \mathbb{E} \left[\sum_{i=1}^{n\Xi} \left\| \left((1 - \varepsilon^{(t)})\beta_P + \varepsilon^{(t)}\beta_Q - \beta_P \right) \mathbf{p}_i + \left(\beta_Q - (1 - \varepsilon^{(t)})\beta_P - \varepsilon^{(t)}\beta_Q \right) \mathbf{q}_i - \epsilon_P + \epsilon_Q \right\|_2^2 \right] \\
&= \Xi \mathbb{E} \left[\sum_{i=1}^{n\Xi} \left\| \left(-\varepsilon^{(t)}\beta_P + \varepsilon^{(t)}\beta_Q \right) \mathbf{p}_i + \left((1 - \varepsilon^{(t)})\beta_Q - (1 - \varepsilon^{(t)})\beta_P \right) \mathbf{q}_i - \epsilon_P + \epsilon_Q \right\|_2^2 \right] \\
&= \Xi \mathbb{E} \left[\sum_{i=1}^{n\Xi} \left\| \varepsilon^{(t)} (\beta_Q - \beta_P) \mathbf{p}_i + (1 - \varepsilon^{(t)}) (\beta_Q - \beta_P) \mathbf{q}_i - \epsilon_P + \epsilon_Q \right\|_2^2 \right]
\end{aligned}$$

Here we can use the fact that $\mathbb{E} [X^2] = \mathbb{E} [X]^2 + \text{Var}(X)$ for a random variable X

Let us define the function $h(\boldsymbol{\theta}_k) = f(\boldsymbol{\theta}_k, S_{k-1}) - f(\boldsymbol{\theta}_k, S_k)$, a strictly non-negative function. From our results above it follows:

$$\begin{aligned}
\phi(\boldsymbol{\theta}_{k+1}) &\leq \phi(\boldsymbol{\theta}_k) - \frac{1}{2\ell} (\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_k, S_k))^2 \\
&= f(\boldsymbol{\theta}_k, S_{k-1}) - h(\boldsymbol{\theta}_k) + \frac{\ell}{2} \|\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_k\|_2^2 \\
&= f(\boldsymbol{\theta}_k, S_{k-1}) - h(\boldsymbol{\theta}_k) - \frac{\ell}{2} \left\| \boldsymbol{\theta}_{k+1} - \frac{1}{\ell} \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_k, S_k) - \boldsymbol{\theta}_k \right\|_2^2 \\
&= f(\boldsymbol{\theta}_k, S_{k-1}) - h(\boldsymbol{\theta}_k) - \frac{\ell}{2} \left(\|\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_k\|_2^2 + \frac{2}{\ell} \|\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_k\| \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_k, S_k) + \frac{1}{\ell^2} \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_k, S_k)^2 \right)
\end{aligned}$$

G STOCHASTIC SUB-QUANTILE OPTIMIZATION

In the age of big data, stochastic methods are necessary for fast training of models to handle large amounts of data. In this section we provide an algorithm for Stochastic Sub-Quantile Optimization.

Algorithm 3: Stochastic Sub-Quantile Minimization Optimization Algorithm

Input: Training iterations T , Quantile p , Corruption Percentage ϵ , Input Parameters d , Batch Size m

Output: Trained Parameters, θ

```
1:  $\theta_1 \leftarrow \mathcal{N}(0, \sigma)^d$ 
2: for  $k \in 1, 2, \dots, T$  do
3:    $I \subseteq [n]$  of size  $m$ 
4:    $\nu = (X_I \theta_k - y_I)^2$ 
5:    $\hat{\nu} = \text{sorted}(\nu)$ 
6:    $t_{k+1} = \hat{\nu}_{mp}$ 
7:    $L := \sum_{i=1}^{mp} \mathbf{x}_i^T \mathbf{x}_i$ 
8:    $\alpha := \frac{1}{2L}$ 
9:    $\theta_{k+1} = \theta_k - \alpha \nabla_{\theta_k} g(t_{k+1}, \theta_k)$ 
10: end
11: return  $\theta_T$ 
```

H ADDITIONAL EXPERIMENTS

H.1 QUADRATIC REGRESSION

Objectives	Test RMSE (Quadratic Regression)		
	$\epsilon = 0$	$\epsilon = 0.2$	$\epsilon = 0.4$
ERM	0.0099 _(0.0002)	2.078 _(0.146)	4.104 _(0.442)
Huber Huber & Ronchetti (2009)	1.000 _(0.0002)	1.000 _(0.0003)	1.13 _(0.087)
RANSAC Fischler & Bolles (1981)	0.010 _(0.0002)	0.011 _(0.0002)	0.061 _(0.053)
TERM Li et al. (2020)	0.010 _(0.0001)	0.012 _(0.0008)	0.017 _(0.0016)
SEVER Diakonikolas et al. (2019)	0.0166 _(0.007)	0.011 _(0.0004)	0.0267 _(0.036)
SubQuantile($p = 0.6$)	0.0099 _(0.0002)	0.00998 _(0.0002)	0.010 _(0.0001)
Genie ERM	0.0099 _(0.0002)	0.00997 _(0.0002)	0.010 _(0.0001)

Table 3: Quadratic Regression Synthetic Dataset. Empirical Risk over \mathbb{P}

H.2 ABALONE

We now provide results on Abalone Dataset introduced in Dua & Graff (2017). This experiment has both feature and label

Objectives	Test RMSE (Abalone Linear Regression)	
	Clean	Noisy
ERM	2.213 _(0.0528)	4845.335 _(117.5557)
CRR Bhatia et al. (2017)	2.345 _(0.0430)	396.872 _(96.5632)
STIR Mukhoty et al. (2019)	2.240 _(0.0473)	931.845 _(32.0864)
Huber Huber & Ronchetti (2009)	5.535 _(0.0665)	971.362 _(28.8863)
RANSAC Fischler & Bolles (1981)	2.522 _(0.1407)	2.621 _(0.1719)
TERM Li et al. (2020)	10.686 _(0.2616)	10.853 _(0.4245)
SEVER Diakonikolas et al. (2019)	2.238 _(0.0901)	2.287 _(0.0757)
SubQuantile($p = 0.8$)	2.292 _(0.0413)	2.261 _(0.0790)
Genie ERM	2.213 _(0.0528)	2.238 _(0.0901)

Table 4: Abalone Regression Real Dataset. Empirical Risk over \mathbb{P}

noise in the Noisy Data. SubQuantile minimization no longer always converges to the \mathbb{P} SubQuantile.

H.3 CAL-HOUSING

We now provide results on Cal-Housing Dataset introduced in Pace & Barry (1997). This experiment has both feature and label noise in the Noisy Data.

In both the Cal-Housing and Abalone datasets there exists feature and label noise that exist with 5% probability. In this the case, the probability is low, however since the noise is very large, even having a few points from \mathbb{Q} in the final subquantile matrix can largely the bias the predictions away from the optimal parameters for \mathbb{P} . Therefore, we reduce p , the size of the subquantile to reduce the probability of obtaining corrupted samples within the subquantile. However, what we get in a decrease in variance, we do increase the bias error, albeit very slightly.

Objectives	Test RMSE (Cal-Housing Linear Regression)	
	Clean	Noisy
ERM	0.598 _(0.0077)	81.758 _(2.6230)
CRR Bhatia et al. (2017)	0.602 _(0.0081)	75.777 _(2.9403)
STIR Mukhoty et al. (2019)	0.604 _(0.0070)	65.555 _(2.1899)
Huber Huber & Ronchetti (2009)	0.601 _(0.0077)	71.813 _(2.0755)
RANSAC Fischler & Bolles (1981)	0.681 _(0.0389)	0.679 _(0.0253)
TERM Li et al. (2020)	0.737 _(0.0070)	0.741 _(0.0155)
SEVER Diakonikolas et al. (2019)	0.640 _(0.0067)	0.642 _(0.0088)
SubQuantile($p = 0.9$)	0.615 _(0.0076)	0.612 _(0.0096)
Genie ERM	0.598 _(0.0077)	0.603 _(0.0068)

Table 5: Cal-Housing Regression Real Dataset. Empirical Risk over \mathbb{P}

I EXPERIMENTAL DETAILS

I.1 ADAPTIVE LINEAR REGRESSION DATASET

We will describe \mathbb{P} and \mathbb{Q} in the Structured Linear Regression Dataset.

$$\mathbf{x} \sim \mathcal{N}(4, 4)^{200}$$

$$\mathbf{m} \sim \mathcal{N}(4, 4)^{200}$$

$$b \sim \mathcal{N}(4, 4)$$

$$\mathbf{m}' \sim \mathcal{N}(4, 4)^{200}$$

$$b' \sim \mathcal{N}(4, 4)$$

$$n_{\text{train}} = 1\text{e}4$$

$$\mathbb{P} : y|\mathbf{x} \sim \mathcal{N}(\mathbf{m}^T \mathbf{x} + b, 0.1)$$

$$\mathbb{Q} : y|\mathbf{x} \sim \mathcal{N}(\mathbf{m}'^T \mathbf{x} + b', 0.1)$$

Please note \mathbf{m} , b , \mathbf{m}' , b' , are all sampled independently. The noise is added after normalization of the dataset to the standard normal $\mathcal{N}(0, 1)$.

I.2 OBLIVIOUS LINEAR REGRESSION DATASET

We will describe \mathbb{P} and \mathbb{Q} in the Noisy Linear Regression Dataset.

$$\mathbf{x} \sim \mathcal{N}(0, 3)^{500}$$

$$\mathbf{m} \sim \mathcal{N}(4, 4)^{500}$$

$$b \sim \mathcal{N}(4, 4)$$

$$\mathbf{m}' = \mathbf{0}$$

$$b' \sim \mathcal{N}(5, 5)$$

$$n_{\text{train}} = 8\text{e}3$$

$$n_{\text{test}} = 2\text{e}3$$

$$\mathbb{P} : y|\mathbf{x} \sim \mathcal{N}(\mathbf{m}^T \mathbf{x} + b, 0.01)$$

$$\mathbb{Q} : y|\mathbf{x} \sim \mathcal{N}(5, 5)$$

Please note \mathbf{m} , b , \mathbf{m}' , b' , are all sampled independently. The noise is added after normalization of the dataset to the standard normal.

I.3 QUADRATIC REGRESSION DATASET

We will describe \mathbb{P} and \mathbb{Q} in the Quadratic Regression dataset.

$$x \sim \mathcal{N}(0, 1)$$

$$n_{\text{train}} = 1\text{e}4$$

$$\mathbb{P} : y|x \sim \mathcal{N}(x^2 - x + 2, 0.01)$$

$$\mathbb{Q} : y|x \sim \mathcal{N}(-x^2 + x + 4, 0.01)$$

I.4 DRUG DISCOVERY DATASET

This dataset is downloaded from Diakonikolas et al. (2019). We utilize the same noise procedure as in Li et al. (2020). \mathbb{P} is given from an 80/20 train test split from the dataset.

\mathbb{Q} is random noise sampled from $\mathcal{N}(5, 5)$.

The noise represents a noisy worker

I.5 FEATURE NOISE

Take 5% of the training data and multiply features by 100 and responses by 10000.