

ROBUST LINEAR REGRESSION BY SUB-QUANTILE OPTIMIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Robust Linear Regression is the problem of fitting data to a distribution, \mathbb{P} when there exists contaminated samples, \mathbb{Q} . We model this as $\hat{\mathbb{P}} = (1 - \epsilon)\mathbb{P} + \epsilon\mathbb{Q}$. Traditional Least Squares Methods fit the empirical risk model to all training data in $\hat{\mathbb{P}}$. 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.

1 INTRODUCTION

Linear Regression is one of the most widely used statistical estimators throughout Science. Although robustness is only a somewhat recent topic in machine learning, it has been a topic in statistics for many decades. Several popular methods have been very popular due to their simplicity and high effectiveness including quantile regression Koenker & Hallock (2001), Theil-Sen Estimator Sen (1968), and Huber Regression Huber & Ronchetti (2009).

Our goal is to provide a theoretic analysis and convergence conditions for sub-quantile optimization and offer practitioners a method for robust linear regression.

In this paper we will show how Sub-Quantile Optimization can address the shortcomings of ERM in the case of corrupted data or imbalanced data, where there exists a majority class and a minority class.

2 RELATED WORK

Least Trimmed Squares (LTS) Mount et al. (2014).

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

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

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.

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

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

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

Super-Quantile Optimization Rockafellar et al. (2014)

Robust Risk Minimization Osama et al. (2020)

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

$$Q_p(p) = \inf\{x \in \mathbb{R} : p \leq F(x)\} \quad (4)$$

Note $Q_p(0.5)$ represents the median of the random variable.

Definition 2. The **Empirical Distribution Function** is defined as follows

$$\hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{X_i \leq t} \quad (5)$$

Definition 3. Let ℓ be the loss function. **Risk** is defined as follows

$$U = \mathbb{E}[\ell(f(\mathbf{x}; \boldsymbol{\theta}), \mathbf{y})] \quad (6)$$

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

$$\mathbb{L}_p(U) = \frac{1}{p} \int_0^p \mathcal{Q}_q(U) dq = \mathbb{E}[U | U \leq \mathcal{Q}_p(U)] = \max_{t \in \mathbb{R}} \left\{ t - \frac{1}{p} \mathbb{E}[(t - U)^+] \right\} \quad (7)$$

In equation 7, 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

$$\boldsymbol{\theta}_{SM} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \max_{t \in \mathbb{R}} \left\{ t - \frac{1}{p} \mathbb{E}(t - \ell(f(\mathbf{x}; \boldsymbol{\theta}), \mathbf{y}))^+ \right\} \quad (8)$$

For the linear regression case, this equation becomes

$$\boldsymbol{\theta}_{SM} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \max_{t \in \mathbb{R}} \left\{ t - \frac{1}{np} \sum_{i=1}^n (t - (\boldsymbol{\theta}^\top \mathbf{x}_i - y_i)^2)^+ \right\} \quad (9)$$

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

$$t_{k+1} = \arg \max_t g(t, \boldsymbol{\theta}_k) \quad (10)$$

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \alpha \nabla_{\boldsymbol{\theta}_k} g(t, \boldsymbol{\theta}_k) \quad (11)$$

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

3.1 MOTIVATION

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

$$\mathbf{y}_P = \mathbf{P}\boldsymbol{\beta}_P + \boldsymbol{\epsilon}_P \quad (12)$$

and similarly

$$\mathbf{y}_Q = \mathbf{Q}\boldsymbol{\beta}_Q + \boldsymbol{\epsilon}_Q \quad (13)$$

where $\boldsymbol{\beta}_P$ and $\boldsymbol{\beta}_Q$ the oracle regressors for \mathbb{P} and \mathbb{Q} .

Since we are interested in learning the optimal model for distributions, our goal is to learn the parameters $\boldsymbol{\beta}_P$ from the distribution $\hat{\mathbb{P}}$. 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 \mathbf{P} represent the data from distribution \mathbb{P} and let \mathbf{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 \mathbf{P} and \mathbf{Q} are sampled from the same multivariate normal distribution.

$$\mathbf{P}_i, \mathbf{Q}_j \sim \mathcal{N}_p(\mathbf{0}, \boldsymbol{\Sigma}) \quad (14)$$

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

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

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

$$\mathbb{E}[\mathbf{X}^\dagger \mathbf{y}] = \beta_P + \epsilon(\beta_Q - \beta_P) \quad (15)$$

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.

Theorem 2. *In the case of \mathbb{P} and \mathbb{Q} being from different Normal Distributions. The expected optimal parameters of the corrupted model $\hat{\mathbb{P}}$*

$$\mathbb{E}[\mathbf{X}^\dagger \mathbf{y}] = \beta_P - n(1 - \epsilon)\Sigma_P^{-1}\beta_P \quad (16)$$

The proof can be found in Appendix B.2. We will show our results hold in Numerical Experiments. As seen in the results in table 1, the theory we provide is supported by Numerical Experimentation.

Dataset	$\epsilon = 0.2$		$\epsilon = 0.4$	
	$\mathbb{E}[\mathbf{X}^\dagger \mathbf{y}]$	Experimental	$\mathbb{E}[\mathbf{X}^\dagger \mathbf{y}]$	Experimental
Quadratic Regression	(0.6, -0.6, 2.4)	0.777 _(0.007)	(0.2, -0.2, 2.8)	7.749 _(0.009)
Drug Discovery	0.895 _(0.009)	0.775 _(0.006)	8.944 _(0.007)	7.742 _(0.006)

Table 1: Verification of Theorem 1 over Quadratic Regression Synthetic Dataset

In equation 15, 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 4.1. *$g(t_{k+1}, \theta_k)$ 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 4.2. *The maximizing value of t in $g(t, \theta)$ in t -update step of optimization as described by Equation 10 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 $\hat{\nu} = \text{sorted}((\theta^\top \mathbf{X} - \mathbf{y})^2)$ 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 C.1. \square

Lemma 4.3. *Let $t = \hat{\nu}_{np}$. The θ -update step described in Equation 9 is equivalent to minimizing the least squares loss of the np elements with the lowest squared loss.*

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

We provide a proof in Appendix C.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 4.4. *$g(t_{k+1}, \theta_k)$ is convex with respect to θ_k .*

Proof. We see by lemma 4.2 and interpretation 1, we are optimizing by the np points with the lowest squared error. Mathematically,

$$g(t_{k+1}, \theta_k) = t_{k+1} - \frac{1}{np} \sum_{i=1}^n (t_{k+1} - (\theta^\top \mathbf{x}_i - y_i)^2)^+ \quad (20)$$

$$= t_{k+1} - \frac{1}{np} \sum_{i=1}^{np} (t_{k+1} - (\theta^\top \mathbf{x}_i - y_i)^2)^+ \quad (21)$$

$$= t - t + \frac{1}{np} \sum_{i=1}^{np} (\theta^\top \mathbf{x}_i - y_i)^2 \quad (22)$$

$$= \frac{1}{np} \sum_{i=1}^{np} (\theta^\top \mathbf{x}_i - y_i)^2 \quad (23)$$

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 4.5. *$g(t, \theta)$ is L -smooth with respect to θ with $L = \left\| \frac{2}{np} \sum_{i=1}^{np} \|\mathbf{x}_i\|^2 \right\|$*

4.2 OPTIMIZATION

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

Definition 4. *(t^*, θ^*) is a **Nash Equilibrium** of g if for any $(t, \theta) \in \mathbb{R} \times \mathbb{R}^d$*

$$g(t^*, \theta) \leq g(t^*, \theta^*) \leq g(t, \theta^*) \quad (24)$$

Definition 5. *(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:*

$$g(t^*, \theta) \leq g(t^*, \theta^*) \leq g(t, \theta^*) \quad (25)$$

Proposition 1. *As g is first-order differentiable, any local Nash Equilibrium satisfies $\nabla_{\theta} g(t, \theta) = 0$ 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 4.2, we have shown $\nabla_t g(t, \theta) =$

0 when $\nu_{np} \leq t < \nu_{np+1}$. Furthermore, by lemma 4.3, we have shown $\nabla_{\theta}(g, \theta) = 0$ when the least squares error is minimized for the np points with lowest squared error. This means that for a subset of np points from \mathbf{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 \mathbb{Q} , as they will introduce error to our predictions on points from \mathbb{P} .

Lemma 4.6. *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.1 \square

4.3 REDUCING ϵ

Let us define two functions for the empirical loss on \mathbb{P} and \mathbb{Q}

$$\phi(\theta) = \frac{1}{np} \sum_{i=1}^m (\theta^\top \mathbf{p}_i - y_i)^2 \quad (26)$$

$$\psi(\theta) = \frac{1}{np} \sum_{i=1}^l (\theta^\top \mathbf{q}_i - y_i)^2 \quad (27)$$

These two functions hold nice properties.

$$\nabla_{\theta} \phi(\theta) = \frac{1}{np} \sum_{i=1}^m 2\mathbf{p}_i (\theta^\top \mathbf{p}_i - y_i) \quad (28)$$

$$\nabla_{\theta} \psi(\theta) = \frac{1}{np} \sum_{i=1}^l 2\mathbf{q}_i (\theta^\top \mathbf{q}_i - y_i) \quad (29)$$

Here we note that the summation of these derivatives is equal to the theta update

$$\nabla_{\theta_k} g(t_{k+1}, \theta_k) = \nabla_{\theta_k} \phi(\theta) + \nabla_{\theta_k} \psi(\theta) \quad (30)$$

Assumption 3. *Training on a subset of the from \mathbb{P} or from \mathbb{Q} generalizes well to the total data.*

From Assumption 3, we can make the assumption that if we take an optimization step with respect to the data in \mathbb{P} or \mathbb{Q} within the lowest np squared errors, it will generalize well to the data from the $n(1-p)$ highest squared errors.

Theorem 3. *In each iteration of the two-step optimization, by expectation the number of elements from \mathbb{P} will increase.*

Corollary 3.1. *If the $\epsilon > 0.5$, i.e., the corruption is the majority class, Sub-Quantile Optimization is still able to converge to the optimal class when p is chosen to be less than 0.5.*

Corollary 3.1 seems like a very unintuitive result, as the whole idea of subquantile minimization is to optimize over the majority class and leave out the tail of the distribution. We provide a proof in

5 EMPIRICAL RESULTS

The first experiment we will run will display the difference of the following two t updates

$$t_{k+1} = \hat{\nu}_{np} \quad (31)$$

$$t_{k+1} = \frac{1}{np} \sum_{i=1}^{np} \hat{\nu}_i \quad (32)$$

In general, if the $\hat{\nu}_1, \hat{\nu}_2, \dots, \hat{\nu}_{np}$ are closely distributed, then $\frac{1}{np} \sum_{i=1}^{np} \hat{\nu}_i \approx \hat{\nu}_{np}$. In Algorithm 1, we display our training method for Sub-quantile Optimization with the t update as described in equation 31. We also compare against the t -update as described in equation 32.

Algorithm 1: Sub-Quantile Minimization Optimization Algorithm**Input:** Training iterations T , Quantile p , Corruption Percentage ϵ , Input Parameters m **Output:** Trained Parameters, θ **Data:** Inliers: $y|x \sim \mathcal{N}(x^2 - x + 2, 0.01)$, Outliers: $y|x \sim \mathcal{N}(-x^2 + x + 4, 0.01)$

```

1:  $\theta_1 \leftarrow \mathcal{N}(0, \sigma)^d$ 
2: for  $k \in 1, 2, \dots, m$  do
3:    $\nu = (X\theta_k - y)^2$ 
4:    $\hat{\nu} = \text{sorted}(\nu)$ 
5:    $t_{k+1} = \hat{\nu}_{np}$ 
6:    $t_{k+1} = \frac{1}{np} \sum_{i=1}^{np} \nu_i$ 
7:    $L := \sum_{i=1}^{np} x_i^\top x_i$ 
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  $\theta_T$ 

```

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.

Algorithm 2: 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 **Data:** Inliers: $y|x \sim \mathcal{N}(x^2 - x + 2, 0.01)$, Outliers: $y|x \sim \mathcal{N}(-x^2 + x + 4, 0.01)$

```

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:    $t_{k+1} = \frac{1}{mp} \sum_{i=1}^{mp} \nu_i$ 
8:    $L := \sum_{i=1}^{mp} x_i^\top x_i$ 
9:    $\alpha := \frac{1}{2L}$ 
10:   $\theta_{k+1} = \theta_k - \alpha \nabla_{\theta_k} g(t_{k+1}, \theta_k)$ 
11: end
12: return  $\theta_T$ 

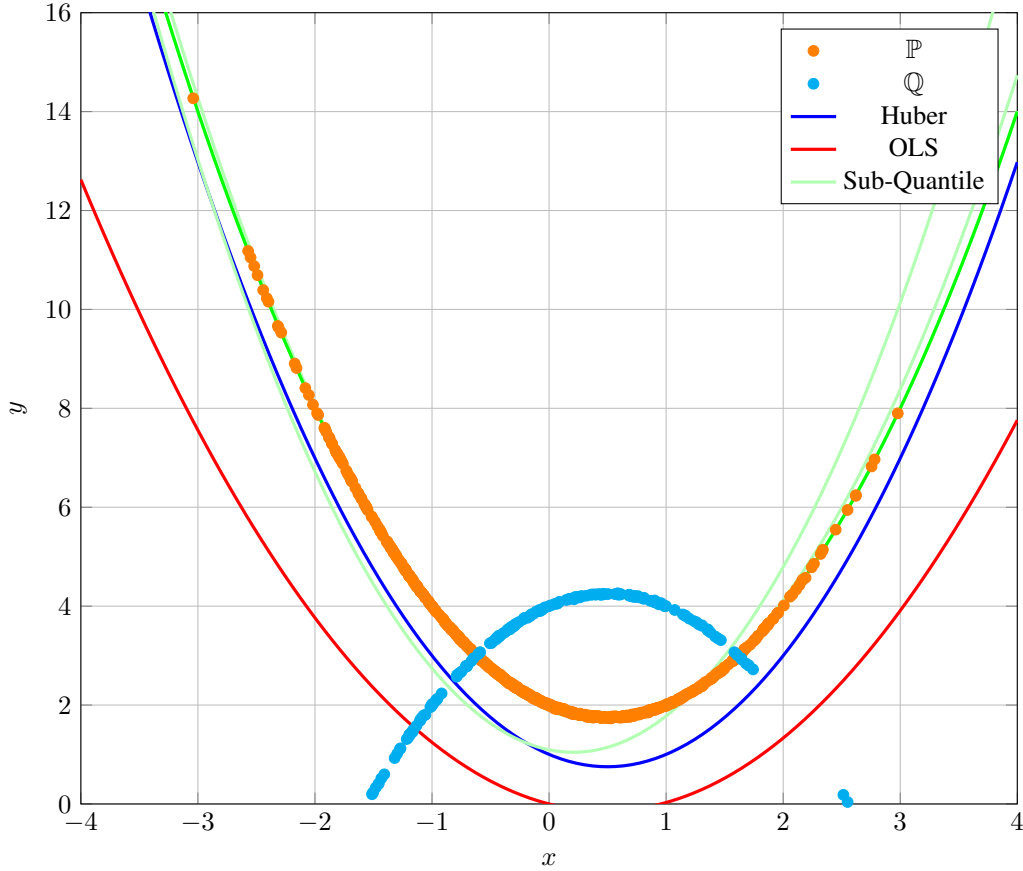
```

5.1 SYNTHETIC DATA

In our first synthetic experiment, we run Algorithm 1 on synthetically generated quadratic data. The results of Sub-Quantile Minimization can be seen in Figure 1. Our results show the near optimal performance of Sub-Quantile Minimization. The results and comparison with other methods can be seen in Table ???. Note we are not interested in $\epsilon \geq 0.5$ as the concept of corruptness becomes unclear. We see in Table ??, 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 F for more details on the Quadratic 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).

Figure 1: Quadratic Regression $n = 1000$ and $\epsilon = 0.2$

Objectives	Test RMSE (Quadratic Regression)		
	$\epsilon = 0$	$\epsilon = 0.2$	$\epsilon = 0.4$
OLS 119	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 2: Qudatic Regression Synthetic Dataset. Empirical Risk over \mathbb{P}

As we can see in Table 3, we obtain state of the art results in the lower range of range of noise, and futher more, we obtain results on par with the current state of the art. This makes our model the strongest among the tested, due to our strength throughout the whole range of noises.

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.

Objectives	Test RMSE (Drug Discovery)			
	$\epsilon = 0$	$\epsilon = 0.2$	$\epsilon = 0.4$	$\epsilon = 0.8$
OLS 119	0.990 _(0.060)	1.969 _(0.118)	2.829 _(0.086)	4.682 _(0.101)
Huber Huber & Ronchetti (2009)	1.326 _(0.096)	1.628 _(0.253)	2.023 _(0.498)	3.442 _(0.581)
RANSAC Fischler & Bolles (1981)	∞	∞	∞	∞
TERM Li et al. (2020)	1.313 _(0.072)	1.334 _(0.105)	1.343 _(0.0740)	1.428 _(0.107)
SEVER Diakonikolas et al. (2019)	1.079 _(0.059)	1.076 _(0.048)	1.067 _(0.091)	3.993 _(0.203)
SubQuantile($p = 1 - \epsilon$)	1.052 _(0.062)	1.060 _(0.065)	1.073 _(0.101)	1.479 _(0.0695)
Genie ERM	0.990 _(0.060)	1.038 _(0.041)	1.037 _(0.086)	∞

Table 3: Drug Discovery Dataset. Empirical Risk over \mathbb{P}

AUTHOR CONTRIBUTIONS

ACKNOWLEDGMENTS

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>.
- 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.
- 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>.
- Gene H. Golub and Charles F. Van Loan. *Matrix Computations (3rd Ed.)*. Johns Hopkins University Press, USA, 1996. ISBN 0801854148.
- 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>.
- 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>.
- 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.
- 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>.
- Steven W Nydick. The wishart and inverse wishart distributions. *Electronic Journal of Statistics*, 6 (1-19), 2012.
- 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.
- 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>.

A	Assumptions	11
B	Proofs in Section 3	11
B.1	Proof of Theorem 1	11
B.2	Proof of Theorem 2	12
C	General Properties of Sub-Quantile Minimization	13
C.1	Proof of Lemma 4.2	13
C.2	Proof of Lemma 4.3	13
C.3	Proof of Lemma 4.5	14
D	Proofs for Convergence	15
D.1	Proof of Lemma 4.6	15
E	Proofs for Reducing ε	16
E.1	Proof of Theorem 3	16
F	Experimental Details	18
F.1	Quadratic Regression Dataset	18
F.2	Drug Discovery Dataset	18
F.3	Baseline Methods in Section 5	18

A ASSUMPTIONS

B PROOFS IN SECTION 3

B.1 PROOF OF THEOREM 1

Proof.

We will first calculate the pseudo-inverse

$$\mathbf{X}^\top \mathbf{X} = (\mathbf{P}^\top \quad \mathbf{Q}^\top) \begin{pmatrix} \mathbf{P} \\ \mathbf{Q} \end{pmatrix} \quad (33)$$

$$= \mathbf{P}^\top \mathbf{P} + \mathbf{Q}^\top \mathbf{Q} \quad (34)$$

Now we can calculate the Moore-Penrose Inverse

$$(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = (\mathbf{P}^\top \mathbf{P} + \mathbf{Q}^\top \mathbf{Q})^{-1} (\mathbf{P}^\top \quad \mathbf{Q}^\top) \quad (35)$$

$$= ((\mathbf{P}^\top \mathbf{P} + \mathbf{Q}^\top \mathbf{Q})^{-1} \mathbf{P}^\top \quad (\mathbf{P}^\top \mathbf{P} + \mathbf{Q}^\top \mathbf{Q})^{-1} \mathbf{Q}^\top) \quad (36)$$

Now we solve for the optimal model

$$\mathbf{X}^\dagger \mathbf{y} = ((\mathbf{P}^\top \mathbf{P} + \mathbf{Q}^\top \mathbf{Q})^{-1} \mathbf{P}^\top \quad (\mathbf{P}^\top \mathbf{P} + \mathbf{Q}^\top \mathbf{Q})^{-1} \mathbf{Q}^\top) \begin{pmatrix} \mathbf{y}_P \\ \mathbf{y}_Q \end{pmatrix} \quad (37)$$

$$= (\mathbf{P}^\top \mathbf{P} + \mathbf{Q}^\top \mathbf{Q})^{-1} \mathbf{P}^\top \mathbf{y}_P + (\mathbf{P}^\top \mathbf{P} + \mathbf{Q}^\top \mathbf{Q})^{-1} \mathbf{Q}^\top \mathbf{y}_Q \quad (38)$$

By assumption 2, all rows of \mathbf{P} and \mathbf{Q} are sampled from a common Normal Distribution. Thus we are able to utilize properties of the Wishart Distribution, Nydick (2012).

$$\mathbf{P}^\top \mathbf{P} = \sum_{\substack{i=1 \\ n \in}}^{n*(1-\epsilon)} \mathbf{P}_i \mathbf{P}_i^\top \quad (39)$$

$$\mathbf{Q}^\top \mathbf{Q} = \sum_{j=1}^n \mathbf{Q}_j \mathbf{Q}_j^\top \quad (40)$$

Thus we can say $\mathbf{P}^\top \mathbf{P}$ and $\mathbf{Q}^\top \mathbf{Q}$ are sampled from the Wishart distribution.

$$\mathbf{P}^\top \mathbf{P} \sim \mathcal{W}(n(1-\epsilon), \Sigma) \quad (41)$$

$$\mathbf{Q}^\top \mathbf{Q} \sim \mathcal{W}(n\epsilon, \Sigma) \quad (42)$$

We can now use the Expected Value of the Wishart Distribution.

$$\mathbb{E}(\mathbf{P}^\top \mathbf{P}) = n(1-\epsilon)\Sigma \quad (43)$$

$$\mathbb{E}(\mathbf{Q}^\top \mathbf{Q}) = n\epsilon\Sigma \quad (44)$$

It thus follows

$$\mathbb{E}[\mathbf{P}^\top \mathbf{P} + \mathbf{Q}^\top \mathbf{Q}] = n\Sigma \quad (45)$$

Since we are interested in the pseudo-inverse, we will utilize the Inverse Wishart Distribution.

$$(\mathbf{P}^\top \mathbf{P} + \mathbf{Q}^\top \mathbf{Q})^{-1} \sim \mathcal{W}^{-1}(n, \Sigma) \quad (46)$$

It thus follows by the expectation of the Inverse Wishart Distribution

$$\mathbb{E}[(\mathbf{P}^\top \mathbf{P} + \mathbf{Q}^\top \mathbf{Q})^{-1}] = n\Sigma^{-1} \quad (47)$$

Now we will plug this into Equation 38:

$$\mathbb{E}[\mathbf{X}^\dagger \mathbf{y}] = (n\Sigma^{-1}) \mathbf{P}^\top \mathbf{y}_P + (n\Sigma^{-1}) \mathbf{Q}^\top \mathbf{y}_Q \quad (48)$$

$$= (n\Sigma^{-1}) \mathbf{P}^\top (\mathbf{P}\beta + \epsilon_P) + (n\Sigma^{-1}) \mathbf{Q}^\top (\mathbf{Q}\beta_Q + \epsilon_Q) \quad (49)$$

$$= (n\Sigma^{-1}) ((\mathbf{P}^\top \mathbf{P})\beta_P + (\mathbf{Q}^\top \mathbf{Q})(\beta_P + (\beta_Q - \beta_P))) \quad (50)$$

$$= (n\Sigma^{-1}) ((n(1-\epsilon)\Sigma)\beta_P + n\epsilon\Sigma(\beta_P + \Psi)) \quad (51)$$

$$= (n\Sigma^{-1}) (n\Sigma\beta_P + n\epsilon\Sigma\Psi) \quad (52)$$

$$= \beta_P + \epsilon(\Psi) \quad (53)$$

This concludes the proof. \square

B.2 PROOF OF THEOREM 2

Proof. The first half of the proof follows from Appendix B.1. We start by noting new notation. Σ_P represents the covariance matrix for \mathbb{P} and Σ_Q represents the covariance matrix for \mathbb{Q} .

$$\mathbb{E} [\mathbf{P}^\top \mathbf{P}] = n(1 - \epsilon) \Sigma_P \quad (54)$$

$$\mathbb{E} [\mathbf{Q}^\top \mathbf{Q}] = n\epsilon \Sigma_Q \quad (55)$$

It thus follows

$$\mathbb{E} [\mathbf{P}^\top \mathbf{P} + \mathbf{Q}^\top \mathbf{Q}] = (n(1 - \epsilon) \Sigma_P + n\epsilon \Sigma_Q) \quad (56)$$

This is where the structure of the proof because we can no longer follow the Inverse Wishart Distribution.

$$\mathbb{E} [(\mathbf{P}^\top \mathbf{P} + \mathbf{Q}^\top \mathbf{Q})^{-1}] = (n(1 - \epsilon) \Sigma_P + n\epsilon \Sigma_Q)^{-1} \quad (57)$$

Now we can use the Woodbury Formula Golub & Van Loan (1996)

$$= n(1 - \epsilon) \Sigma_P^{-1} - n(1 - \epsilon) \Sigma_P^{-1} (n\epsilon \Sigma_Q^{-1}) \quad (58)$$

We will now calculate the expected optimal parameters by plugging this into Equation 38:

$$\mathbb{E} [\mathbf{X}^\dagger \mathbf{y}] = n(1 - \epsilon) \Sigma_P^{-1} (\mathbf{P}^\top \mathbf{P}) \beta_P - n(1 - \epsilon) \Sigma_P^{-1} (n\epsilon \Sigma_Q^{-1}) (\mathbf{Q}^\top \mathbf{Q}) \beta_P \quad (59)$$

$$= n(1 - \epsilon) \Sigma_P^{-1} (n(1 - \epsilon) \Sigma_P) \beta_P - n(1 - \epsilon) \Sigma_P^{-1} (n\epsilon \Sigma_Q^{-1}) (n\epsilon \Sigma_Q) \beta_P \quad (60)$$

$$= \beta_P - n(1 - \epsilon) \Sigma_P^{-1} \beta_P \quad (61)$$

This concludes the proof. \square

C GENERAL PROPERTIES OF SUB-QUANTILE MINIMIZATION

C.1 PROOF OF LEMMA 4.2

Proof. 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 $\hat{\nu} = \text{sorted}((\theta^\top \mathbf{X} - \mathbf{y})^2)$ and note $0 < p < 1$

$$\partial(-g(t, \theta)) = \partial\left(-t + \frac{1}{np} \sum_{i=1}^n (t - \hat{\nu}_i)^+\right) \quad (62)$$

$$= \partial(-t) + \partial\left(\frac{1}{np} \sum_{i=1}^n (t - \hat{\nu}_i)^+\right) \quad (63)$$

$$= -1 + \frac{1}{np} \sum_{i=1}^n \partial(t - \hat{\nu}_i)^+ \quad (64)$$

$$= -1 + \frac{1}{np} \sum_{i=1}^n \begin{cases} 1, & \text{if } t > \hat{\nu}_i \\ 0, & \text{if } t < \hat{\nu}_i \\ [0, 1], & \text{if } t = \hat{\nu}_i \end{cases} \quad (65)$$

$$= 0 \text{ when } t = \hat{\nu}_{np} \quad (66)$$

This is the p -quantile of ν . Not necessarily the p -quantile of $Q_p(U)$ \square

C.2 PROOF OF LEMMA 4.3

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

$$\nabla_{\theta_k} g(t_{k+1}, \theta_k) = \nabla_{\theta_k} \left(\nu_{np} - \frac{1}{np} \sum_{i=1}^n (\nu_{np} - (\theta_k^\top \mathbf{x}_i - y_i)^2)^+ \right) \quad (67)$$

$$= \nabla_{\theta_k} \left((\theta_k^\top \mathbf{x}_{np} - y_{np})^2 - \frac{1}{np} \sum_{i=1}^n ((\theta_k^\top \mathbf{x}_{np} - y_{np})^2 - (\theta_k^\top \mathbf{x}_i - y_i)^2)^+ \right) \quad (68)$$

$$= \nabla_{\theta_k} (\theta_k^\top \mathbf{x}_{np} - y_{np})^2 - \frac{1}{np} \sum_{i=1}^n \nabla_{\theta_k} ((\theta_k^\top \mathbf{x}_{np} - y_{np})^2 - (\theta_k^\top \mathbf{x}_i - y_i)^2)^+ \quad (69)$$

$$= 2\mathbf{x}_{np}(\theta_k^\top \mathbf{x}_{np} - y_{np}) - \frac{1}{np} \sum_{i=1}^n 2\mathbf{x}_{np}(\theta_k^\top \mathbf{x}_{np} - y_{np}) - 2\mathbf{x}_i(\theta_k^\top \mathbf{x}_i - y_i) \begin{cases} 1, & \text{if } t > v_i \\ 0, & \text{if } t < v_i \\ [0, 1], & \text{if } t = v_i \end{cases} \quad (70)$$

$$= 2\mathbf{x}_{np}(\theta_k^\top \mathbf{x}_{np} - y_{np}) - \frac{1}{np} \sum_{i=1}^{np} 2\mathbf{x}_{np}(\theta_k^\top \mathbf{x}_{np} - y_{np}) - 2\mathbf{x}_i(\theta_k^\top \mathbf{x}_i - y_i) \quad (71)$$

$$= 2\mathbf{x}_{np}(\theta_k^\top \mathbf{x}_{np} - y_{np}) - 2\mathbf{x}_{np}(\theta_k^\top \mathbf{x}_{np} - y_{np}) + \frac{1}{np} \sum_{i=1}^{np} 2\mathbf{x}_i(\theta_k^\top \mathbf{x}_i - y_i) \quad (72)$$

$$= \frac{1}{np} \sum_{i=1}^{np} 2\mathbf{x}_i(\theta_k^\top \mathbf{x}_i - y_i) \quad (73)$$

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

C.3 PROOF OF LEMMA 4.5

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

$$\|\nabla_{\boldsymbol{\theta}} g(\boldsymbol{\theta}', t) - \nabla_{\boldsymbol{\theta}} g(\boldsymbol{\theta}, t)\| \leq L \|\boldsymbol{\theta}' - \boldsymbol{\theta}\| \quad (74)$$

$$\left\| \nabla_{\boldsymbol{\theta}} g(\boldsymbol{\theta}', t) - \nabla_{\boldsymbol{\theta}} g(\boldsymbol{\theta}, t) \right\| = \left\| \frac{1}{np} \sum_{i=1}^{np} 2\mathbf{x}_i(\boldsymbol{\theta}'^T \mathbf{x}_i - y_i) - \frac{1}{np} \sum_{i=1}^{np} 2\mathbf{x}_i(\boldsymbol{\theta}_k^T \mathbf{x}_i - y_i) \right\| \quad (75)$$

$$= \left\| \frac{1}{np} \sum_{i=1}^{np} 2\mathbf{x}_i(\boldsymbol{\theta}'^T \mathbf{x}_i - \boldsymbol{\theta}_k^T \mathbf{x}_i) \right\| \quad (76)$$

$$= \left\| \frac{1}{np} \sum_{i=1}^{np} 2\mathbf{x}_i^T \mathbf{x}_i (\boldsymbol{\theta}'^T - \boldsymbol{\theta}_k^T) \right\| \quad (77)$$

$$\stackrel{\text{Cauchy-Schwarz}}{\leq} \left\| \frac{2}{np} \sum_{i=1}^{np} \|\mathbf{x}_i\|^2 \right\| \left\| \boldsymbol{\theta}'^T - \boldsymbol{\theta}_k^T \right\| \quad (78)$$

$$= L \left\| \boldsymbol{\theta}'^T - \boldsymbol{\theta}_k^T \right\| \quad (79)$$

where $L = \left\| \frac{2}{np} \sum_{i=1}^{np} \|\mathbf{x}_i\|^2 \right\|$

This concludes the proof.

D PROOFS FOR CONVERGENCE

D.1 PROOF OF LEMMA 4.6

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 $t_k \geq \hat{\nu}_k$

$$g(t_k, \theta_k) = t_k - \frac{1}{np} \sum_{i=1}^n (t_k - \nu_i)^+ \quad (80)$$

$$= t_k - \frac{1}{np} (np)t_k + \frac{1}{np} \sum_{i=1}^{np} \nu_i + \frac{1}{np} \sum_{i=np}^n (t_k - \nu_i)^+ \quad (81)$$

$$= \frac{1}{np} \sum_{i=1}^{np} \nu_i + \frac{1}{np} \sum_{i=np}^n (t_k - \nu_i)^+ \quad (82)$$

$$g(t_{k+1}, \theta_k) - g(t_k, \theta_k) = \frac{1}{np} \sum_{i=1}^{np} \nu_i - \left(\frac{1}{np} \sum_{i=1}^{np} \nu_i + \frac{1}{np} \sum_{i=np}^n (t_k - \nu_i)^+ \right) \quad (83)$$

$$= -\frac{1}{np} \sum_{i=np}^n (t_k - \nu_i)^+ \quad (84)$$

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 $\delta \in (0, p)$.

$$g(t_k, \theta_k) = t_k - \frac{1}{np} \sum_{i=1}^n (t_k - \nu_i)^+ \quad (85)$$

$$= t_k - \frac{1}{np} \sum_{i=1}^{n(p-\delta)} (t_k - \nu_i)^+ \quad (86)$$

$$= t_k - \frac{1}{np} (n(p-\delta))t_k + \frac{1}{np} \sum_{i=1}^{n(p-\delta)} \nu_i \quad (87)$$

$$g(t_k, \theta_{k+1}) - g(t_k, \theta_k) = \frac{1}{np} \sum_{i=1}^{np} \nu_i - \left(\delta t_k + \frac{1}{np} \sum_{i=1}^{n(p-\delta)} \nu_i \right) \quad (88)$$

$$= \left(\frac{1}{np} \sum_{i=n(p-\delta)}^n \nu_i \right) - \delta t_k \quad (89)$$

□

E PROOFS FOR REDUCING ε

E.1 PROOF OF THEOREM 3

Proof. Note Assumption 2 tells us the whole distribution is sampled from the same distribution.

$$\mathbb{E} [\nabla_{\theta} \phi(\theta)] = \frac{1}{np} \sum_{i=1}^{np(1-\epsilon)} 2\mathbb{E} [\mathbf{p}_i] \mathbb{E} [\theta^\top \mathbf{p}_i - y_i] \quad (90)$$

$$= \frac{1}{np} \sum_{i=1}^{np(1-\epsilon)} 2\mathbb{E} [\mathbf{p}_i] \mathbb{E} [(\theta^\top - \beta_P^\top) \mathbf{p}_i] \quad (91)$$

$$= \frac{2}{np} \sum_{i=1}^{np(1-\epsilon)} \mu(\theta^\top - \beta_P^\top) \mu \quad (92)$$

$$= 2(1-\epsilon) \mu \mu^\top (\theta - \beta_P) \quad (93)$$

We will use similar logic to find the Expectation of the derivative w.r.t \mathbb{Q}

$$\mathbb{E} [\nabla_{\theta} \psi(\theta)] = 2\epsilon \mu \mu^\top (\theta - \beta_Q) \quad (94)$$

Thus we can calculate the expected derivative

$$\mathbb{E} [\nabla_{\theta} g(t_{k+1}, \theta_k)] = \mathbb{E} [\nabla_{\theta} \phi(\theta_k)] + \mathbb{E} [\nabla_{\theta} \psi(\theta_k)] \quad (95)$$

$$= 2(1-\epsilon) \mu \mu^\top (\theta - \beta_P) + 2\epsilon \mu \mu^\top (\theta_k - \beta_Q) \quad (96)$$

$$= 2\mu \mu^\top (\theta_k - (1-\epsilon)\beta_P - \epsilon\beta_Q) \quad (97)$$

Definition 6. A function f is L smooth if

$$\|f(\mathbf{x}) - f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y}$$

We will now find the expected change in loss with respect to the distributions of \mathbb{P} and \mathbb{Q} in each iteration. Furthermore, let us note $g(t, \theta)$ is L -smooth with respect to θ . Thus we can use the properties of Definition 6.

$$\mathbb{E} [g(t_{k+1}, \theta_{k+1})] \leq \mathbb{E} [g(t, \theta_k)] + \langle \nabla_{\theta} \mathbb{E} [g(t_k, \theta_k)], \theta_{k+1} - \theta_k \rangle + \frac{L}{2} \|\theta_{k+1} - \theta_k\|_2^2 \quad (98)$$

$$= \mathbb{E} [g(t, \theta_k)] + \langle \nabla_{\theta} \phi(\theta) + \nabla_{\theta} \psi(\theta), \theta_{k+1} - \theta_k \rangle + \frac{L}{2} \|\theta_{k+1} - \theta_k\|_2^2 \quad (99)$$

We can use the θ update rule described in Equation 11

$$\begin{aligned} &= \mathbb{E} [g(t, \theta_k)] + \langle \nabla_{\theta} \phi(\theta) + \nabla_{\theta} \psi(\theta), -\frac{1}{L} (\nabla_{\theta} \phi(\theta_k) + \nabla_{\theta} \psi(\theta_k)) \rangle \\ &\quad + \frac{L}{2} \|\phi(\theta) + \nabla_{\theta} \psi(\theta)\|_2^2 \end{aligned} \quad (100)$$

We can now use the expected derivative described in Equation 97

$$\begin{aligned} &= \mathbb{E} [g(t, \theta_k)] + \langle 2\mu \mu^\top (\theta^\top - (1-\epsilon)\beta_P^\top - \epsilon\beta_Q^\top), -\frac{1}{L} 2\mu \mu^\top (\theta^\top - (1-\epsilon)\beta_P^\top - \epsilon\beta_Q^\top) \rangle \\ &\quad + \frac{1}{2L} (2\mu \mu^\top (\theta^\top - (1-\epsilon)\beta_P^\top - \epsilon\beta_Q^\top))^\top (2\mu \mu^\top (\theta^\top - (1-\epsilon)\beta_P^\top - \epsilon\beta_Q^\top)) \end{aligned} \quad (101)$$

$$= \mathbb{E} [g(t, \theta_k)] - \frac{1}{2L} (2\mu \mu^\top (\theta^\top - (1-\epsilon)\beta_P^\top - \epsilon\beta_Q^\top))^\top (2\mu \mu^\top (\theta^\top - (1-\epsilon)\beta_P^\top - \epsilon\beta_Q^\top)) \quad (102)$$

$$(103)$$

Let us now calculate the expected change in θ in each iteration. We will start with the θ -update as described in Equation 11.

$$\mathbb{E} [\theta_{k+1}] = \theta_k - \frac{1}{2L} \mathbb{E} [\nabla_{\theta} g(t_{k+1}, \theta_k)] \quad (104)$$

We can now use the expected derivative in equation 97

$$= \theta_k - \frac{1}{2L} (2\mu\mu^\top (\theta_k - (1-\varepsilon)\beta_P - \varepsilon\beta_Q)) \quad (105)$$

$$= \theta_k - \frac{1}{L} (\mu\mu^\top (\theta_k - (1-\varepsilon)\beta_P - \varepsilon\beta_Q)) \quad (106)$$

$$= \left(I - \frac{1}{L} \mu\mu^\top \right) \theta_k + \frac{1-\varepsilon}{L} \mu\mu^\top \beta_P + \frac{\varepsilon}{L} \mu\mu^\top \beta_Q \quad (107)$$

Let us now calculate $\mathbb{E} [\theta_{k+1}\mu]$

$$\mathbb{E} [\theta_{k+1}^\top \mu] = \left(\left(I - \frac{1}{L} \mu\mu^\top \right) \theta_k + \frac{1-\varepsilon}{L} \mu\mu^\top \beta_P + \frac{\varepsilon}{L} \mu\mu^\top \beta_Q \right)^\top \mu \quad (108)$$

$$= \left(\theta_k^\top \left(I - \frac{1}{L} \mu\mu^\top \right) + \frac{1-\varepsilon}{L} \beta_P^\top \mu\mu^\top + \frac{\varepsilon}{L} \beta_Q^\top \mu\mu^\top \right) \mu \quad (109)$$

$$= \left(\theta_k^\top \mu - \frac{1}{L} \theta_k^\top \mu\mu^\top \mu + \frac{1-\varepsilon}{L} \beta_P^\top \mu\mu^\top \mu + \frac{\varepsilon}{L} \beta_Q^\top \mu\mu^\top \mu \right) \quad (110)$$

let $\mathbb{E} [\mu^\top \mu] = C$ for simplicity

$$= \theta_k^\top \mu - \frac{C}{L} \theta_k^\top \mu + \frac{C(1-\varepsilon)}{L} \beta_P^\top \mu + \frac{C\varepsilon}{L} \beta_Q^\top \mu \quad (111)$$

$$\theta_{k+1}^\top \mu - \theta_k^\top \mu = \left(-\frac{C}{L} \theta_k^\top + \frac{C(1-\varepsilon)}{L} \beta_P^\top + \frac{C\varepsilon}{L} \beta_Q^\top \right) \mu \quad (112)$$

It thus follows nicely

$$\theta_{k+1} - \theta_k = -\frac{C}{L} \theta_k + \frac{C(1-\varepsilon)}{L} \beta_P + \frac{C\varepsilon}{L} \beta_Q \quad (113)$$

This gives us insight into how θ_{k+1} changes. We can now calculate a telescopic sum to see the change in θ after T iterations. We will first assume ε does not change. Furthermore, let us note we start at $\theta_0 = \mathbf{0}$. Let us display a couple of iterations to show how θ_k changes

$$\theta_0 = \mathbf{0} \quad (114)$$

$$\theta_1 = \frac{C(1-\varepsilon)}{L} \beta_P + \frac{C\varepsilon}{L} \beta_Q \quad (115)$$

$$\theta_2 = -\frac{C}{L} \left(\frac{C(1-\varepsilon)}{L} \beta_P + \frac{C\varepsilon}{L} \beta_Q \right) + \frac{C(1-\varepsilon)}{L} \beta_P + \frac{C\varepsilon}{L} \beta_Q \quad (116)$$

$$= \frac{CL(1-\varepsilon)}{L^2} \beta_P + \frac{CL\varepsilon}{L^2} \beta_Q - \frac{C^2(1-\varepsilon)}{L^2} \beta_P - \frac{C^2\varepsilon}{L^2} \beta_Q \quad (117)$$

$$= \frac{(CL - C^2)(1-\varepsilon)}{L^2} \beta_P + \frac{(CL - C^2)\varepsilon}{L^2} \beta_Q \quad (118)$$

This only converges to β_P when ε is decreasing. Otherwise it will simply converge to the OLS optimal solution we proved in theorem 1. \square

F EXPERIMENTAL DETAILS

F.1 QUADRATIC REGRESSION DATASET

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

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

$$n_{\text{train}} \in \{10\text{e}4, 10\text{e}6\}$$

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

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

F.2 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 $\mathcal{N}(5, 5)$

F.3 BASELINE METHODS IN SECTION 5

Here we will describe the objective functions used in the synthetic data experiments.

Ordinary Least Squares (OLS) can be solved utilizing the Moore Penrose Inverse.

$$\mathbf{X}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (119)$$

Huber Regression is solved with the following objective function.

$$L_\delta(y, f(\mathbf{x})) = \begin{cases} \frac{1}{2}(y - f(\mathbf{x}))^2 \\ \delta \cdot (|y - f(\mathbf{x})| - \frac{1}{2}\delta) \end{cases} \quad \text{otherwise} \quad (120)$$

RANSAC