

Pattern Information Processing:¹²⁴ Robust Method

Masashi Sugiyama
(Department of Computer Science)

Contact: W8E-505

sugi@cs.titech.ac.jp

<http://sugiyama-www.cs.titech.ac.jp/~sugi/>

Outliers

125

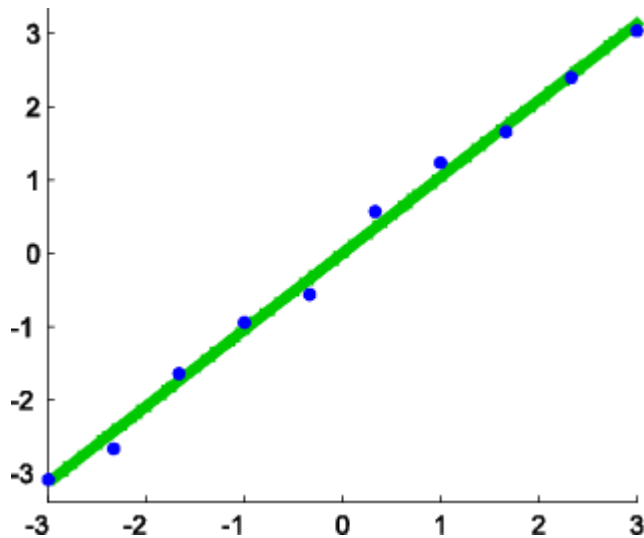
- In practice, very large noise sometimes appears.
- Furthermore, irregular values can be observed by measurement trouble or by human error.
- Samples with such irregular values are called **outliers**.

Outliers (cont.)

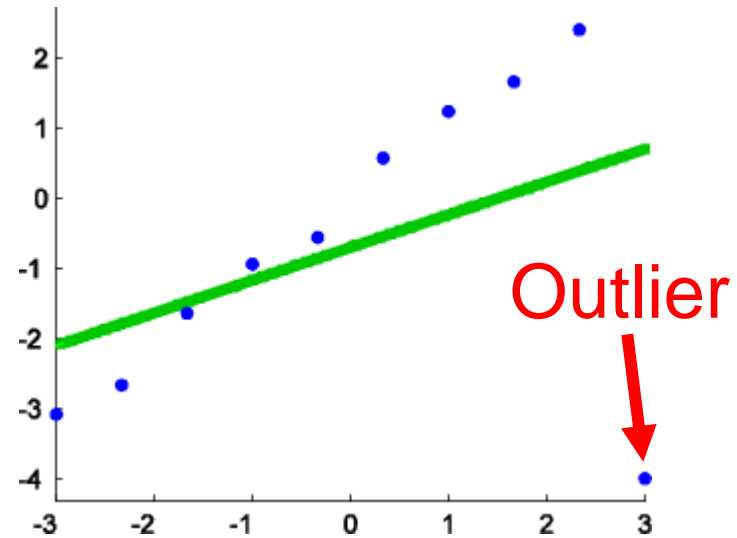
126

- LS criterion is sensitive to outliers.

$$f_{\alpha}(x) = \alpha_1 + \alpha_2 x$$



LS (without outlier)



LS (with outlier)

- Even a single outlier can corrupt the learning result!

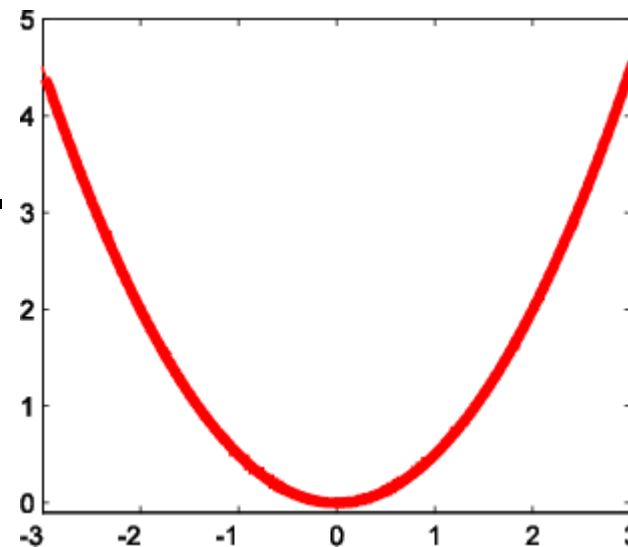
Today's Plan

- Robust learning with ℓ_1 -loss
- Robustness and convexity
- Robustness and efficiency
- Robust learning with Huber's loss
- Robustness and sparsity

Quadratic Loss

$$J_{LS}(\boldsymbol{\alpha}) = \sum_{i=1}^n \left(f_{\boldsymbol{\alpha}}(\boldsymbol{x}_i) - y_i \right)^2$$

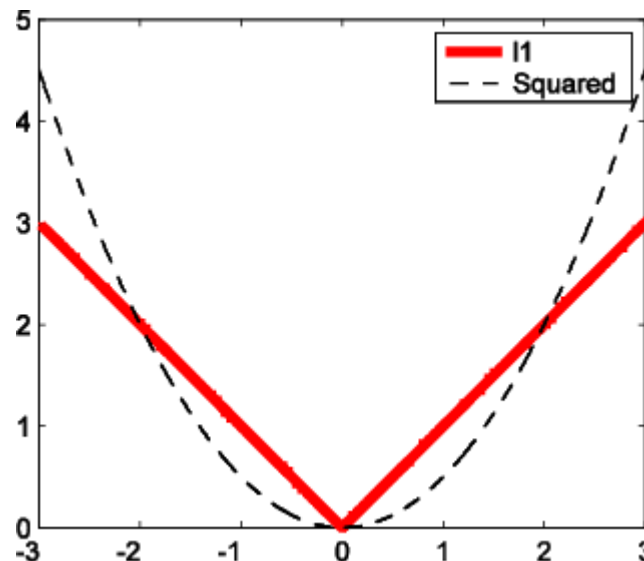
- In LS, goodness-of-fit is measured by the squared loss.
- Therefore, even a single outlier has quadratic power to “pull” the learned function.
- The solution will be robust if outliers are deemphasized.



- Use ℓ_1 -loss for measuring goodness-of-fit:

$$\hat{\alpha}_{\ell_1} = \operatorname{argmin}_{\alpha \in \mathbb{R}^b} \left[\sum_{i=1}^n |f_{\alpha}(x_i) - y_i| \right]$$

- Outliers influence only linearly!



How to Obtain a Solution

130

$$\hat{\alpha}_{\ell_1} = \operatorname{argmin}_{\alpha \in \mathbb{R}^b} \left[\sum_{i=1}^n \left| f_{\alpha}(\mathbf{x}_i) - y_i \right| \right] \quad f_{\alpha}(\mathbf{x}) = \sum_{i=1}^b \alpha_i \varphi_i(\mathbf{x})$$

■ Use the ℓ_1 -trick:

$$|\epsilon| = \min_{v \in \mathbb{R}} v \quad \text{subject to} \quad -v \leq \epsilon \leq v$$

■ $\hat{\alpha}_{\ell_1}$ is given as the solution of the following linearly-constrained linear program:

$$\operatorname{argmin}_{\alpha \in \mathbb{R}^b, \mathbf{v} \in \mathbb{R}^n} \left[\sum_{i=1}^n v_i \right]$$

$$\text{subject to} \quad -\mathbf{v} \leq \mathbf{X}\alpha - \mathbf{y} \leq \mathbf{v}$$

Linearly-Constrained Linear Program (LP)

- Standard optimization software can solve LP:

$$\min_{\beta} \langle \beta, q \rangle \quad \text{subject to } \begin{aligned} H\beta &\leq h \\ G\beta &= g \end{aligned}$$

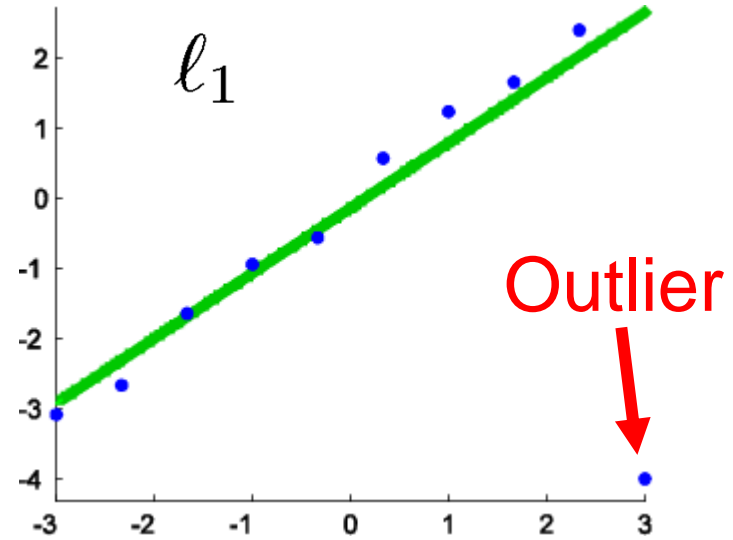
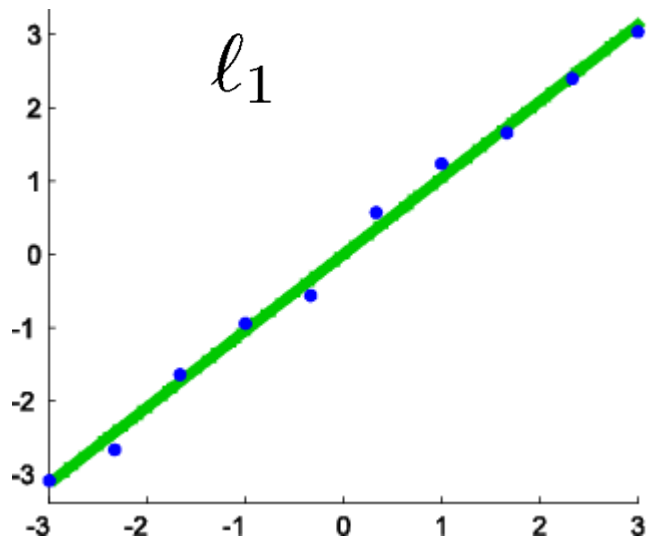
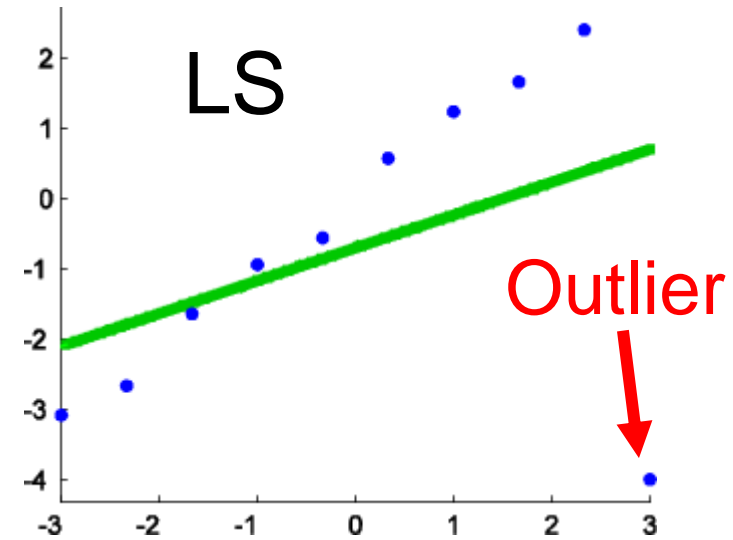
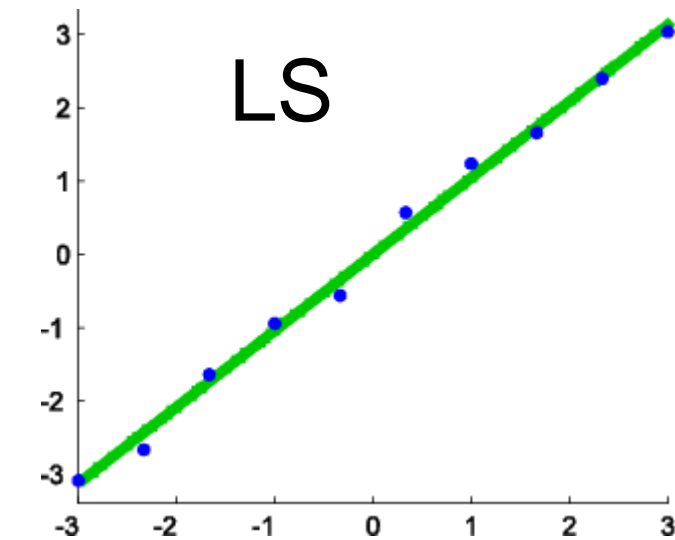
- Let $\beta = \begin{pmatrix} \alpha \\ v \end{pmatrix}$ $\begin{aligned} \Gamma_{\alpha} &= (I_b, O_{b \times n}) \\ \Gamma_v &= (O_{n \times b}, I_n) \end{aligned}$ $\Rightarrow \begin{aligned} \alpha &= \Gamma_{\alpha} \beta \\ v &= \Gamma_v \beta \end{aligned}$

- $\sum_{i=1}^n v_i \Rightarrow \langle \beta, \Gamma_v^{\top} \mathbf{1}_n \rangle$

- $-v \leq X\alpha - y \leq v \Rightarrow \begin{pmatrix} -X\Gamma_{\alpha} - \Gamma_v \\ X\Gamma_{\alpha} - \Gamma_v \end{pmatrix} \beta \leq \begin{pmatrix} -y \\ y \end{pmatrix}$

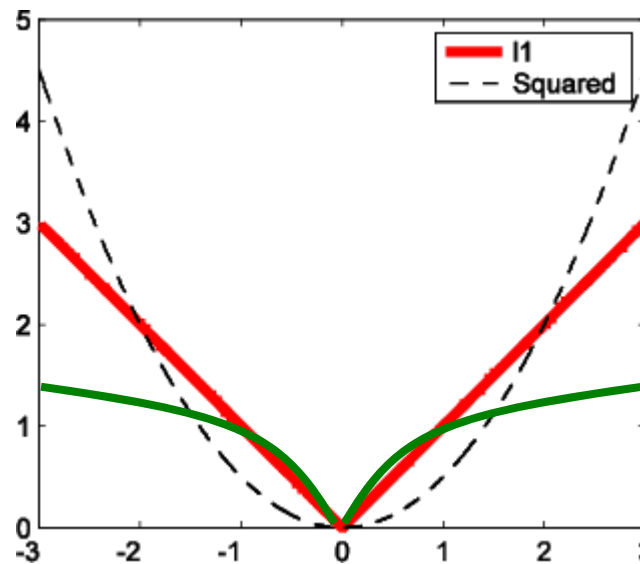
Examples

132



Robustness and Convexity ¹³³

- Influence of outliers can be further reduced by using a **sub-linear** loss:



- However, such a sub-linear loss is **non-convex**.
- Obtaining a global optimal solution is difficult.

Statistical Interpretation

134

- **Data:** Observation = True value + Noise

$$\{y_i \mid y_i = \mu^* + \epsilon_i\}_{i=1}^n$$

- **Goal:** Estimate μ^* from $\{y_i\}_{i=1}^n$.

- ℓ_2 -loss: Sample **mean** is the solution.

$$\hat{\mu}_{\ell_2} = \operatorname{argmin}_{\mu} \left[\sum_{i=1}^n (y_i - \mu)^2 \right] = \operatorname{mean}(\{y_i\}_{i=1}^n)$$

- ℓ_1 -loss: Sample **median** is the solution.

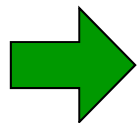
$$\hat{\mu}_{\ell_1} = \operatorname{argmin}_{\mu} \left[\sum_{i=1}^n |y_i - \mu| \right] = \operatorname{median}(\{y_i\}_{i=1}^n)$$

Proof: Homework!

Robustness and Efficiency 135

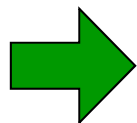
- We move $\alpha\%$ of samples to infinity.
- **Breakdown point**: The maximum α with which a learned function still stays finite.

- ℓ_2 -loss: 0%

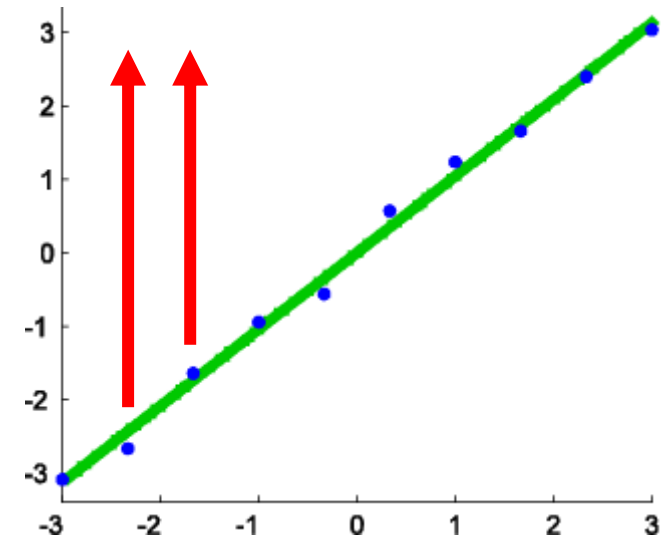


Not at all robust

- ℓ_1 -loss: 50%



Most robust



- However, ℓ_1 -loss is not statistically efficient for Gaussian noise (i.e., having larger variance)

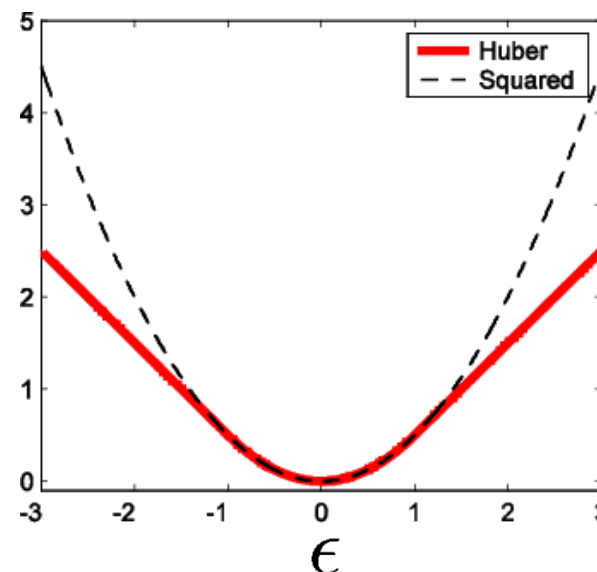
Huber's Robust Learning

136

$$\hat{\alpha}_{Huber} = \operatorname{argmin}_{\alpha \in \mathbb{R}^b} \sum_{i=1}^n \rho(f_{\alpha}(x_i) - y_i)$$

$$\rho(\epsilon) = \begin{cases} \frac{1}{2}\epsilon^2 & (|\epsilon| \leq t) \\ t|\epsilon| - \frac{1}{2}t^2 & (|\epsilon| > t) \end{cases}$$

$$t \geq 0$$



- ℓ_2 -loss for **inliers** (samples with small errors).
- ℓ_1 -loss for **outliers** (samples with large errors).

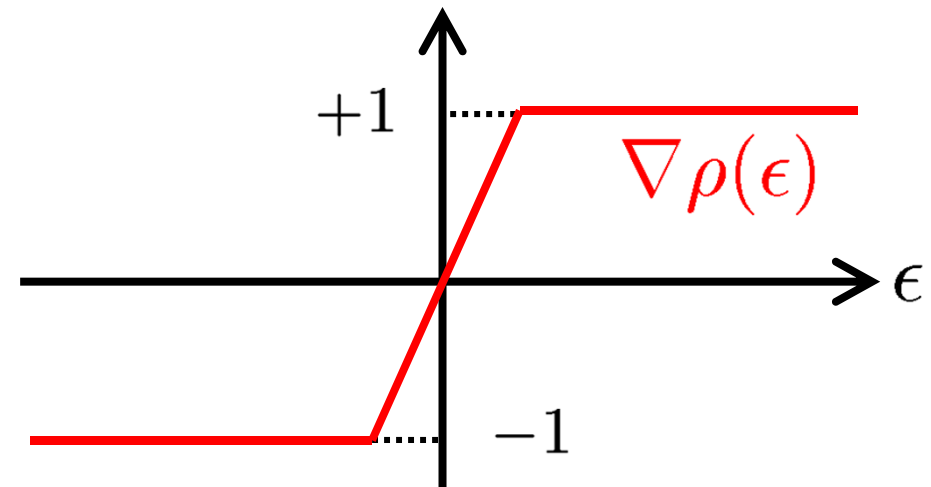
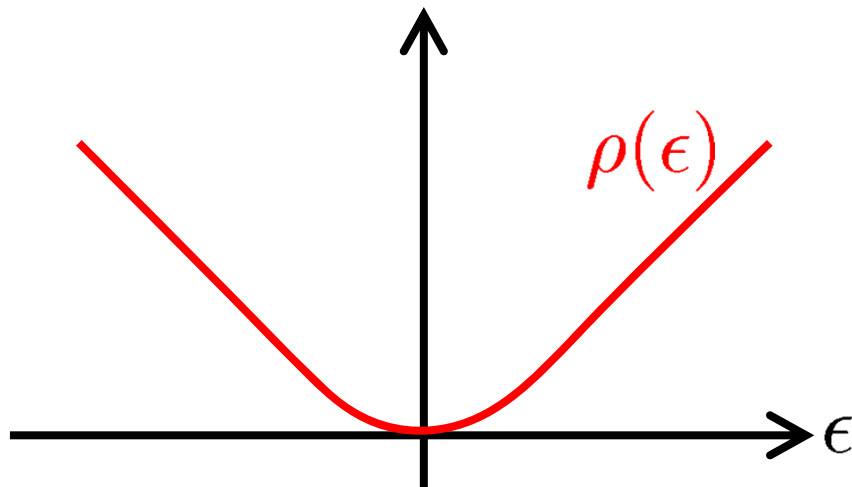
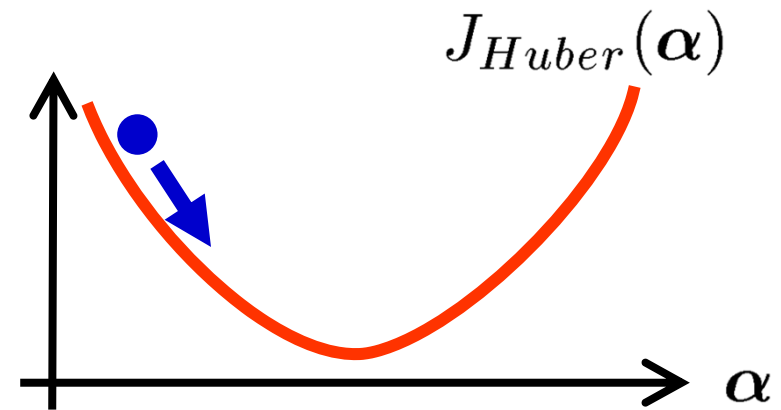
P. J. Huber, Robust Statistics, Wiley, 1981.

How to Obtain A Solution: Gradient Descent

137

$$\alpha \leftarrow \alpha - \epsilon \nabla J_{Huber}(\alpha)$$

$$J_{Huber}(\alpha) = \sum_{i=1}^n \rho(f_{\alpha}(x_i) - y_i)$$



■ A quasi-Newton method may also be used.

Quadratic Program (QP)

138

- Another expression of Huber's loss:

$$\rho(y) = \min_{v \in \mathbb{R}} g(v) \quad g(v) = \frac{1}{2}v^2 + t|y - v|$$

- Then $\hat{\alpha}_{Huber}$ can be obtained as the solution of

$$\min_{\alpha \in \mathbb{R}^b, v \in \mathbb{R}^n} \left[\frac{1}{2} \|v\|^2 + t \|X\alpha - y - v\|_1 \right]$$

- Using the ℓ_1 -trick, this is expressed as QP:

$$\min_{\alpha \in \mathbb{R}^b, u, v \in \mathbb{R}^n} \left[\frac{1}{2} \|v\|^2 + t \sum_{i=1}^n u_i \right]$$

subject to $-u \leq X\alpha - y - v \leq u$

Transforming into Standard Form¹³⁹

$$\min_{\beta} \left[\frac{1}{2} \langle Q\beta, \beta \rangle + \langle \beta, q \rangle \right] \quad \text{subject to } H\beta \leq h$$

$$G\beta = g$$

■ Let $\beta = \begin{pmatrix} \alpha \\ u \\ v \end{pmatrix}$ $\begin{matrix} \Gamma_{\alpha} = (I_b, O_{b \times n}, O_{b \times n}) \\ \Gamma_u = (O_{n \times b}, I_n, O_{n \times n}) \\ \Gamma_v = (O_{n \times b}, O_{n \times n}, I_n) \end{matrix}$ \Rightarrow $\begin{matrix} \alpha = \Gamma_{\alpha}\beta \\ u = \Gamma_u\beta \\ v = \Gamma_v\beta \end{matrix}$

■ $\frac{1}{2} \|v\|^2 + t \sum_{i=1}^n u_i$ \Rightarrow $\frac{1}{2} \langle \Gamma_v^{\top} \Gamma_v \beta, \beta \rangle + \langle \beta, t \Gamma_u^{\top} \mathbf{1}_n \rangle$

■ $-u \leq X\alpha - y - v \leq u$

$\Rightarrow \begin{pmatrix} -X\Gamma_{\alpha} - \Gamma_u + \Gamma_v \\ X\Gamma_{\alpha} - \Gamma_u - \Gamma_v \end{pmatrix} \beta \leq \begin{pmatrix} -y \\ y \end{pmatrix}$

Robustness and Sparseness¹⁴⁰

- Huber's method does not generally provide a sparse solution.
- Combining Huber's loss with ℓ_1 -penalty:

$$\hat{\alpha}_{SparseHuber} = \operatorname{argmin}_{\alpha \in \mathbb{R}^b} \left[\sum_{i=1}^n \rho\left(f_{\alpha}(\mathbf{x}_i) - y_i\right) + \lambda \|\alpha\|_1 \right]$$

- An approximate solution $\hat{\alpha}_{SparseHuber}$ can be obtained by approximate gradient descent.

Linear Programming Learning¹⁴¹

- Combine ℓ_1 -loss and ℓ_1 -constraint:

$$\hat{\alpha}_{\ell_1} = \operatorname{argmin}_{\alpha \in \mathbb{R}^b} \left[\sum_{i=1}^n |f_{\alpha}(x_i) - y_i| + \lambda \|\alpha\|_1 \right]$$

- Using the ℓ_1 -trick, we can obtain $\hat{\alpha}_{LP}$ as the solution of the following LP:

$$\operatorname{argmin}_{\alpha, u \in \mathbb{R}^b, v \in \mathbb{R}^n} \left[\sum_{i=1}^n v_i + \lambda \sum_{i=1}^b u_i \right]$$

$$\text{subject to } -v \leq X\alpha - y \leq v$$

$$-u \leq \alpha \leq u$$

Transforming into Standard Form ¹⁴²

$$\min_{\beta} \langle \beta, q \rangle \quad \text{subject to } H\beta \leq h$$

$$G\beta = g$$

■ Let $\beta = \begin{pmatrix} \alpha \\ u \\ v \end{pmatrix}$ $\Gamma_{\alpha} = (I_b, O_{b \times b}, O_{b \times n})$ $\Gamma_u = (O_{b \times b}, I_b, O_{b \times n})$ $\Gamma_v = (O_{n \times b}, O_{n \times b}, I_n)$ \Rightarrow $\alpha = \Gamma_{\alpha}\beta$ $u = \Gamma_u\beta$ $v = \Gamma_v\beta$

■ $\sum_{i=1}^n v_i + \lambda \sum_{i=1}^b u_i \Rightarrow \langle \beta, \Gamma_v^{\top} \mathbf{1}_n + \lambda \Gamma_u^{\top} \mathbf{1}_b \rangle$

■ $-v \leq X\alpha - y \leq v$ $-u \leq \alpha \leq u$ \Rightarrow $\begin{pmatrix} -X\Gamma_{\alpha} - \Gamma_v \\ X\Gamma_{\alpha} - \Gamma_v \\ -\Gamma_{\alpha} - \Gamma_u \\ \Gamma_{\alpha} - \Gamma_u \end{pmatrix} \beta \leq \begin{pmatrix} -y \\ y \\ 0_b \\ 0_b \end{pmatrix}$

Combinations of Various Losses and Penalties

<div>Penalty</div> <div>Loss</div>	None	ℓ_2	ℓ_1
		Smooth	Smooth & Sparse
ℓ_2 -loss Efficient	Analytic	Analytic	QP, AGD
Huber <div>↕</div>	QP, GD	QP, GD	QP, AGD
ℓ_1 -loss Robust	LP, AGD	QP, AGD	LP, AGD

- QP: Quadratic Program, LP: Linear Program, GD: Gradient Descent, AGD: Approximate GD.

Homework

1. Prove

$$\hat{\mu}_{\ell_2} = \operatorname{argmin}_{\mu} \left[\sum_{i=1}^n (y_i - \mu)^2 \right] = \operatorname{mean}(\{y_i\}_{i=1}^n)$$

$$\hat{\mu}_{\ell_1} = \operatorname{argmin}_{\mu} \left[\sum_{i=1}^n |y_i - \mu| \right] = \operatorname{median}(\{y_i\}_{i=1}^n)$$

under $\{y_i \mid y_i = \mu^* + \epsilon_i\}_{i=1}^n$.

Homework (cont.)

2. For your own toy 1-dimensional data, perform simulations using

- Linear/Gaussian kernel models
- Huber learning

and analyze the results, e.g., by changing

- Target functions
- Number of samples
- Noise level

Including outliers in the dataset would be essential for this homework.