Pattern Information Processing^{1,25} Robust Method

Masashi Sugiyama (Department of Computer Science)

Contact: W8E-505 <u>sugi@cs.titech.ac.jp</u> http://sugiyama-www.cs.titech.ac.jp/~sugi/

Outliers

- 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.)

127

LS criterion is sensitive to outliers.



Even a single outlier can corrupt the learning result!

Today's Plan

Robust learning with l₁ -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.



129

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

$$\hat{oldsymbol{lpha}}_{\ell_1} = \operatorname*{argmin}_{oldsymbol{lpha} \in \mathbb{R}^b} \left[\sum_{i=1}^n \left| f_{oldsymbol{lpha}}(oldsymbol{x}_i) - y_i
ight|
ight]$$

Outliers influence only linearly!



How to Obtain a Solution ¹³¹

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

Use the ℓ_1 -trick:

$$|y| = \min_{v \in \mathbb{R}} v$$
 subject to $-v \le y \le v$

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

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

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

Linearly-Constrained Linear Program (LP)

132

Standard optimization software can solve LP:

$$rac{\min \langle oldsymbol{eta}, oldsymbol{q}
angle}{oldsymbol{eta}} ~~ rac{ ext{subject to } oldsymbol{H}oldsymbol{eta} \leq oldsymbol{h}}{oldsymbol{G}oldsymbol{eta} = oldsymbol{g}}$$

$$\begin{array}{c} \textbf{Let} \quad \boldsymbol{\beta} = \begin{pmatrix} \boldsymbol{\alpha} \\ \boldsymbol{v} \end{pmatrix} \begin{pmatrix} \boldsymbol{\Gamma}_{\boldsymbol{\alpha}} = (\boldsymbol{I}_{b}, \boldsymbol{O}_{b \times n}) \\ \boldsymbol{\Gamma}_{\boldsymbol{v}} = (\boldsymbol{O}_{n \times b}, \boldsymbol{I}_{n}) \end{pmatrix} \quad \boldsymbol{\rho} \quad \boldsymbol{\alpha} = \boldsymbol{\Gamma}_{\boldsymbol{\alpha}} \boldsymbol{\beta} \\ \boldsymbol{v} = \boldsymbol{\Gamma}_{\boldsymbol{v}} \boldsymbol{\beta} \\ \bullet \quad \sum_{i=1}^{n} v_{i} \quad \boldsymbol{\rho} \quad \boldsymbol{\beta} \quad \boldsymbol{\beta}, \boldsymbol{\Gamma}_{\boldsymbol{v}}^{\top} \boldsymbol{1}_{n} \\ \bullet \quad -\boldsymbol{v} \leq \boldsymbol{X} \boldsymbol{\alpha} - \boldsymbol{y} \leq \boldsymbol{v} \quad \boldsymbol{\rho} \quad \boldsymbol{\rho} \quad \boldsymbol{\gamma} \quad \boldsymbol{\beta} \leq \begin{pmatrix} -\boldsymbol{y} \\ \boldsymbol{y} \end{pmatrix} \\ \end{array}$$



Robustness and Convexity ¹³⁴

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.

135
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.

$$\hat{\mu}_{\ell_2} = \underset{\mu}{\operatorname{argmin}} \left[\sum_{i=1}^n (y_i - \mu)^2\right] = \operatorname{mean}\left(\{y_i\}_{i=1}^n\right)$$
 ℓ_1 -loss: Sample median is the solution.

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

Proof: Homework!

Robustness and Efficiency ¹³⁶

We move α % of samples to infinity.

Breakdown point: The maximum α with which a learned function still stays finite.



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

Huber's Robust Learning ¹³⁷

$$\hat{\boldsymbol{\alpha}}_{Huber} = \underset{\boldsymbol{\alpha} \in \mathbb{R}^{b}}{\operatorname{argmin}} \sum_{i=1}^{p} \rho \left(f_{\boldsymbol{\alpha}}(\boldsymbol{x}_{i}) - y_{i} \right)$$
Huber, Robust Statistics (Wiley, 1981)
$$\int \frac{1}{2} y^{2} \quad (|y| \leq t)$$

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

$$t \ge 0$$

$$t \ge 0$$

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



A quasi-Newton method may also be used.

Quadratic Program (QP) ¹³⁹

Another expression of Huber's loss:

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

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

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

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

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

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

Transforming into Standard For¹⁴⁰

$$\min_{oldsymbol{eta}} \left[rac{1}{2} \langle oldsymbol{Q}oldsymbol{eta}, oldsymbol{eta}
ight
angle + \langle oldsymbol{eta}, oldsymbol{q}
ight
angle \left[rac{1}{2} \langle oldsymbol{Q}oldsymbol{eta}, oldsymbol{eta}
ight
angle + \langle oldsymbol{eta}, oldsymbol{q}
ight
angle \left[rac{1}{2} \langle oldsymbol{Q}oldsymbol{eta}, oldsymbol{eta}
ight
angle + \langle oldsymbol{eta}, oldsymbol{q}
ight
angle \left[rac{1}{2} \langle oldsymbol{Q}oldsymbol{eta}, oldsymbol{eta}
ight
angle + \langle oldsymbol{eta}, oldsymbol{q}
ight
angle \left[rac{1}{2} \langle oldsymbol{Q}oldsymbol{eta}, oldsymbol{eta}
ight
angle + \langle oldsymbol{eta}, oldsymbol{q}
ight
angle \left[rac{1}{2} \langle oldsymbol{Q}oldsymbol{eta}, oldsymbol{eta}
ight
angle + \langle oldsymbol{eta}, oldsymbol{q}
ight
angle \left[rac{1}{2} \langle oldsymbol{Q}oldsymbol{eta}, oldsymbol{eta}
ight
angle + \langle oldsymbol{eta}, oldsymbol{q}
ight
angle \left[rac{1}{2} \langle oldsymbol{Q} oldsymbol{eta}, oldsymbol{eta}
ight
angle + \langle oldsymbol{eta}, oldsymbol{q}
ight
angle \left[rac{1}{2} \langle oldsymbol{Q} oldsymbol{eta}, oldsymbol{eta}
ight
angle + \langle oldsymbol{eta}, oldsymbol{q}
ight
angle
ight
angle \left[rac{1}{2} \langle oldsymbol{Q} oldsymbol{eta}, oldsymbol{eta}
ight
angle
ight
angle \left[rac{1}{2} \langle oldsymbol{B} oldsymbol{eta}, oldsymbol{eta}
ight
angle
ight
angle
ight
angle
ight
angle
ight
angle \left[rac{1}{2} \langle oldsymbol{B} oldsymbol{eta}, oldsymbol{eta}
ight
angle
ight
angle \left[rac{1}{2} \langle oldsymbol{A} oldsymbol{eta}, oldsymbol{eta}
ight
angle
ight
ang$$

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

$$egin{aligned} -u \leq Xlpha - y - v \leq u \ & igcap_{eta} & igcap_{eta} & igcap_{eta} & igcap_{eta} - \Gamma_{eta} + \Gamma_{eta} \ X\Gamma_{eta} - \Gamma_{eta} - \Gamma_{eta} & igcap_{eta} & igcap_{eta} \leq igg(egin{aligned} -y \ y \end{pmatrix} eta \leq igcap_{eta} & igcap_{eta} \end{pmatrix}$$

Robustness and Sparseness¹⁴¹

Huber's method does not generally provide a sparse solution.

Combining Huber's loss with ℓ_1 -penalty:

$$\hat{oldsymbol{lpha}}_{SparseHuber} = \operatorname*{argmin}_{oldsymbol{lpha} \in \mathbb{R}^b} \left[\sum_{i=1}^n
ho \Big(f_{oldsymbol{lpha}}(oldsymbol{x}_i) - y_i \Big) + \lambda \|oldsymbol{lpha}\|_1
ight]$$

Quadratic program (QP)

Linear Programming Learning¹⁴² Combine ℓ_1 -loss and ℓ_1 -constraint:

$$\hat{oldsymbol{lpha}}_{\ell_1} = \operatorname*{argmin}_{oldsymbol{lpha} \in \mathbb{R}^b} \left| f_{oldsymbol{lpha}}(oldsymbol{x}_i) - y_i
ight| + \lambda \|oldsymbol{lpha}\|_1$$

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

$$egin{argmin} rgmin \ egin{argmin} \lambda, oldsymbol{u} \in \mathbb{R}^{b}, oldsymbol{v} \in \mathbb{R}^{n} \ egin{argmin} \sum_{i=1}^{n} v_{i} + \lambda \sum_{i=1}^{b} u_{i} \ egin{array} \sum_{i=1}^{n} v_{i} & b \ egin{array} \sum_{i=1}^{n} v_{i} \ egin{array} \sum_{i=1}^{n} v_{i} & b \ egin{array} \sum_{i=1}^{n} v_{i} \ egin{ar$$

Transforming into Standard Form³ $\min_{oldsymbol{eta}} \langle oldsymbol{eta}, oldsymbol{q} angle \ ext{ subject to } oldsymbol{H}oldsymbol{eta} \leq oldsymbol{h} \ oldsymbol{G}oldsymbol{eta} = oldsymbol{g}$ Let $\beta = \begin{pmatrix} \alpha \\ u \\ v \end{pmatrix} \begin{pmatrix} \Gamma_{\alpha} = (I_b, O_{b \times b}, O_{b \times n}) \\ \Gamma_{u} = (O_{b \times b}, I_b, O_{b \times n}) \\ \Gamma_{v} = (O_{n \times b}, O_{n \times b}, I_n) \end{pmatrix} \Rightarrow \begin{pmatrix} \alpha = \Gamma_{\alpha} \beta \\ u = \Gamma_{u} \beta \\ v = \Gamma_{v} \beta \end{pmatrix}$ $egin{aligned} -v &\leq Xlpha - y \leq v \ -u &\leq lpha \leq u \end{aligned} egin{aligned} -X\Gamma_{oldsymbol lpha} - \Gamma_{oldsymbol v} \ -\Gamma_{oldsymbol lpha} - \Gamma_{oldsymbol u} \ -\Gamma_{oldsymbol lpha} \ -\Gamma_{oldsymbol lpha} - \Gamma_{oldsymbol u} \ -\Gamma_{oldsymbol a} - \Gamma_{oldsymbol u} \ -\Gamma_{oldsymbol lpha} \ -\Gamma_{oldsymbol lpha} \ -\Gamma_{oldsymbol a} \$

Combinations of ¹⁴⁴ Various Losses and Penalties

	Penalty	None	ℓ_2	ℓ_1
Loss			Smooth	Smooth & Sparse
ℓ_2 -loss	Efficient	Analytic	Analytic	QP, AGD
Huber		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.

Schedule

June 17: Special lecture by Dr. Song Liu

- Change detection
- June 24: Special lecture by Dr. Marthinus Christoffel du Plessis
 - Learning under class-balance change
- July 1: Lecture on Support Vector Machines
 - Application deadline to Mini-Workshop
- July 8: Lecture on Density Ratio Estimation
- July 15 & July 22: Mini-Workshop
- August 1: Final report deadline

Homework

1. Prove

$$\widehat{\mu}_{\ell_2} = \underset{\mu}{\operatorname{argmin}} \left[\sum_{i=1}^n (y_i - \mu)^2 \right] = \operatorname{mean} \left(\{y_i\}_{i=1}^n \right) = \frac{1}{n} \sum_{i=1}^n y_i$$
$$\widehat{\mu}_{\ell_1} = \underset{\mu}{\operatorname{argmin}} \left[\sum_{i=1}^{2m} |y_i - \mu| \right] = \operatorname{median} \left(\{y_i\}_{i=1}^{2m} \right)$$
$$= \frac{1}{2} (y_m + y_{m+1})$$

for
$$y_1 \leq \cdots \leq y_m \leq y_{m+1} \leq \cdots \leq y_{2m}$$

Homework (cont.)

147

- 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.

If matrix Q in the QP standard form is ill-conditioned, you may add a small positive constant to diagonals.