### Robust statistics

Tomáš Pevný

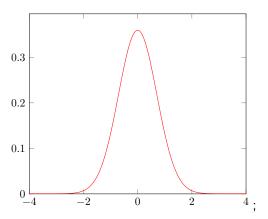
November 25, 2024

#### Goals of robust statistics

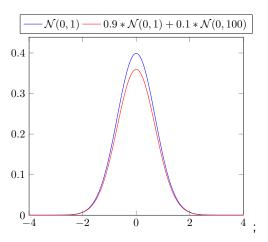
It should not be affected by

- the presence of outliers (even malicious)
- or in-correctness of assumed probability distribution.

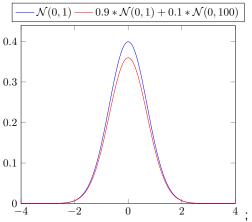
## Which distribution is this?



## Motivation

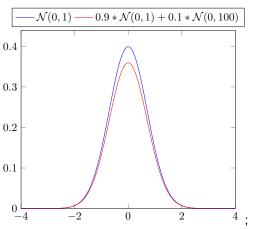


#### Motivation



mean estimated from 1000 samples:  $-5 \cdot 10^{-3}$ , 0.49

#### Motivation



median estimated from 1000 samples: -0.012, -0.013

#### Plan

#### How to compare estimators

Estimators of location

Estimators of scale

M-estimators

Robust regression

Measuring (testing) correlation between variables

Non-parametric tests

#### Breakdown Point

Breakdown Point: the largest proportion of sample observations which may be given arbitrary values without taking the estimator to a limit uninformative about the parameter being estimated.

# Example: Breakdown point

### Breakdown point of

- mean is 0,
- ▶ median is 50%.

# Gross error sensitivity

Influence function

$$\operatorname{IF}(x|p,\eta) = \lim_{\varepsilon \to 0} \frac{\eta((1-\varepsilon)p + \varepsilon \delta_x) - \eta(p)}{\varepsilon}$$

- p probability distribution
- $\eta$  estimator
- $\delta_x$  dirac function at x

# Gross error sensitivity

Influence function

$$\operatorname{IF}(x|p,\eta) = \lim_{\varepsilon \to 0} \frac{\eta((1-\varepsilon)p + \varepsilon \delta_{x}) - \eta(p)}{\varepsilon}$$

Gross error sensitivity

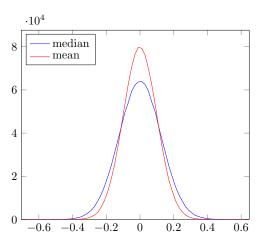
$$GES(p, \eta) = \sup_{x} |IF(x)|$$

- p probability distribution
- $\eta$  estimator
- $\delta_{\mathsf{x}}$  dirac function at x

# How to measure efficiency

How is the sampling distribution of the estimator spread about the true value?

## How to measure efficiency



Distribution of mean and median estimates from 100 samples from  $\mathcal{N}(0,1)$ .

# Asymptotic relative efficiency

Asymptotic relative efficiency (ARE) is defined as

$$ARE(\hat{\eta}_1, \hat{\eta}_2, p) = \frac{V_2}{V_1},$$

where  $\frac{V_1}{n}$ ,  $\frac{V_2}{n}$  are variances of estimators  $\hat{\eta}_1$ ,  $\hat{\eta}_2$  of a parameter  $\mu$  of probability distribution p.

# Example of comparison of efficiency:

Assuming  $x_i$  are i.i.d samples from  $N(\mu, \sigma)$ 

# Example of comparison of efficiency:

Assuming  $x_i$  are i.i.d samples from  $N(\mu, \sigma)$ 

- ►  $Med(X_{n_2}) = med\{x_1, ..., x_{n_2}\} \sim \mathcal{N}\left(\mu, \frac{1}{4p^2(\mu)n_2}\right)$
- Median and mean estimates are equally precise, iff

$$n_1 = \sigma_p^2 4 p^2(\mu) n_2$$

ARE of median to mean for  $\mathcal{N}(0,1)$  is ARE(Med,  $\bar{X}$ ) =  $\frac{2}{\pi}$  = 0.6366.

#### Plan

How to compare estimators

#### Estimators of location

Estimators of scale

M-estimators

Robust regression

Measuring (testing) correlation between variables

Non-parametric tests

### Estimators of location

- mean
- median
- ▶ *q*%-trimmed
- ▶ *q*%-winsorized
- ► Hodges-Lehmann

#### Mean

$$\{-39.61, -26.29, -1.07, -0.92, -0.85, -0.16, 0.93, 1.91, 2.18, 133.65\}$$

- mean  $\frac{1}{n}\sum_{i} x_{i} = 6.97$
- Zero breakdown
- Optimal if samples follows Normal distribution.

#### Median

$$\{-39.61, -26.29, -1.07, -0.92, -0.85, -0.16, 0.93, 1.91, 2.18, 133.65\}$$

- median  $median\{x_1,...,x_{10}\} = -0.51$
- ► 50% breakdown
- ► ARE = 0.637 for Normal distribution

## q%-trimmed

$$\{ \textcolor{red}{-39.61}, \textcolor{red}{-26.29}, \textcolor{blue}{-1.07}, \textcolor{blue}{-0.92}, \textcolor{blue}{-0.85}, \textcolor{blue}{-0.16}, 0.93, 1.91, \textcolor{blue}{2.18}, \textcolor{blue}{133.65} \}$$

- lacktriangle calculate mean from samples  $\left\{x|x_{q\%} \le x \le x_{1-q\%}\right\}$
- ightharpoonup mean  $\frac{1}{|\mathscr{X}_q|}\sum_{X\in\mathscr{X}_q}x_i=-0.41$
- ▶ q% breakdown
- ► ARE = 0.943 for Normal distribution

### q%-Windsorized

$$\{-1.07, -1.07, -1.07, -0.92, -0.85, -0.16, 0.93, 1.91, 1.91, 1.91\}$$

- ▶ replace samples outside  $\langle x_{q\%}, x_{1-q\%} \rangle$  by bounds, return mean.
- q% breakdown
- Robust alternative to mean, more dependant on the distribution then median.

## Hodges-Lehman

$$\{-39.61, -26.29, -1.07, -0.92, -0.85, -0.16, 0.93, 1.91, 2.18, 133.65\}$$

- $\blacktriangleright \text{ HL} = \text{med}\left\{\frac{x_i + x_j}{2} | i, j \in N\right\} = -0.03$
- 0.29 breakdown
- ► ARE = 0.955 for Normal distribution

# Comparison of location estimators

location
6.97
-0.51
-0.41
0.33
-0.03

#### Plan

How to compare estimators

Estimators of location

#### Estimators of scale

M-estimators

Robust regression

Measuring (testing) correlation between variables

Non-parametric tests

### Estimators of scale

- sample standard deviation
- ▶ median absolute deviation
- $\triangleright$   $S_n$
- ► Q

# Sample standard deviation

- (unbiassed) formula:  $\sigma^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i \bar{x})^2$
- (biassed) formula:  $\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i \bar{x})^2$
- breakdown point 0
- ► ARE=1 optimal for Normal distribution



#### Median absolute deviation

- ▶ formula:  $MAD = med\{|x_i med\{x_i\}|\}$
- ▶ breakdown point 50%
- ► For Normal distribution
  - ► ARE=0.37
  - $\hat{\sigma} = 1.4826 \cdot \text{MAD}$

$$S_n$$

- ▶ formula:  $S_n = \text{med}_i\{\text{med}_j|x_i x_j|\}$
- ▶ breakdown point 29%
- ► For Normal distribution
  - ► ARE=0.86
  - $\hat{\sigma} = 1.0483 \cdot S_n$

- ightharpoonup sample standard deviation  $Q = \{|x_i x_j||i < j\}_{q_{25}}$
- ▶ breakdown point 50%
- ► For Normal distribution
  - ► ARE=0.82
  - $\hat{\boldsymbol{\sigma}} = 2.2219 \cdot \boldsymbol{Q}$

#### Plan

How to compare estimators

Estimators of location

Estimators of scale

#### M-estimators

Robust regression

Measuring (testing) correlation between variables

Non-parametric tests

## Question

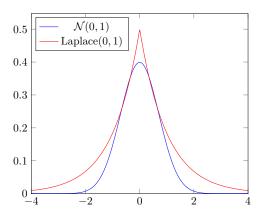
Why mean is so popular?

$$\arg\max_{\mu}\mathscr{L} = \prod_{i} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x_i - \mu)^2}$$

$$\arg\max_{\mu} \mathscr{L} = \prod_{i} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^{2}}(x_{i}-\mu)^{2}}$$

$$\arg\max_{\mu} \log L = -\frac{1}{2\sigma^{2}} \sum_{i} (x_{i}-\mu)^{2} - \log\sqrt{2\pi}\sigma$$

$$\arg\max_{\mu} \mathscr{L} = \prod_{i} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^{2}}(x_{i}-\mu)^{2}}$$
 
$$\arg\max_{\mu} \log L = -\frac{1}{2\sigma^{2}} \sum_{i} (x_{i}-\mu)^{2} - \log\sqrt{2\pi}\sigma$$
 
$$\mu = \frac{1}{n} \sum_{i} x_{i}$$



Assuming  $x \sim \text{Laplace}(\mu, \sigma)$ , then maximum likelihood estimate (ML) of  $\mu$  from  $\{x_1, \ldots, x_n\}$  is

Assuming  $x \sim \text{Laplace}(\mu, \sigma)$ , then maximum likelihood estimate (ML) of  $\mu$  from  $\{x_1, \ldots, x_n\}$  is

$$\arg\max_{\mu}\mathscr{L} = \prod_{i} \frac{1}{2\sigma} e^{-\frac{1}{\sigma}|x_{i} - \mu|}$$

Assuming  $x \sim \text{Laplace}(\mu, \sigma)$ , then maximum likelihood estimate (ML) of  $\mu$  from  $\{x_1, \dots, x_n\}$  is

$$\arg\max_{\mu} \mathscr{L} = \prod_i \frac{1}{2\sigma} e^{-\frac{1}{\sigma}|x_i - \mu|}$$
 
$$\arg\max_{\mu} \log \mathscr{L} = -\frac{1}{\sigma} \sum_i |x_i - \mu| - \log 2\sigma$$

Assuming  $x \sim \text{Laplace}(\mu, \sigma)$ , then maximum likelihood estimate (ML) of  $\mu$  from  $\{x_1, \dots, x_n\}$  is

$$\arg\max_{\mu} \mathscr{L} = \prod_{i} \frac{1}{2\sigma} e^{-\frac{1}{\sigma}|x_{i} - \mu|}$$
 
$$\arg\max_{\mu} \log \mathscr{L} = -\frac{1}{\sigma} \sum_{i} |x_{i} - \mu| - \log 2\sigma$$
 
$$0 = \sum_{i} \operatorname{sgn}(x_{i} - \mu)$$

### Question

Can we generalize this?

$$\arg\max_{\mu} \mathscr{L} = \prod_{i} \frac{1}{Z} e^{-\rho(\frac{x_{i}-\mu}{\sigma})}$$

$$\arg\max_{\mu}\mathscr{L} = \prod_{i} \frac{1}{Z} e^{-\rho(\frac{x_{i}-\mu}{\sigma})}$$

$$\arg\max_{\mu}\log\mathcal{L} = -\sum_{i}\rho(\frac{x_{i}-\mu}{\sigma}) - \log Z$$

$$\arg\max_{\mu} \mathscr{L} = \prod_{i} \frac{1}{Z} e^{-\rho(\frac{x_{i}-\mu}{\sigma})}$$

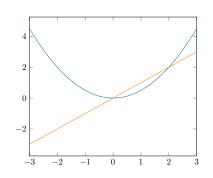
$$\arg\max_{\mu} \log \mathscr{L} = -\sum_{i} \rho(\frac{x_{i}-\mu}{\sigma}) - \log Z$$

$$0 = \sum_{i} \rho'(\frac{x_{i}-\mu}{\sigma}).$$

### Normal distribution

$$ightharpoonup \rho' = x$$

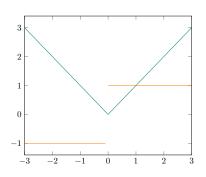
$$\blacktriangleright \mu = \frac{1}{n} \sum_i x_i$$



# Laplace distribution

$$\rho' = \operatorname{sgn}(x)$$

$$0 = \sum_{i} \operatorname{sgn}(x_{i} - \mu)$$



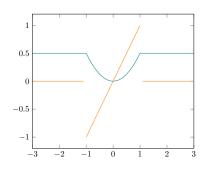
# Huber loss (not called Huber)

$$\rho = \begin{cases} \frac{x^2}{2} & |x| < a \\ \frac{a^2}{2} & |x| \ge a \end{cases}$$

$$\rho' = \begin{cases} x & |x| < a \\ 0 & |x| \ge a \end{cases}$$

$$\blacktriangleright \mu = \frac{1}{n_{< a}} \sum_{i|abs(x_i) < a} x_i$$

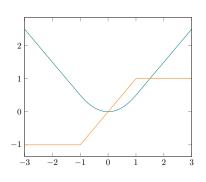
trimming



# Huber loss (called Huber)

$$\mu = \frac{1}{n} \left[ \sum_{i|abs(x_i) < a} x_i + n_{>a} \cdot a \right]$$

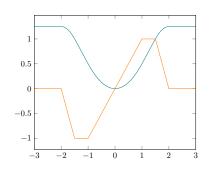
Windsorizing



# Hampel loss

$$\rho = \begin{cases}
\frac{x^2}{2} & 0 \le x < a \\
ax - \frac{a^2}{2} & a \le x < b \\
\frac{a(x-c)^2}{2(b-c)} + \frac{1}{2}a(b+c-a) & b \le x < c \\
\frac{1}{2}a(b+c-a) & c \le x
\end{cases}$$

$$\rho' = \begin{cases} x & 0 \le x < a \\ a & a \le x < b \\ \frac{a(x-c)}{b-c} & b \le x < c \\ 0 & c \le x \end{cases}$$



#### Caveats of robust losses

- ▶ To use them you need to set a scale use robust estimate.
- Robust losses might have unfavourable efficiency.
- Sometimes you need to select parameters.
- Hampel loss is not convex difficult to optimize.

#### Plan

How to compare estimators

Estimators of location

Estimators of scale

M-estimators

#### Robust regression

Measuring (testing) correlation between variables

Non-parametric tests

# Least-square regression is an M-estimator

Generative model behind OLS is

$$y = x^{\mathrm{T}} \beta + \varepsilon, \varepsilon \sim \mathcal{N}(0, \sigma^2).$$

# Least-square regression is an M-estimator

Generative model behind OLS is

$$y = x^{T}\beta + \varepsilon, \varepsilon \sim \mathcal{N}(0, \sigma^{2}).$$

Therefore

$$p(y|x,\beta,\sigma^2) \sim \mathcal{N}(x^T\beta,\sigma^2)$$

## Least-square regression is an M-estimator

#### Generative model behind OLS is

$$y = x^{\mathrm{T}}\beta + \varepsilon, \varepsilon \sim \mathcal{N}(0, \sigma^2).$$

Therefore

$$p(y|x, \beta, \sigma^2) \sim \mathcal{N}(x^T \beta, \sigma^2)$$

and

$$\hat{\beta} = \arg\min_{\beta} \sum_{i} (x_i^{\mathrm{T}} \beta - y_i)^2.$$

### Robust regression

Assume different distribution of noise

$$y = x^{\mathrm{T}} \boldsymbol{\beta} + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim \mathrm{Laplace}(0, \boldsymbol{\sigma}).$$

### Robust regression

Assume different distribution of noise

$$y = x^{\mathrm{T}}\beta + \varepsilon, \varepsilon \sim \mathrm{Laplace}(0, \sigma).$$

and obtain median absolute regression

$$\hat{\beta} = \arg\min_{\beta} \sum_{i} |x_{i}^{\mathrm{T}} \beta - y_{i}|.$$

# Robust regression

Replace the mean estimate by robust alternatives

► Huber A: least median of squares (LMS)

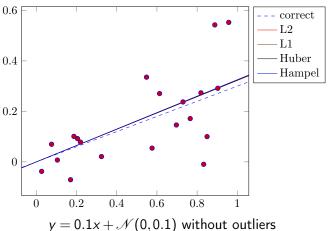
$$\hat{\beta} = \arg\min_{\beta} \operatorname{med} \left\{ (x_i^{\mathrm{T}} \beta - y_i)^2 \right\}$$

Huber B: least trimmed squares (LTS)

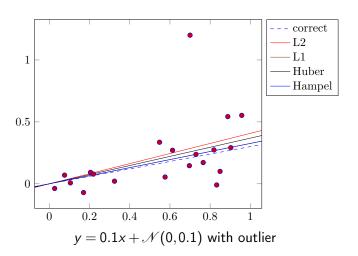
$$\hat{\beta} = \arg\min_{\beta} \sum_{i} (x_i^{\mathrm{T}} \beta - y_i)_{(j)}^2$$



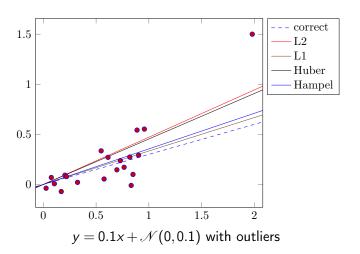
# Examples of robust regression



# Examples of robust regression



# Examples of robust regression



#### Plan

How to compare estimators

Estimators of location

Estimators of scale

M-estimators

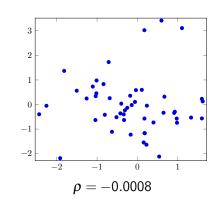
Robust regression

Measuring (testing) correlation between variables

Non-parametric tests

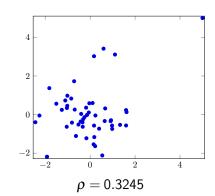
### Pearson's correlation

- ► Assume pairs of samples  $\{(x_i, y_i)\}_{i=1}^n$



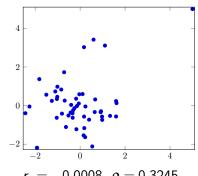
#### Pearson's correlation

- ► Assume pairs of samples  $\{(x_i, y_i)\}_{i=1}^n$
- ► Breakdown point is zero



# Spearman's correlation

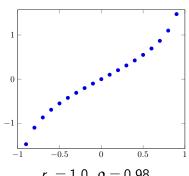
- ▶ Replaces  $\{x_i, y_i\}_i$  by ranks  $\{r_i^x, r_i^y\}_i$



$$r_s = -0.0008, \, \rho = 0.3245$$

# Spearman's correlation

- $\triangleright$  Replaces  $\{x_i, y_i\}_i$  by their ranks  $\{r_i^x, r_i^y\}_i$
- Statistic  $r_s \sqrt{\frac{n-2}{1-r^2}}$  follows Student-t



$$r_s = 1.0, \, \rho = 0.98$$

#### Kendall correlation

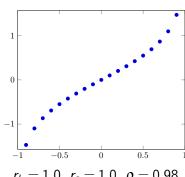
- $\triangleright$  Kendalls' $\tau$  removes all quantities and uses order
- Samples are concordant if

$$ightharpoonup x_i < x_i$$
 and  $y_i < y_i$ 

$$\triangleright$$
  $x_i > x_i$  and  $y_i > y_i$ 

$$r_k = \frac{1}{\binom{n}{2}} (n_c - n_d)$$

$$\qquad \qquad \tau \sim \mathcal{N}\left(0, \frac{2(2N+5)}{9N(N-1)}\right)$$



$$r_k = 1.0, r_s = 1.0, \rho = 0.98$$

#### Kendall correlation

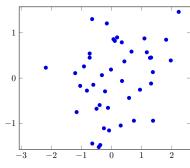
- ightharpoonup Kendalls'au removes all quantities and uses order
- ► Samples are concordant if

$$ightharpoonup x_i < x_i$$
 and  $y_i < y_i$ 

$$\triangleright$$
  $x_i > x_i$  and  $y_i > y_i$ 

$$r_k = \frac{1}{\binom{n}{2}} (n_c - n_d)$$

$$\qquad \qquad \tau \sim \mathcal{N}\left(0, \frac{2(2N+5)}{9N(N-1)}\right)$$



$$r_k = 0.34, \; r_s = 0.48, \; \rho = 0.55$$

#### Plan

How to compare estimators

Estimators of location

Estimators of scale

M-estimators

Robust regression

Measuring (testing) correlation between variables

Non-parametric tests

### Sign test

Tests if differences between pairs of observations are consistent.

# Sign test

Population of pairs  $\{(x_i, y_i)\}_i$ 

- 1. discard samples for which  $|y_i x_i| = 0$
- 2. test statistic

$$W = \sum_{i=1}^{N_r} I(y_i > x_i)$$

3. under null hypothesis W follows binomial distribution  $\operatorname{Bi}(N,0.5)$ 



# Wilcoxon-signed rank test

Tests if population of two related (matched) samples have equal mean rank.

## Test hypothesis of Wilcoxon-signed rank test

Difference between pairs follows a symmetric distribution around zero.

# Wilcoxon-signed rank test

### Population of pairs $\{(x_i, y_i)\}_i$

- 1. calculate  $|y_i x_i|$  and discard those with  $|y_i x_i| = 0$
- 2. rank remaining samples according to  $|y_i x_i|$
- 3. test statistic

$$W = \sum_{i=1}^{N} [\operatorname{sgn}(y_i - x_i) \cdot R_i]$$

- 4. under null hypothesis W has
  - zero mean
  - variance  $\sigma_w^2 = \frac{N(N+1)(2N+1)}{6}$
- 5. For small N critical values are tabulated.
- 6. For large N with  $z = \frac{W}{\sigma_W}, \sigma_W$

## Discussion of sign and signed-rank test

- Sign test have less assumptions needs only order relationship
- Signed rank test have higher power: ARE is 0.67.
- ▶ Would differences follows normal distribution, paired t-test is more appropriate; ARE is 0.95.
- ▶ Generalization of a sign test to *n*-tuples is a Friedman test.

### Mann-Whitney U-test

Tests, whether a probability that a value from population X is greater than a value from population Y (and vice versa) is greater than 0.5.

Tests, whether the distributions of both populations are equal.

## Mann-Whitney U-test

- 1. Assume we have  $\{(x_i)\}_{i=1}^{n_1}$
- 2. Calculate ranks of all samples together.
- 3. Sum ranks of samples from the first population,  $R_1$ .
- 4. Sum ranks of samples from the second population,  $R_2$ .
- 5. Calculate  $U_1 = R_1 \frac{n_1(n_1+1)}{2}$  and  $U_2 = R_2 \frac{n_2(n_2+1)}{2}$ .
- 6.  $U = \min\{U_1, U_2\}$
- 7. For small  $n_1, n_2$  critical values are tabulated, for large  $n_1, n_2$   $U \sim \mathcal{N}\left(\frac{n_1 n_2}{2}, \frac{n_1 n_2 (n_1 + n_2 + 1)}{12}\right)$ .