Robustness: Main ideas

Ruben Zamar Dept. of Stats, UBC

September 9, 2015

Fitting well all the data

versus

Fitting well most of the data

A NUMERICAL

EXAMPLE

æ

-

-

Image: A matrix and a matrix



P opulation (N = 1000)

Image: A matrix



Sample (n=50)

A Prediction Model

• Linear regression model:

$$y_i = \alpha + \beta x_i + \epsilon_i, \quad i = 1, 2, ..., 50$$

< □ > < ---->

A Prediction Model

• Linear regression model:

$$y_i = \alpha + \beta x_i + \epsilon_i, \quad i = 1, 2, ..., 50$$

• Regression residuals

$$r_i(b_0, b_1) = y_i - b_0 - b_1 x_i$$

• Least Squares (LS)

$$\left(\widehat{\alpha}^{LS}, \widehat{\beta}^{LS}\right) = \arg\min_{b_1, b_2} \sum_{i=1}^{n=50} r_i^2 (b_0, b_1)$$

3

• • • • • • • •

• Least Squares (LS)

$$\left(\widehat{lpha}^{LS}, \widehat{eta}^{LS}
ight) = rg\min_{b_1, b_2} \sum_{i=1}^{n=50} r_i^2 \left(b_0, b_1
ight)$$

• For our example, using the function lm in R, we obtain:

$$\widehat{lpha}^{LS}=-0.9237$$
 and $\widehat{eta}^{LS}=3.8808$

Least Squares Fit



LSFit

A Robust Alternative

• Sorted squared residuals

$$r_{(1)}^{2}(b_{0}, b_{1}) \le r_{(2)}^{2}(b_{0}, b_{1}) \le \cdots \le r_{(n)}^{2}(b_{0}, b_{1})$$

A Robust Alternative

• Sorted squared residuals

$$r_{(1)}^{2}(b_{0}, b_{1}) \leq r_{(2)}^{2}(b_{0}, b_{1}) \leq \cdots \leq r_{(n)}^{2}(b_{0}, b_{1})$$

• Least Trimmed Squares (LTS)

$$\left(\widehat{lpha}^{LTS}, \widehat{eta}^{LTS}
ight) = rg\min_{b_1, b_2} \sum_{i=1}^{n(1-lpha)} r_{(i)}^2 \left(b_0, b_1
ight)$$

• For our example we set $\alpha = 0.4$, so

$$\left(\widehat{\alpha}^{LTS}, \widehat{\beta}^{LTS}\right) = \arg\min_{b_1, b_2} \sum_{i=1}^{30} r_{(i)}^2 (b_0, b_1)$$

э

Image: Image:

• For our example we set $\alpha = 0.4$, so

$$\left(\widehat{lpha}^{LTS},\widehat{eta}^{LTS}
ight) = rgmin_{b_1,b_2}\sum_{i=1}^{30}r^2_{(i)}\left(b_0,b_1
ight)$$

~ ~

• Using the function ltsReg in the package robust in R, we obtain:

$$\widehat{lpha}^{LTS}$$
 = 11.662 and \widehat{eta}^{LTS} = 1.875



LTS Fit

< A

æ

• Predict the value of y_i for a new case with covariate value x_i

- Predict the value of y_i for a new case with covariate value x_i
- Use the prediction model and the estimated parameters $\left(\widehat{lpha},\widehat{eta}
 ight)$:

$$\hat{y}_i = \hat{\alpha} + \hat{\beta} x_i$$

Out of Sample Prediction (continued)

• Prediction error

$$e\left(x_{i},\widehat{\alpha},\widehat{\beta}
ight) = y_{i}-\widehat{y}_{i}$$

 $= y_{i}-\widehat{\alpha}-\widehat{\beta}x_{i}$

-

- 一司

æ

Out of Sample Prediction (continued)

• Prediction error

$$e\left(x_{i},\widehat{lpha},\widehat{eta}
ight) = y_{i} - \widehat{y}_{i}$$

 $= y_{i} - \widehat{lpha} - \widehat{eta}x_{i}$

• Absolute prediction error (APE)

$$APE = |y_i - \hat{y}_i| = \left| e\left(x_i, \widehat{\alpha}, \widehat{\beta}\right) \right|$$

• We consider the distribution of APE for the *N* - *n* out-of-sample cases

Image: Image:

3

Distribution of APE

- We consider the distribution of APE for the *N n* out-of-sample cases
- In our example

N - n = 950

3

Distribution of APE

- We consider the distribution of APE for the *N n* out-of-sample cases
- In our example

$$N - n = 950$$

 Use numerical summaries and plots to compare APE's based on LS and LTS

Comparison of α -Quantiles

δ	$q_{ extsf{LTS}}\left(\delta ight)$	$q_{LS}\left(\delta ight)$
0.01	0.05	0.22
0.05	0.32	0.55
0.25	1.69	2.58
0.50	3.55	5.72
0.75	6.31	10.41
0.95	14.87	17.98
0.99	43.21	24.42

APE Densities (LTS vs LS)



APEDistribution: LTSvsLS

APE

æ

APE Boxplots (LTS vs LS)



Absolute Predition Errors

< 17 ▶

э

APE q-q Plots (LTS vs LS)



QQ-Plot: LTSResiduals vs LSResiduals

Absolute LTS residual quantiles

- 一司

3

Classical methods results

versus

Robust methods results

"... just which robust/resistant methods you use is not important – what is important is that you use some..."

"... It is perfectly proper to use both classical and robust/resistant methods routinely, and only worry when they differ enough to matter." "... It is perfectly proper to use both classical and robust/resistant methods routinely, and only worry when they differ enough to matter."

"...when they differ, you should think hard."

A REAL DATA

EXAMPLE

Ruben Zamar Dept. of Stats, UBC

Robustness

September 9, 2015 22 / 45

æ

メロト メポト メヨト メヨト

• Data first published by Brownlee (1965)

Image: Image:

æ

- Data first published by Brownlee (1965)
- Available in R (dataset name = stackloss)

- Data first published by Brownlee (1965)
- Available in R (dataset name = stackloss)
- 21 daily observations of the oxidation of ammonia to nitric acid

- Data first published by Brownlee (1965)
- Available in R (dataset name = stackloss)
- 21 daily observations of the oxidation of ammonia to nitric acid
- Extensively studied in the statistical literature

- Data first published by Brownlee (1965)
- Available in R (dataset name = stackloss)
- 21 daily observations of the oxidation of ammonia to nitric acid
- Extensively studied in the statistical literature
 - Daniel and Wood, 1980, Chapters 5 and 7
- Data first published by Brownlee (1965)
- Available in R (dataset name = stackloss)
- 21 daily observations of the oxidation of ammonia to nitric acid
- Extensively studied in the statistical literature
 - Daniel and Wood, 1980, Chapters 5 and 7
 - Atkinson, 1985, pp. 129-136, 267-8

- Data first published by Brownlee (1965)
- Available in R (dataset name = stackloss)
- 21 daily observations of the oxidation of ammonia to nitric acid
- Extensively studied in the statistical literature
 - Daniel and Wood, 1980, Chapters 5 and 7
 - Atkinson, 1985, pp. 129-136, 267-8
 - Venables and Ripley, 1997

æ

Image: A matrix

Air Flow

The rate flow of cooling air

Air Flow The rate flow of cooling air

Water Temperature Inlet cooling water temperature

Air Flow The rate flow of cooling air

Water Temperature Inlet cooling water temperature

Acid Concentration Conc

Concentration of acid

Air Flow The rate flow of cooling air

Water Temperature Inlet cooling water temperature

Acid Concentration Concentration of acid

Output Variable

Air Flow The rate flow of cooling air

Water Temperature Inlet cooling water temperature

Acid Concentration Concentration of acid

Output Variable

Stack Loss

Inverse measure for the overall efficiency of the plant

Regression Coefficient Estimate	LS	MM
Intercept	-39.92	-37.65
Air Flow	0.716	0.798
Water Temperature	1.300	0.577
Acid Concentration	-0.152	-0.067
Residual SE	3.243	1.837

Big differences!

"We must think hard..."



LS Standardized Residuals

æ

Robust Residual Plot



Standardized Robust Residuals

- 一司

э

• Daniel and Wood (1971, Chapter 5, page 81) noticed a different behavior in the response variable whenever the water temperature was over 60 degrees.

- Daniel and Wood (1971, Chapter 5, page 81) noticed a different behavior in the response variable whenever the water temperature was over 60 degrees.
- The plant needs to stabilize after the water temperature reaches 60 degrees.

- Daniel and Wood (1971, Chapter 5, page 81) noticed a different behavior in the response variable whenever the water temperature was over 60 degrees.
- The plant needs to stabilize after the water temperature reaches 60 degrees.
- They concluded that observations obtained with Water Temperature 2 60 degrees require special attention, and should be removed from the analysis.

- Daniel and Wood (1971, Chapter 5, page 81) noticed a different behavior in the response variable whenever the water temperature was over 60 degrees.
- The plant needs to stabilize after the water temperature reaches 60 degrees.
- They concluded that observations obtained with Water Temperature 2 60 degrees require special attention, and should be removed from the analysis.
- These correspond to cases 1, 3, 4 and 21 directly uncovered by the robust fit.

ANOTHER REAL DATA

EXAMPLE

イロト イポト イヨト イヨト

3

• Socioeconomic-Demographic data on 506 urban districts in the Boston area, USA

- Socioeconomic-Demographic data on 506 urban districts in the Boston area, USA
- Data downloaded from the R package "spdep" (dataset name = boston)

- Socioeconomic-Demographic data on 506 urban districts in the Boston area, USA
- Data downloaded from the R package "spdep" (dataset name = boston)
- Number of Variables = 12

- Socioeconomic-Demographic data on 506 urban districts in the Boston area, USA
- Data downloaded from the R package "spdep" (dataset name = boston)
- Number of Variables = 12
- Number of Cases = 506

Variable	Brief Description
CMEDV	median values of owner-occupied housing in USD 1000
CRIM	per capita crime
INDUS	proportions of non-retail business acres per town
NOX	nitric oxides concentration (parts per 10 million) per town
RM	average numbers of rooms per dwelling
AGE	proportion of owner-occupied units built prior to 1940
DIS	weighted distances to five Boston employment centres
RAD	an index of accessibility to radial highways per town
TAX	full-value property-tax rate per USD 10,000 per town
PTRATIO	pupil-teacher ratios per town (constant for all Boston tracts)
В	$1000*(Bk - 0.63)^2$ where Bk is the proportion of blacks
LSTAT	percentage values of lower status population

æ

< □ > < ---->

Exploratory Data Analysis

• Compute estimates of multivariate location $\pmb{\mu}$ and scatter matrix $\hat{\Sigma}$

Exploratory Data Analysis

- Compute estimates of multivariate location $\pmb{\mu}$ and scatter matrix $\hat{\Sigma}$
- Compute (squared) Mahalanobis distances (MD) for each data point:

$$d_i^2 = (\mathbf{x}_i - \mathbf{\hat{\mu}})' \hat{\Sigma}^{-1} (\mathbf{x}_i - \mathbf{\hat{\mu}})$$

Exploratory Data Analysis

• Compute estimates of multivariate location $\pmb{\mu}$ and scatter matrix $\hat{\Sigma}$

• Compute (squared) Mahalanobis distances (MD) for each data point:

$$d_i^2 = (\mathbf{x}_i - \mathbf{\hat{\mu}})' \hat{\Sigma}^{-1} (\mathbf{x}_i - \mathbf{\hat{\mu}})$$

• Threshold: Compare the MD's with a quantile of a $\chi^2_{(p)}$

$$q = F_{\chi^2_{(p)}}^{-1}(0.99999) = 45.07615$$

Mahalanobis Distances: Using the Sample Mean and Covariance Matrix



Mahalanobis Distances: Using Robust Location-Scatter MVES-Estimates



- Using a robust multivariate-location estimate (called MVE-S) we find 171 outliers
- Using classical sample mean and covariance matrix we find 7 outliers
- Big difference ⇒ Must think hard!
- We will be back to this example

NOISE AND DATA QUALITY

-

イロト イ団ト イヨト イ

3

• Many data can be modeled as follows:

OUTPUT DATA = SIGNAL (INPUT DATA, θ) + NOISE

• We distinguish two types of noise

• We distinguish two types of noise

"TYPICAL" NOISE

0

• We distinguish two types of noise

"TYPICAL" NOISE

ATYPICAL NOISE

2

• Typical noise comes from

NATURAL FLUCTUATIONS

MEASUREMENT ERRORS

ITEM TO ITEM VARIABILITY, ETC

• Typical noise comes from

NATURAL FLUCTUATIONS

MEASUREMENT ERRORS

ITEM TO ITEM VARIABILITY, ETC

• Not necessarily "Gaussian Noise"

• Typical noise comes from

NATURAL FLUCTUATIONS

MEASUREMENT ERRORS

ITEM TO ITEM VARIABILITY, ETC

- Not necessarily "Gaussian Noise"
- Other classical parametric models such as Gamma, Weibull, Poisson, etc
• Atypical noise comes from

OUTLIERS AND GROSS ERRORS MEASUREMENTS OF UNEVEN QUALITY (mixture) DATA CONTAMINATION (mixture) MISSING DATA (declared or unsuspected) DUPLICATIONS, ETC

• Filter noise (both typical and atypical noise)

- Filter noise (both typical and atypical noise)
- Extract the signal (point estimation)

- Filter noise (both typical and atypical noise)
- Extract the signal (point estimation)
- Measure the strength of the noise (statistical inference)

- Filter noise (both typical and atypical noise)
- Extract the signal (point estimation)
- Measure the strength of the noise (statistical inference)
- Assess uncertainty in the estimates (statistical inference)

- Filter noise (both typical and atypical noise)
- Extract the signal (point estimation)
- Measure the strength of the noise (statistical inference)
- Assess uncertainty in the estimates (statistical inference)
- Predict future data (prediction)

• Point Estimates

 $\widehat{\boldsymbol{\theta}}$

• Point Estimates

 $\widehat{\pmb{\theta}}$

• Confidence Regions



Point Estimates

 $\widehat{\boldsymbol{\theta}}$

• Confidence Regions

$$\mathit{Cov}\left(\widehat{oldsymbol{ heta}}
ight)$$
 , Confidence Region for $oldsymbol{ heta}$

• Prediction / Interpolation

$$\widehat{SIGNAL} \pm 2 \times SE\left(\widehat{SIGNAL}\right)$$





∃ → (∃ →

Image: A matrix and a matrix

2

Typically

$$\widehat{\boldsymbol{ heta}}
ightarrow \boldsymbol{ heta}$$

• and

$$Cov\left(\widehat{\theta}\right) = \frac{1}{n}C_{\widehat{\theta}} \to 0$$

э.

Image: A matrix and a matrix

2

• Typically

$$\widehat{\boldsymbol{ heta}}
ightarrow \boldsymbol{ heta}$$

and

$$Cov\left(\widehat{\theta}\right) = \frac{1}{n}C_{\widehat{\theta}} \to 0$$

• Better when $C_{\hat{\theta}}$ is small \implies use efficient procedures

< 67 ▶

3

Typically

$$\widehat{\boldsymbol{ heta}}
ightarrow \boldsymbol{ heta}$$

and

$$Cov\left(\widehat{\theta}\right) = \frac{1}{n}C_{\widehat{\theta}} \to 0$$

- Better when $C_{\widehat{ heta}}$ is small \implies use efficient procedures
- Much attention has been given to the problem of minimizing $C_{\hat{H}}$

• Atypical noise tends to produce asymptotic bias

- Atypical noise tends to produce asymptotic bias
- That is

 $\widehat{\boldsymbol{\theta}} \rightarrow \Delta, \quad \Delta \neq \boldsymbol{\theta}$

- Atypical noise tends to produce asymptotic bias
- That is

$$\widehat{\boldsymbol{ heta}}
ightarrow \Delta$$
, $\Delta \neq \boldsymbol{ heta}$

ullet The difference between Δ and ${m heta}\,$ is called "contamination bias" (cb)

- Atypical noise tends to produce asymptotic bias
- That is

$$\widehat{\boldsymbol{ heta}}
ightarrow \Delta$$
, $\Delta \neq \boldsymbol{ heta}$

• The difference between Δ and θ is called "contamination bias" (cb) • $cb\left(\widehat{\theta}\right)$ is of order 1 while $Cov\left(\widehat{\theta}\right)$ of order 1/n

- Atypical noise tends to produce asymptotic bias
- That is

$$\widehat{oldsymbol{ heta}}
ightarrow \Delta, \quad \Delta
eq oldsymbol{ heta}$$

- The difference between Δ and θ is called "contamination bias" (cb) • $cb\left(\widehat{\theta}\right)$ is of order 1 while $Cov\left(\widehat{\theta}\right)$ of order 1/n
- Therefore, for large n, $cb\left(\widehat{\theta}\right)$ should be a leading concern