# Robust regression

- Recall: $Y=X\beta+\varepsilon$, usually assume error $\varepsilon$ is Normal $\Rightarrow$ ordinary least square (OLS) approach best. **Q**: what if error not Normally distributed?

- Recall: particular concern when errors not Normal $\Rightarrow$ long-tailed error
  $\Rightarrow$ large errors are expected to appear more often
  - OLS not necessary best when large errors exist (**Q**: why? $RSS = \sum(y_i - x_i^T\beta)^2$ )

  - Previous approach: check and remove observations with large residuals, i.e., regard them as outliers, use OLS *after* removing them
    $\Rightarrow$ not effective when there are many outliers because:
    - "leave-out-one" nature in outlier tests
    - not statistically efficient for the estimation of $\beta$
  - Two ways of handling outliers or large errors:
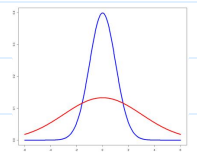    (a) change data, keep model    (b) keep data, change model

- Statistical modeling: $Y=X\beta+\varepsilon$, where error $\varepsilon$ can be modeled as
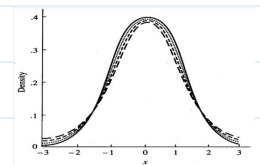  - $\varepsilon \sim$ a mixture distribution, e.g.,
    $\varepsilon \sim \pi\, N(0, \sigma^2) + (1-\pi)\, N(0, c\,\sigma^2)$, $0<\pi<1$ and $c>1$

  

  $t_5$ (long dash), $t_{10}$ (short dash), $t_{30}$ (dot), $N(0,1)$ (solid)

  - $\varepsilon \sim \sigma\, t_d$ distribution with a small $d$

  

  - $\varepsilon \sim$ any distribution with median=0

---

- alternative approach: robust regression (observations are weighted unequally)
  - M-estimators: find $\beta$ to minimize $\Sigma_i\, \rho((y_i - x_i^T\beta)/\sigma)$
    - choice of $\rho$ :
      - $\rho(z) = z^2$ is OLS
      - $\rho(z) = |z|$ is called least absolute deviations (LAD) regression
      - Huber method: $\rho(z) = z^2$, if $|z|\leq c$, and $2c|z|-c^2$, if $|z|>c$. It's a compromise between OLS and LAD.
        [**Q**: how to pick $c$? suggestion: $c\in[1, 2]$ ]
      - many other choices, such as Tukey's biweight, Hampel, ...

    

    

    - compute M-estimates (related to iteratively re-weighted least square, IRWLS):
      - for LS with weights, estimate $\beta$ by solving $(X^TWX)\beta=X^TWY$
        $\Rightarrow X^TW(Y-X\beta)=0$, i.e., $\Sigma_i\, w_i\, x_{ij}\, (y_i-\Sigma_k\, x_{ik}\beta_k) = 0$    for all $j=1,...,p$
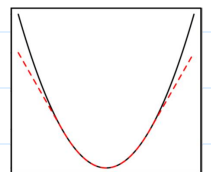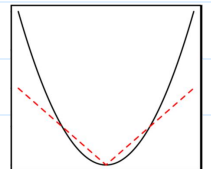      - for robust estimates, differentiating the M-estimate criterion w.r.t. $\beta_j$ and setting to zero, we get:
        $\Sigma_i\, \rho'((y_i-\Sigma_k\, x_{ik}\beta_k)/\sigma)\, x_{ij} = 0$ for all $j=1,...,p$
      - let $u_i = (y_i-\Sigma_k\, x_{ik}\beta_k)/\sigma$, we get $\Sigma_i\, (\rho'(u_i)/u_i)\, x_{ij}\, (y_i-\Sigma_k\, x_{ik}\beta_k) = 0$
        $\Rightarrow$ set $w_i = \rho'(u_i)/u_i$, can use WLS to estimate $\beta$
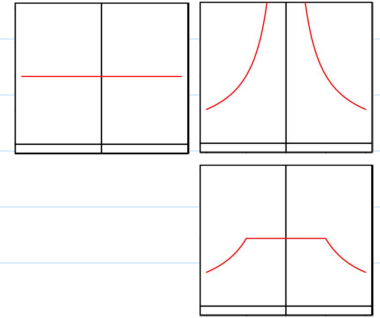        $\Rightarrow$ but, $u_i$ depends on the residuals $\hat{\varepsilon}_i \leftarrow \hat{\beta} \leftarrow w_i \leftarrow u_i \leftarrow \hat{\varepsilon}_i \leftarrow \ldots$
        $\Rightarrow$ IRWLS

- weights for various $\rho$ [note: larger residuals cause smaller weights in robust method]
  - OLS: $w(u)=2$ is a constant
  - LAD: $w(u)=1/|u|$ --- note the asymptote at 0, it may make a weighting approach difficult
  - Huber: $w(u)=2$, if $|z|\leq c$, and $2c/|u|$, if $|z|>c$.
- procedure: IRWLS for M-estimator
  - (1) start with any estimate of $\boldsymbol{\beta}$, say OLS
  - (2) compute residuals $\hat{\varepsilon}_i$
  - (3) compute $u_i$, may use median $|\hat{\varepsilon}_i - \text{median}(\hat{\varepsilon}_i)|/0.6745$ to estimate $\sigma$
  - (4) compute $w_i = \rho'(u_i)/u_i$
  - (5) do WLS to get a new estimate of $\boldsymbol{\beta}$, then go to step (2) until converge

➢ resistant regression (more resistant to outliers than M-estimators):

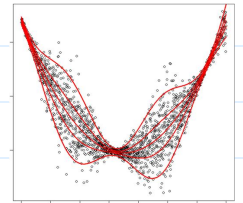- least trimmed squares (LTS): find $\boldsymbol{\beta}$ to minimize $\sum_{i=1}^{q} |y_i - \boldsymbol{x}_i^T\boldsymbol{\beta}|^2_{(i)}$, where $(i)$ indicates sorting, and $q<n$ [$q\approx(n+p+1)/2$ is recommended]
- least median of squares (LMS): find $\boldsymbol{\beta}$ to minimize median $|y_i - \boldsymbol{x}_i^T\boldsymbol{\beta}|^2$

- $\boldsymbol{\beta}$ estimated by S-estimation method. [see Rousseeuw and Leroy, 1987]
- resistant regression will do well even if a substantial proportion of data is "bad" (see an example in Lab)

➢ quantile regression

• **Q**: Why not always use robust estimates?

➢ if errors are (close to) normally distributed, robust estimators are less efficient

➢ very little distribution theory for robust estimator: can estimate $\boldsymbol{\beta}$ and (possibly) their standard errors, but, methodology and software for inference, such as testing, is not easy to come by. [$\Rightarrow$ may try bootstrap method]

➢ recommendation: use robust estimates as a check on OLS estimates. If they are close, use OLS theory. If not, try to find out why.

• Note: robust estimators provide protection against long-tailed errors, but they cannot overcome problems with non-constant variance or curvature in the mean of residuals.

❖ **Reading**: Faraway (1st ed.), 6.4,     ❖ **Further reading**: D&S, chapter 25

## Incomplete data

• Some values of some cases are missing. **Q**: When this happened, what can be done?

➢ find them --- may not be possible

➢ ask <u>why</u> the data are <u>missing</u>, i.e., *what is the <u>missing mechanism</u>?*

- <u>missing completely at random</u> (<u>MCAR</u>): <u>missing probability</u> is the <u>same</u> for <u>all cases</u> ⇒ <u>non-informative missing</u>

- <u>missing at random</u> (<u>MAR</u>): <u>missing probability</u> is <u>not constant</u>, but <u>depends on</u> a *<u>known</u>* mechanism, say some <u>observed variables</u> $T$ ⇒ <u>non-informative missing</u> if $T$ are included <u>in the model</u>

  
  Sampled cases

- <u>missing not at random</u> (<u>MNAR</u>): <u>missing probability</u> is <u>not constant</u>, and <u>depends on</u> some *<u>unknown</u>* mechanism ⇒ <u>informative missing</u>, e.g.:

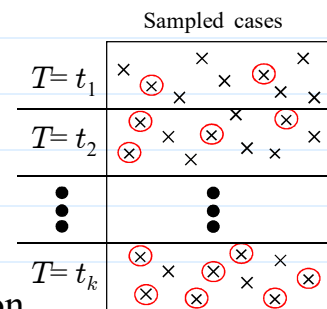  - <u>People</u> having something to <u>hide</u> are typically <u>less likely</u> to <u>provide information</u>

  - <u>Patients</u> <u>drop out</u> a drug study <u>more often</u> when they feel <u>treatment</u> is <u>not working</u>

  <u>MNAR</u> data require <u>special assumptions</u> and <u>modeling</u> [see <u>Little and Rubin, 2019</u>] ⇒ <u>Analyses</u> <u>without</u> considering the <u>information</u> in <u>missingness</u> may cause <u>biased conclusion</u>.

  
  Sampled cases
  $T = t_1$
  $T = t_2$
  $T = t_k$

- some <u>fix-up methods</u> for <u>non-informative missing</u>

  ➢ <u>approach 1</u>: <u>deletion</u>, i.e., ignore and <u>delete cases</u> with <u>missing value</u> ⇒ <u>no bias</u> but <u>lose information</u>. It is <u>OK</u> if % of <u>missing data</u> is <u>small</u>.

➢ <u>approach 2</u>: <u>single imputation</u> (<u>SI</u>), i.e., fill-in or <u>impute</u> a <u>missing value</u>, e.g.,

- <u>replace</u> missing value by <u>average of predictor</u>, often causing a <u>bias</u> of $\boldsymbol{\beta}$ toward $\boldsymbol{0}$.

- use a <u>regression</u> model to <u>predict</u> $x_i$ using <u>other predictors</u>,

  ⇒ how much <u>trouble</u> to take in <u>building these models</u>?

  ⇒ may be <u>difficult</u> with <u>multiple missing values</u>

  ⇒ cause some <u>bias</u>, but <u>filled-in case</u> will have <u>lower leverage</u>

  ⇒ **Q**: Is <u>inference</u> <u>valid</u> after <u>estimating</u> the <u>coefficients</u>?

- A <u>SI</u> value tends to be <u>less variable</u> than the <u>missing value</u> because the <u>imputed value</u> does <u>not include</u> the <u>error variation</u>.

➢ <u>approach 3</u>: <u>multiple imputation</u> (<u>MI</u>), i.e., <u>impute</u> a <u>missing value</u> $m$ times by multiple <u>draws</u> from <u>predictive distribution</u>

- <u>MI</u> re-includes <u>error variation</u>, which <u>reflects uncertainty</u> about <u>imputed values</u> and yields <u>valid estimates</u> of <u>standard errors</u>.

- <u>MI</u> may better <u>mitigate bias</u>

- Let $\hat{\beta}_{ij}$ and $s_{ij}$ be the <u>estimate</u> and <u>standard error</u> of the coefficient $\beta_i$ of $x_i$ for the *j*th <u>imputed</u> result, $j = 1, \ldots, m$.

  - The <u>combined</u> estimate of $\beta_i$ is: $\widehat{\beta}_i = \frac{1}{m} \sum_{j=1}^{m} \widehat{\beta}_{ij}$

□ The combined standard errors $s_i$ of $\widehat{\beta}_i$ is given by:

$$s_i^2 = \frac{1}{m} \sum_{j=1}^{m} s_{ij}^2 + \left(1 + \frac{1}{m}\right) \mathsf{var}(\widehat{\beta}_i),$$

where $\mathsf{var}(\widehat{\beta}_i)$ is the (unbiased) sample variance over the imputed $\hat{\beta}_{ij}$'s .

➢ approach 4: maximum likelihood method

Assuming complete data $\underline{D} = (\underline{D}_{\mathrm{obs}}, \underline{D}_{\mathrm{mis}})$, both observed and missing, are from a family of distribution with parameters $\boldsymbol{\theta}$, say multivariate normal, then it is possible to compute maximum likelihood estimates using:

▪ (if available) the likelihood of $\boldsymbol{\theta}$ based on $\underline{D}_{\mathrm{obs}}$ :

$$\mathcal{L}(\boldsymbol{\theta}|\underline{D}_{\mathsf{obs}}) = \int \underline{f}_{\underline{D}}(D_{\mathsf{obs}}, \underline{D}_{\mathsf{mis}} \,|\, \boldsymbol{\theta}) \, d\underline{D}_{\mathsf{mis}}$$

▪ the EM algorithm

But,

▪ the distribution assumption might not be tenable

▪ tests, inferences, and diagnostics are not easy to come by

❖ **Reading**: Faraway (1st ed.), chapter 12; W, 5.6