

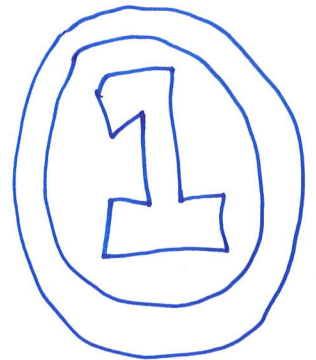
# ROBUST REGRESSION

(when errors are non-normally heavy-tailed, extreme)

## M-ESTIMATION

- Chooses the  $\hat{\beta}$  value to minimize

$$\sum_{i=1}^n \rho(y_i - x_i^T \hat{\beta}) \quad (*)$$



where  $\rho(\cdot)$  is a function. Common choices:

①  $\rho(x) = x^2$  (least squares regression)

②  $\rho(x) = |x|$  least absolute deviation regression (LAD)  
or  $L_1$  regression ( $|x| = \sum_{i=1}^n |x_i|$  is the  $L_1$ -norm)

③ HUBER'S METHOD:  
$$\rho(x) = \begin{cases} x^2/2 & \text{for } |x| \leq c \leftarrow \text{some threshold;} \\ c|x| - c^2/2 & \text{otherwise} \end{cases}$$

Here,  $c$  is a robust estimator of  $\sigma$ , such as the median of the  $|\hat{\epsilon}_i|$  values.

## RELATION TO WEIGHTED LEAST SQUARES

Normal Equations for LEAST SQUARES (what you get when differentiating w.r.t.  $\beta$  and setting to 0)

$$X^T X \hat{\beta} = X^T Y$$

Normal Equations for WEIGHTED LEAST SQUARES:

$$X^T \Sigma^{-1} X \hat{\beta} = X^T \Sigma^{-1} Y, \text{ or in summations: } \sum_{i=1}^n w_i x_{ij} y_i = \sum_{i=1}^n w_i x_{ij} \sum_{j=1}^p w_i x_{ij} \hat{\beta}_j$$

... or with summations ...

~~$$\sum_{i=1}^n w_i x_{ij} \sum_{k=1}^p x_{ik} \hat{\beta}_k = \sum_{i=1}^n w_i x_{ij} y_i, \quad j=1, 2, \dots, p$$~~

(2)

$$\sum_{i=1}^n w_i x_{ij} \sum_{k=1}^p x_{ik} \hat{\beta}_k = \sum_{i=1}^n w_i x_{ij} y_i, \quad j \in \{1, 2, \dots, p\} \quad (**)$$

Differentiating the M-estimate criterion (\*) w.r.t.  $\hat{\beta}_j$  and setting to 0:

$$\frac{d}{d\hat{\beta}_j} \sum_{i=1}^n \rho(y_i - x_i^T \hat{\beta}) = \sum_{i=1}^n \rho'(y_i - x_i^T \hat{\beta}) \cdot (-x_{ij}) = 0$$

(in summation form)  $\Rightarrow$

$$\sum_{i=1}^n \rho'(y_i - \underbrace{\sum_{k=1}^p x_{ik} \hat{\beta}_k}_{i^{th} \text{ residual} = u_i}) x_{ij} = 0, \quad j \in \{0, 1, \dots, n\}$$

Set Or

$$\sum_{i=1}^n \underbrace{\frac{\rho'(u_i)}{u_i}}_{w(u_i)} x_{ij} \left( y_i - \sum_{k=1}^p x_{ik} \hat{\beta}_k \right) = 0, \quad j \in \{1, 2, \dots, p\}$$

or 
$$\sum_{i=1}^n \underbrace{w(u_i)}_{\text{weight function}} x_{ij} \left( y_i - \sum_{k=1}^p x_{ik} \hat{\beta}_k \right) = 0, \quad j \in \{1, 2, \dots, p\}$$

SAME as (\*\*)  $\Rightarrow$

1.  $w(u) = \text{constant} \Rightarrow$  ordinary least squares.

2.  $w(u) = \frac{1}{|u|} \Rightarrow$  L.A.D. Note the weight decreases as the residual increases: more extreme observations have smaller weight.

3. 
$$w = \begin{cases} 1 & \text{for } |u| \leq c \\ c/|u| & \text{otherwise} \end{cases}$$

3

is the HUBER method, which is a sort-of compromise between least-squares and L.A.D.

\* There are many other potentially good choices for  $\rho(\cdot)$ .

\* M-estimation can require significant computing time because it is carried out iteratively: the weights depend on the residuals: get weights, get  $\hat{\beta}$ 's, get residuals, get the weights, get  $\hat{\beta}$ 's, get residuals, get....

Get standard errors using weighted-least-squares with

$$\widehat{\text{Var}}(\hat{\beta}) = \hat{\sigma}^2 (X^T W X)^{-1}, \text{ w/ a robust estimate of } \sigma^2.$$