

Statistics 203: Introduction to Regression and Analysis of Variance

Robust methods

Jonathan Taylor



Today's class

● Today's class

- Heteroskedasticity
- MLE for one sample problem
- Weighted least squares
- Estimating σ^2
- Weighted regression example
- Robust methods
- Example
- M -estimators
- Huber's ψ
- Hampel's ψ
- Tukey's ψ
- Solving for $\hat{\beta}$
- Iteratively reweighted least squares (IRLS)
- Robust estimate of scale
- Other resistant fitting methods
- Why not always use robust regression?

- Weighted regression.
- Robust methods.
- Robust regression.



Heteroskedasticity

- Today's class
- Heteroskedasticity
- MLE for one sample problem
- Weighted least squares
- Estimating σ^2
- Weighted regression example
- Robust methods
- Example
- M -estimators
- Huber's ψ
- Hampel's ψ
- Tukey's ψ
- Solving for $\hat{\beta}$
- Iteratively reweighted least squares (IRLS)
- Robust estimate of scale
- Other resistant fitting methods
- Why not always use robust regression?

- In our standard model, we have assumed that

$$\varepsilon \sim N(0, \sigma^2 I).$$

That is, that the errors are independent and have the same variance (*homoskedastic*).

- We have discussed graphical checks for non-constant variance (*heteroskedasticity*) but not “remedies” for heteroskedasticity.

- Suppose that

$$\varepsilon \sim N(0, \sigma^2 D)$$

for some known diagonal matrix D .

- Where does D come from? Suppose that we see that variance increases like $f(X_j)$, then we might choose $D_i = f(X_{ij})$.
- What is the “maximum likelihood” thing to do?



MLE for one sample problem

- Today's class
- Heteroskedasticity
- MLE for one sample problem
- Weighted least squares
- Estimating σ^2
- Weighted regression example
- Robust methods
- Example
- M -estimators
- Huber's ψ
- Hampel's ψ
- Tukey's ψ
- Solving for $\hat{\beta}$
- Iteratively reweighted least squares (IRLS)
- Robust estimate of scale
- Other resistant fitting methods
- Why not always use robust regression?

■ Consider the simpler problem

$$Y_i \sim N(\mu, \sigma^2 D_i)$$

with σ^2 and D 's known.



$$-2 \log L(\mu|Y, \sigma) = \sum_{i=1}^n \frac{(Y_i - \mu)^2}{\sigma^2 D_i} + n \log(2\pi\sigma^2) + \sum_{i=1}^n \log(D_i)$$

■ Differentiating

$$-2 \sum_{i=1}^n \frac{(Y_i - \hat{\mu})}{\sigma^2 D_i} = 0$$

implying

$$\hat{\mu} = \sum_{i=1}^n \frac{Y_i}{D_i} / \sum_{i=1}^n \frac{1}{D_i}.$$

■ Observations are weighted inversely proportional to their variance.



Weighted least squares

- Today's class
- Heteroskedasticity
- MLE for one sample problem
- Weighted least squares
- Estimating σ^2
- Weighted regression example
- Robust methods
- Example
- M -estimators
- Huber's ψ
- Hampel's ψ
- Tukey's ψ
- Solving for $\hat{\beta}$
- Iteratively reweighted least squares (IRLS)
- Robust estimate of scale
- Other resistant fitting methods
- Why not always use robust regression?



$$\begin{aligned} -2 \log L(\beta, \sigma | Y, D) &= \sum_{i=1}^n \frac{(Y_i - \beta_0 - \sum_{j=1}^{p-1} \beta_j X_{ij})^2}{\sigma^2 D_i} \\ &= \frac{1}{\sigma^2} (Y - X\beta)^t D^{-1} (Y - X\beta) \\ &= \frac{1}{\sigma^2} (Y - X\beta)^t W (Y - X\beta). \end{aligned}$$

with $W = D^{-1}$.

■ Normal equations:

$$-2X^t W (Y - X\hat{\beta}_W) = 0$$

or,

$$\hat{\beta}_W = (X^t W X)^{-1} X^t W Y.$$

■ Distribution of $\hat{\beta}_W$

$$\hat{\beta}_W \sim N(\beta, \sigma^2 (X^t W X)^{-1}).$$



Estimating σ^2

- Today's class
- Heteroskedasticity
- MLE for one sample problem
- Weighted least squares
- Estimating σ^2
- Weighted regression example
- Robust methods
- Example
- M -estimators
- Huber's ψ
- Hampel's ψ
- Tukey's ψ
- Solving for $\hat{\beta}$
- Iteratively reweighted least squares (IRLS)
- Robust estimate of scale
- Other resistant fitting methods
- Why not always use robust regression?

■ What are the right residuals?

■ If we knew β exactly

$$Y_i - \beta_0 - \sum_{j=1}^{p-1} X_{ij}\beta_j \sim N(0, \sigma^2 / W_i).$$

■ Suggests that the natural residual is

$$e_{W,i} = \sqrt{W_i} \frac{Y_i - \hat{Y}_{W,i}}{=} \sqrt{W_i} e_i$$

where

$$\hat{Y}_W = (X^t W X)^{-1} X^t W Y.$$

■ Estimate of σ^2

$$\hat{\sigma}_W^2 = \frac{1}{n-p} \sum_{i=1} e_{W,i}^2 = \frac{1}{n-p} \sum_{i=1} w_i e_i^2 \sim \sigma^2 \frac{\chi_{n-p}^2}{n-p}$$



Weighted regression example

- Today's class
- Heteroskedasticity
- MLE for one sample problem
- Weighted least squares
- Estimating σ^2
- Weighted regression example
- Robust methods
- Example
- M -estimators
- Huber's ψ
- Hampel's ψ
- Tukey's ψ
- Solving for $\hat{\beta}$
- Iteratively reweighted least squares (IRLS)
- Robust estimate of scale
- Other resistant fitting methods
- Why not always use robust regression?



Robust methods

- Today's class
- Heteroskedasticity
- MLE for one sample problem
- Weighted least squares
- Estimating σ^2
- Weighted regression example
- Robust methods
- Example
- M -estimators
- Huber's ψ
- Hampel's ψ
- Tukey's ψ
- Solving for $\hat{\beta}$
- Iteratively reweighted least squares (IRLS)
- Robust estimate of scale
- Other resistant fitting methods
- Why not always use robust regression?

- We also discussed outlier detection but no specific remedies.
- One alternative is to discard potential outliers – not always a good idea.
- Outliers can really mess up the sample mean, but have relatively little effect on the sample median.
- Could also “downweight” outliers: basis of robust techniques.
- Another “cause” of outliers may be that the data is not really normally distributed.



Example

- Today's class
- Heteroskedasticity
- MLE for one sample problem
- Weighted least squares
- Estimating σ^2
- Weighted regression example
- Robust methods
- **Example**
- M -estimators
- Huber's ψ
- Hampel's ψ
- Tukey's ψ
- Solving for $\hat{\beta}$
- Iteratively reweighted least squares (IRLS)
- Robust estimate of scale
- Other resistant fitting methods
- Why not always use robust regression?

- Suppose that we have a sample $(Y_i)_{1 \leq i \leq n}$ from

$$f(y|\mu, \sigma) = \frac{1}{2\sigma} e^{-|y-\mu|/\sigma}.$$

This has heavier tails than the normal distribution.

- MLE for μ

$$\hat{\mu} = \operatorname{argmin}_{\mu} \sum_{i=1}^n |Y_i - \mu|.$$

- It can be shown that $\hat{\mu}$ is the sample median (exercise in STATS116).
- Take home message: if errors are not really normally distributed then least squares is not MLE and the MLE downweights large residuals relative to least squares.



M-estimators

- Today's class
- Heteroskedasticity
- MLE for one sample problem
- Weighted least squares
- Estimating σ^2
- Weighted regression example
- Robust methods
- Example
- **M-estimators**
- Huber's ψ
- Hampel's ψ
- Tukey's ψ
- Solving for $\hat{\beta}$
- Iteratively reweighted least squares (IRLS)
- Robust estimate of scale
- Other resistant fitting methods
- Why not always use robust regression?

- Depending on the error distribution of ε (assuming i.i.d.) we get different optimization problems.

- Suppose that

$$f(y|\mu, s) \propto e^{-\rho((y-\mu)/s)}.$$

- MLE:

$$\hat{\mu} = \operatorname{argmin}_{\mu} \sum_{i=1}^n \rho \left(\frac{Y_i - \mu}{s} \right).$$

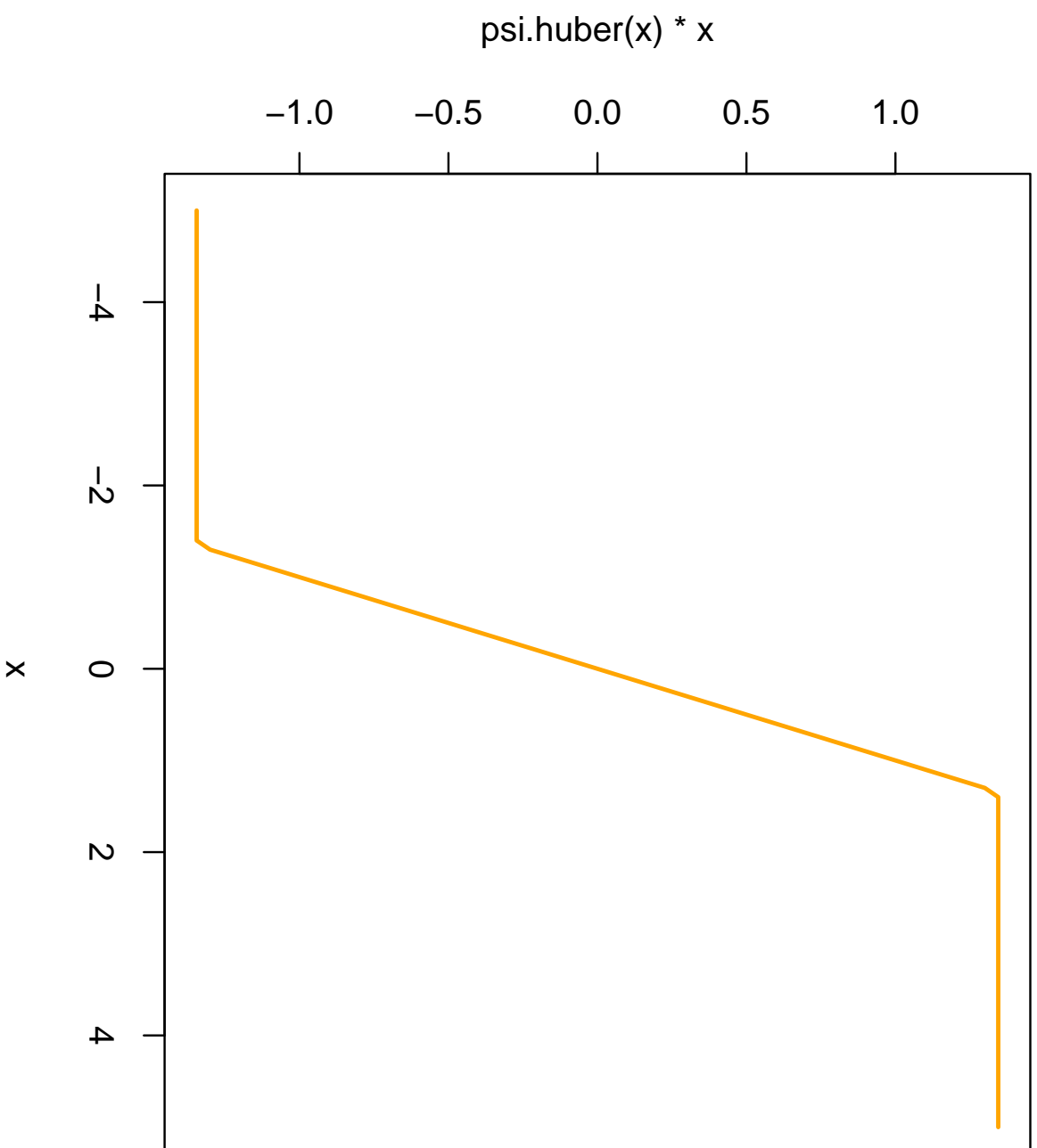
- For every ρ , we get a different estimator: an M -estimator. Let $\psi = \rho'$.

- Generalizes easily to regression problem

$$\hat{\beta} = \operatorname{argmin}_{\beta} \sum_{i=1}^n \rho \left(\frac{Y_i - \beta_0 - \sum_{j=1}^{p-1} X_{i,j} \beta_j}{s} \right).$$

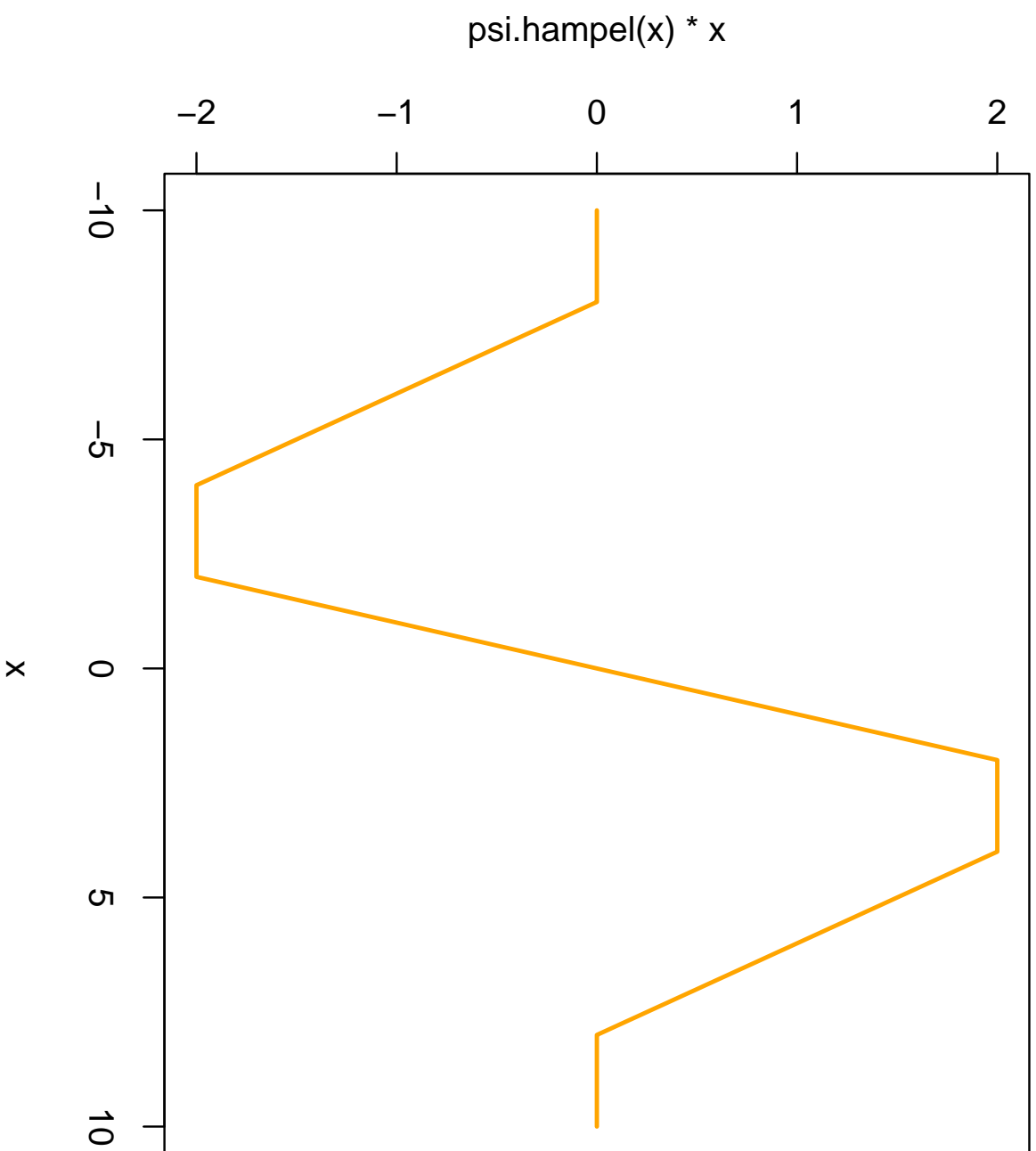


Huber's ψ



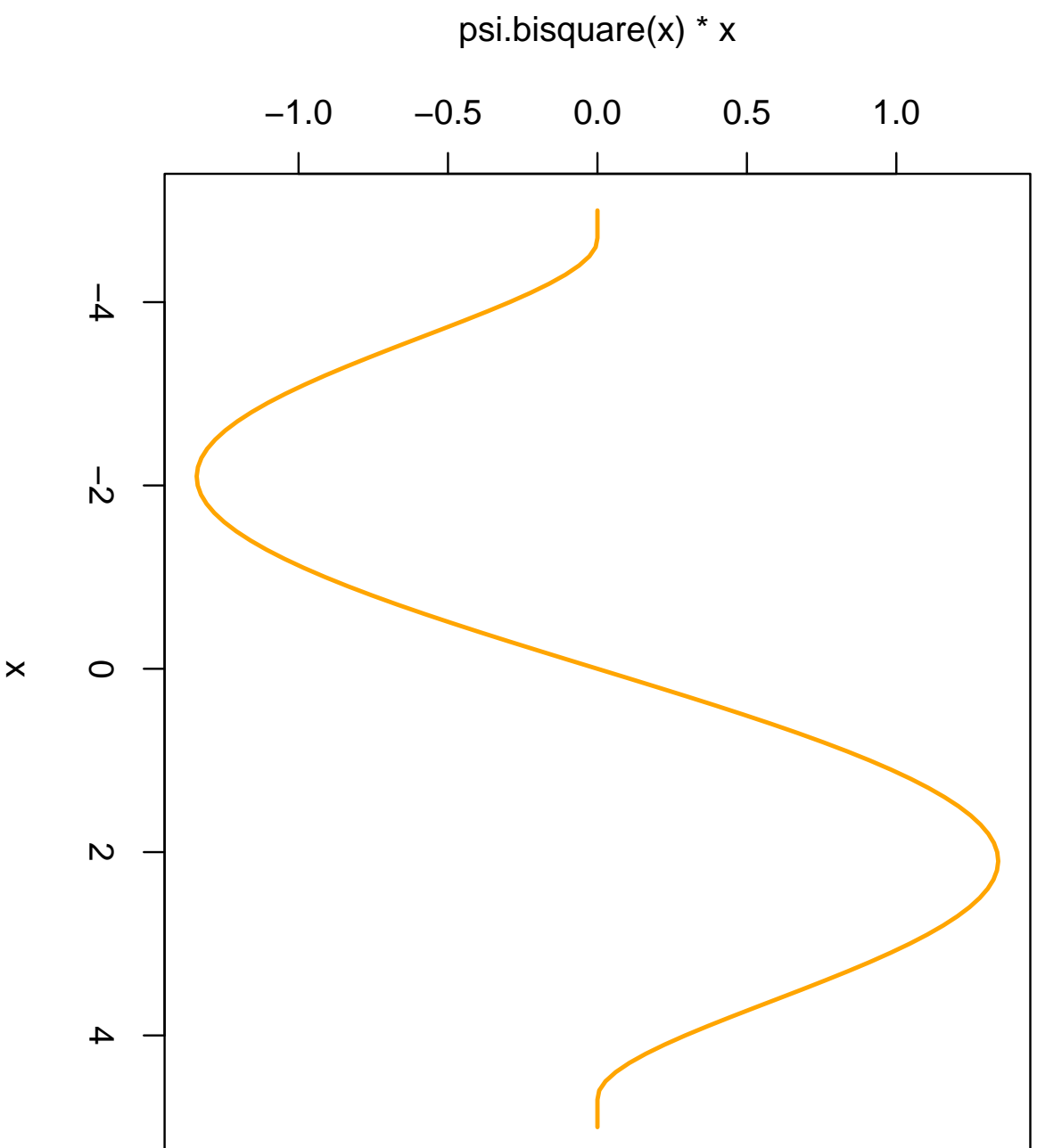


Hampel's ψ





Tukey's ψ





Solving for $\hat{\beta}$

- Today's class
- Heteroskedasticity
- MLE for one sample problem
- Weighted least squares
- Estimating σ^2
- Weighted regression example
- Robust methods
- Example
- M -estimators
- Huber's ψ
- Hampel's ψ
- Tukey's ψ
- Solving for $\hat{\beta}$
- Iteratively reweighted least squares (IRLS)
- Robust estimate of scale
- Other resistant fitting methods
- Why not always use robust regression?

- Assume for now that s is known in the M -estimator.

- Normal equations:

$$\sum_{i=1}^n X_{ij} \psi \left(\frac{Y_i - \sum_{j=0}^{p-1} X_{ij} \hat{\beta}_j}{s} \right) = 0, \quad 0 \leq j \leq p-1$$

where $\psi = \rho'$.

- Set

$$W_i = \frac{\psi \left(\frac{Y_i - \sum_{j=0}^{p-1} X_{ij} \hat{\beta}_j}{s} \right)}{Y_i - \sum_{j=0}^{p-1} X_{ij} \hat{\beta}_j}.$$

- Then

$$\sum_{i=1}^n X_{ij} W_i (Y_i - \sum_{j=0}^{p-1} X_{ij} \hat{\beta}_j) = 0, \quad 0 \leq j \leq p-1.$$

Or

$$\hat{\beta} X^t W X - X^t W Y = 0.$$



Iteratively reweighted least squares (IRLS)

- Today's class
- Heteroskedasticity
- MLE for one sample problem
- Weighted least squares
- Estimating σ^2
- Weighted regression example
- Robust methods
- Example
- M -estimators
- Huber's ψ
- Hampel's ψ
- Tukey's ψ
- Solving for $\hat{\beta}$
- Iteratively reweighted least squares (IRLS)
- Robust estimate of scale
- Other resistant fitting methods
- Why not always use robust regression?

- Although we are trying to estimate β above, given an initial estimate $\hat{\beta}^0$ we can compute initial weights

$$W_i^0 = \psi \left(\frac{Y_i - \sum_{j=0}^{p-1} X_{ij} \hat{\beta}_j^0}{s} \right) / \left(Y_i - \sum_{j=0}^{p-1} X_{ij} \hat{\beta}_j^0 \right)$$

ans solve for $\hat{\beta}^1$

$$\hat{\beta}^1 = (X^t W^0 X)^{-1} X^t W^0 Y.$$

- Now we can recompute weights, and reestimate $\hat{\beta}$.
- In general, given weights W^j we can solve

$$\hat{\beta}^{j+1} = (X^t W^j X)^{-1} X^t W^j Y.$$

- This is very similar to a Newton-Raphson technique.
- Used over and over again in statistics.



Robust estimate of scale

- Today's class
- Heteroskedasticity
- MLE for one sample problem
- Weighted least squares
- Estimating σ^2
- Weighted regression example
- Robust methods
- Example
- M -estimators
- Huber's ψ
- Hampel's ψ
- Tukey's ψ
- Solving for $\hat{\beta}$
- Iteratively reweighted least squares (IRLS)
- Robust estimate of scale
- Other resistant fitting methods
- Why not always use robust regression?

- In general s needs to be estimated. A popular estimate

$$s = MAD(e_1, \dots, e_n)/0.6745.$$

where

$$MAD(e_1, \dots, e_n) = \text{Median}|e_i - \text{Median}(e_1, \dots, e_n)|.$$

- Scale can be estimated at each stage of IRLS procedure based on current residuals, or based on some “very resistant” fit. (LMS or LTS below)
- The constant 0.6745 is chosen so that s is asymptotically unbiased for σ if the e_i 's are $N(0, \sigma^2)$.



Other resistant fitting methods

- Today's class
- Heteroskedasticity
- MLE for one sample problem
- Weighted least squares
- Estimating σ^2
- Weighted regression example
- Robust methods
- Example
- M -estimators
- Huber's ψ
- Hampel's ψ
- Tukey's ψ
- Solving for $\hat{\beta}$
- Iteratively reweighted least squares (IRLS)
- Robust estimate of scale
- Other resistant fitting methods
- Why not always use robust regression?

- If ψ is not redescending (i.e. does not return to 0) then robust regression is still susceptible to outliers.
- Least median of squares (LMS)

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \operatorname{Median}(Y_i - \sum_{j=0}^{p-1} X_{ij}\beta_j)^2$$

- Least trimmed squares (LTS): fix some $q < n$

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^q (Y_i - \sum_{j=0}^{p-1} X_{ij}\beta_j)_{(i)}^2.$$

Here the $\cdot_{(i)}$ represents “order statistic”.



Why not always use robust regression?

- Today's class
- Heteroskedasticity
- MLE for one sample problem
- Weighted least squares
- Estimating σ^2
- Weighted regression example
- Robust methods
- Example
- M -estimators
- Huber's ψ
- Hampel's ψ
- Tukey's ψ
- Solving for $\hat{\beta}$
- Iteratively reweighted least squares (IRLS)
- Robust estimate of scale
- Other resistant fitting methods
- Why not always use robust regression?

- Inference: seems reasonable to treat estimate covariance of $\hat{\beta}$ as

$$\hat{\sigma}^2 (X^t \hat{W} X)^{-1}, \quad \hat{\sigma}^2 = \frac{\sum_{i=1}^n \hat{W}_i e_i^2}{n - p}.$$

What about degrees of freedom?

- R uses a different estimate.
- Efficiency: this robustness comes at a cost. Even asymptotically, confidence intervals for robust estimates are wider than least squares (if least squares model is applicable).
- *Asymptotic Relative Efficiency:*

$$ARE(\hat{\beta}_{LS}, \hat{\beta}_{robust}) = \lim_{n \rightarrow \infty} \frac{\text{Var}(\hat{\beta}_{LS})}{\text{Var}(\hat{\beta}_{robust})} < 1.$$

- Biggest advantage: can have higher breakdown. Median has 50% breakdown, sample mean has 0%.