

STA 302 / 1001

Note Title

10/18/2011

See Blackboard announcements for test information

Continuing from last class:

Properties of residuals : \hat{e}_i

$$E(\hat{e}_i) = 0$$

$$\text{Var}(\hat{e}_i) = \sigma^2(1 - h_{ii}) \quad \text{not constant}$$

$$\text{Cov}(\hat{e}_i, \hat{e}_j) = -h_{ij}\sigma^2, \quad i \neq j$$

(not uncorrelated)

$\hat{\epsilon}_i$'s are "more" normally distributed than ϵ_i 's (if ϵ_i 's distribution deviates from normal)
by CLT

If standardizing residuals,

$$r_i = \frac{\hat{\epsilon}_i}{s \sqrt{1 - h_{ii}}}$$

Two advantages

- r_i 's no longer have unequal variances
- easier to find outliers

that are an
indicator of lack
of model fit
 $|r_i| > 3$ unusual
 $|r_i| > 4$ very unusual

(Small data sets $|r_i| > 2$?)

Watch for $|r_i|$ that are
much larger than all of the
others.

Typically, don't see much ^{qualitative} difference between r_i and
 \hat{e}_i , except for leverage points

$$h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{S_{xx}} \rightarrow 0 \text{ for large } n$$

as long as x_i is not too far from \bar{x}

$$h_{ij} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{S_{xx}} \rightarrow 0 \text{ for large } n$$

as long as x_i, x_j not too far from \bar{x}

So correlation and lack of constant variance in the residuals is often ignored.

Residual Plots for Checking Model Assumptions

Looking for:

- evidence that straight line model is not appropriate
 - curvature
 - outliers
 - influential points

is for some

or all of data

- constant variance
- normally distributed errors

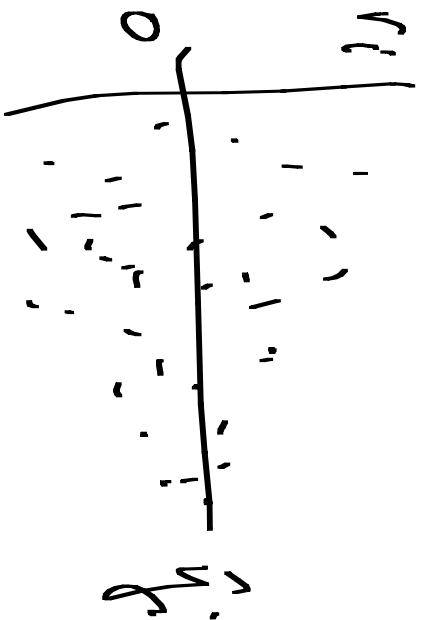
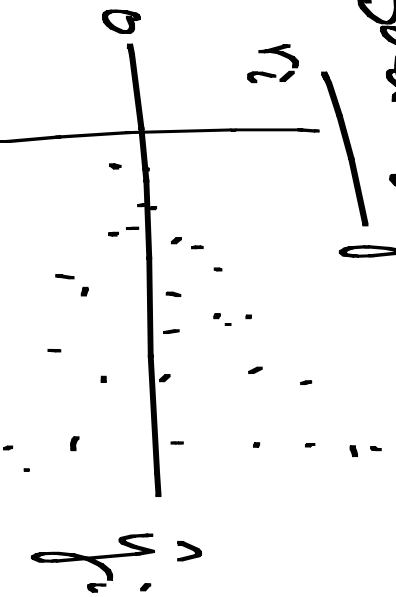
Also (if can)

- lack of independence in errors
- potential missing predictor variables

Recommended Plots

① r_i vs \hat{y}_i (Standardized residuals versus fitted values)

Good sign!



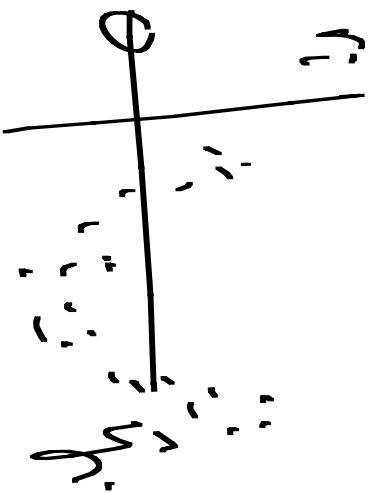
evidence of non-constant variance

$$\text{Kann } \sum \hat{e}_i \hat{y}_i = 0$$

(Is $\sum r_i \hat{y}_i = 0$?)

\Rightarrow correlation coefficient (r) of \hat{e}_i, \hat{y}_i is 0
So shouldn't be a pattern in \hat{e}_i, \hat{y}_i plot

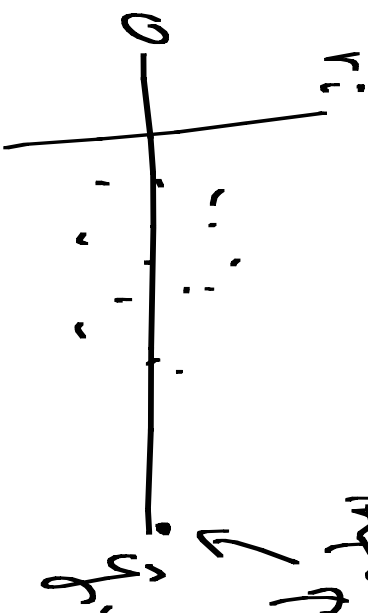
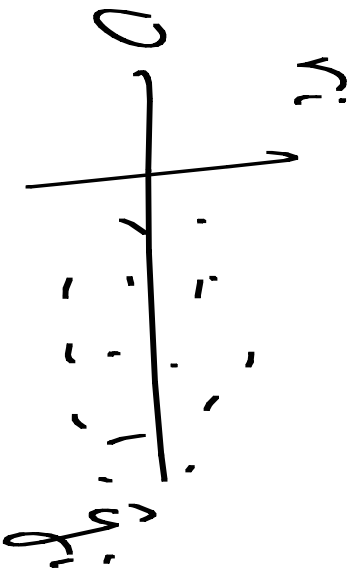
Bad sign



vidence of overfit

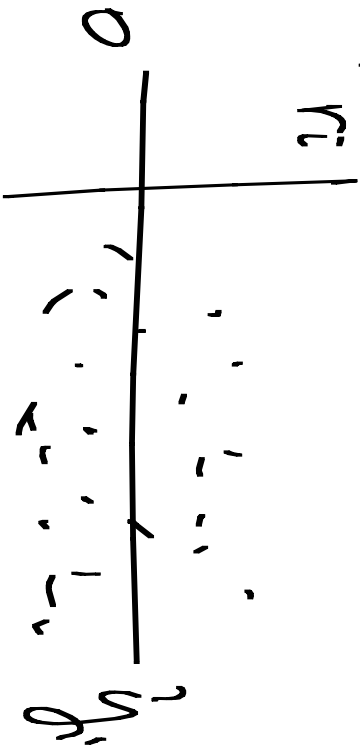
Good sign

← outlier



leverage
~~Adjusted~~
point

'Good' plot - random scatter about 0



① r_i vs x_i Residuals versus explanatory variable.

- for simple linear regression, looks exactly like plot of r_i vs \hat{y}_i ($\hat{y}_i = b_0 + b_1 x_i$) except scale is different and plot is flipped $b_1 < 0$
- look for same things as r_i vs \hat{y}_i plot

Typically, use r_i vs x_i for curvature, r_i vs \hat{y}_i for non-constant, r_i vs \hat{y}_i for unusual observations

- ③ Normal quantile plot (come back to this)
- assess normality
 - don't bother to assess normality in presence of other problems

Other residual analysis

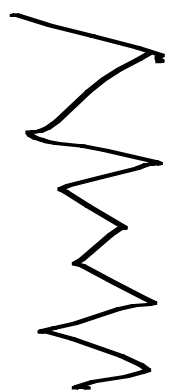
- univariate analysis of residuals
 - histogram should be bell-shaped
 - should satisfy 68/95/99.7% rule
 - look for outliers, skew

- Plot \sum Square root of absolute value of residuals vs \sum

- increasing pattern \Rightarrow non-constant variance

- Residuals versus time or other sequence in observation
- a pattern indicates residuals may be correlated


positive
auto-correlation


negative auto-correlation

- Residuals versus other potential predictors
 - Most random scatter
 - Any pattern indicates other predictor should be in the model

Assignment 1 marking

- Marks will be adjusted so each TA will have same distribution
- 3 may not be perfect
- look at solutions
- look at posted TA comments
- Re-mark policy look at solutions; write down concerns; give to me or Andrew