

Regression Diagnostics: Ch 9

15.1

- To detect problems with a model, like
- outliers
 - Curvilinear trends
 - Non-constant variance
- Not too scary as non-linear

Scatterplots are great if there is only one independent variable, but not practical for 3D + 4

Can look at many plots including

- \hat{Y} vs \hat{Y}
- \hat{Y} vs $\hat{\epsilon}$
- X vs $\hat{\epsilon}$

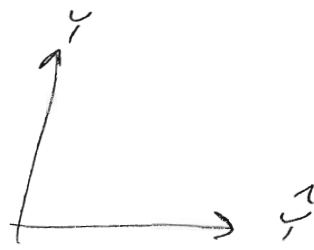
- Time sequence vs $\hat{\epsilon}$

Will need

$$r_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

Look up

\hat{Y} vs Y



15.2

What should we see?

Squared sample correlation between \hat{Y} & Y equals $R^2 = \frac{SSR}{SST}$

Proof: Recall in decomposition of SS , that

$$SST = SSE + 0 + SSR$$
$$\sum (Y_i - \bar{Y})^2 = \sum (Y_i - \hat{Y}_i)^2 + \sum (\hat{Y}_i - \bar{Y})^2$$

$$2 \sum (Y_i - \hat{Y}_i)(\hat{Y}_i - \bar{Y})$$

And, since $\sum \hat{e}_i = \sum (Y_i - \hat{Y}_i) = 0$,

$$\hat{Y} = \bar{Y}$$

, so sample correlation between \hat{Y} & Y is

15.3

$$R^2 = \left(\frac{\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{y})}{\sqrt{\sum (y_i - \bar{y})^2} \sqrt{\sum (\hat{y}_i - \bar{y})^2}} \right)^2$$

$$= \left(\frac{\sum (y_i - \hat{y}_i + \hat{y}_i - \bar{y})(\hat{y}_i - \bar{y})}{\sqrt{SST \times SSR}} \right)^2$$

$$= \frac{\left(\sum (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) + \sum (\hat{y}_i - \bar{y})^2 \right)^2}{SST \times SSR}$$

$$= \frac{SSR \times SSR}{SST \times SSR} = R^2$$

So the plot of \hat{y}_i vs y_i is the closest thing we have to a single scatter plot

Could reveal

outliers

curvilinear trends,

non-constant variance

What if it looks like there are outliers?

15.4

When is a residual "too big" in absolute value?

→ One problem is that while ϵ_i all have same variance, $\hat{\epsilon}_i$ do NOT.

$$\text{Recall } \hat{\epsilon} = (y - \hat{y}) = (I - H)y$$

Where $H = X(X'X)^{-1}X'$ is the "Hat" matrix, so called because it puts a hat on y , by

$$\hat{y} = X\hat{\beta} = X(X'X)^{-1}X'y = Hy$$

$$\text{So } \text{cov}(\hat{\epsilon}) = (I - H)\sigma^2 I_n (I - H)' \\ = \sigma^2 (I - H)$$

Writing $H = [h_{ij}]$, $\text{Var}(\hat{\epsilon}_i) = \sigma^2(1 - h_{ii})$

Could standardize it

$$r_i = \frac{\hat{\epsilon}_i}{\sqrt{\sigma^2(1 - h_{ii})}} \approx \frac{\epsilon_i}{\sqrt{\sigma^2(1 - h_{ii})}} = \frac{\epsilon_i}{\sigma\sqrt{1 - h_{ii}}}$$

The "Studentized" Residual

Studentized residual has (15.5)
 (for large samples) an approximate standard normal dist,
 but not t , because $NUM \neq DEN$ are not quite χ^2 independent.

Better is the STUDENTIZED DELETED RESIDUAL.

Do a regression omitting observation i , obtaining $\hat{\beta}_{(i)}$. Calculate (see 9.3', p. 234)

Let $\hat{E}_{(i)} = y_i - x_i' \hat{\beta}_{(i)}$
 Notice parallel to prediction interval

$$\hat{E}_{(i)} \sim N(0, \sigma^2 + \sigma^2 x_i' (X_{(i)}' X_{(i)})^{-1} x_i)$$

$$Z = \frac{\hat{E}_{(i)}}{\sqrt{\sigma^2 (1 + x_i' (X_{(i)}' X_{(i)})^{-1} x_i)}}$$

$$\text{And } \lambda_i = \frac{Z_i}{\sqrt{W / (n-2)}} = \frac{Z_i}{\Delta_{(i)}}, \Delta_{(i)}^2 = \frac{SSE_{(i)}}{n-2}$$

$$W = \frac{SSE_{(i)}}{\sigma^2} \sim \chi^2(n-2)$$

ind of $y_i \neq \hat{\beta}_{(i)}$

$$\begin{aligned}
 t_i &= \frac{\hat{\varepsilon}(i)}{\sqrt{\sigma^2(1 + x_i'(X_{(i)}' X_{(i)})^{-1} x_i)}} \sim A(n-k-2) \\
 &= \frac{\hat{\varepsilon}(i)}{\sigma(i) \sqrt{1 + x_i'(X_{(i)}' X_{(i)})^{-1} x_i}} \sim A(n-k-2)
 \end{aligned}$$

If true model is correct.



This is a good definition of "too big"

We are conducting n tests

(not independent)

If model is correct,

- Pr rejecting each H_0 is $0.05 = \alpha$
- Pr of rejecting at least one is much greater
- How much greater?

How much greater?

15.7

Suppose $\alpha = 0.05$, all k null hypotheses are true, and tests are independent (which they are not, here). Then letting A_j represent event reject $H_0 j$,

$$P(\text{Reject at least one}) = P\left(\bigcup_{j=1}^k A_j\right)$$

$$= 1 - P\left\{\left(\bigcup_{j=1}^k A_j\right)^c\right\}$$

De Morgan

$$\stackrel{\text{De Morgan}}{=} 1 - P\left\{\bigcap_{j=1}^k A_j^c\right\} \stackrel{\text{ind}}{=} 1 - \prod_{j=1}^k P(A_j^c)$$

$$= 1 - \prod_{j=1}^k (1 - P(A_j)) = 1 - (1 - \alpha)^k$$

So for $\alpha = 0.05$ \neq 50 tests

(remember, n tests here), JOINT

significance level of the 50 tests is

$$1 - (0.95)^{50} = 0.923$$

But tests are not independent

Bonferroni's Inequality

15.8

$$P\left(\bigcup_{j=1}^k A_j\right) \leq \sum_{j=1}^k P(A_j)$$

- A_j is event that Test j rejects H_0
- Applies to any collection of tests
- Do two tests as usual, but with significance level α/k . Then

$$P\left(\bigcup_{j=1}^k A_j\right) \leq \sum_{j=1}^k \alpha/k = k \alpha/k = \alpha$$

A bit conservative, but flexible & SIMPLE

Use the critical values for α/k instead of α

For using Studentized deleted residuals, $k = n$. Use two α/n critical values.

Plot of \hat{Y} vs $\hat{\epsilon}$.

if the model is correct

15.9

What should you see

Nothing, because they are independent.

USE • $E(\hat{\beta}\hat{\beta}') = \sigma^2(X'X)^{-1} + \beta\beta'$

• $E(Y'Y) = \text{cov}(Y) + (EY)(EY)'$
 $= \sigma^2 I_n + X\beta\beta'X'$

$$C(\hat{Y}, \hat{\epsilon}) = E(\hat{Y} - X\beta)(\hat{\epsilon} - 0)'$$

$$= E(\hat{Y}\hat{\epsilon}') = E\{(X\hat{\beta})(Y - X\hat{\beta})'\}$$

$$= E\{X\hat{\beta}Y'\} - E(X\hat{\beta}\hat{\beta}'X')$$

$$= X(X'X)^{-1}X'E(Y'Y) - XE(\hat{\beta}\hat{\beta}')X'$$

$$= X(X'X)^{-1}X'(\sigma^2 I_n + X\beta\beta'X')$$

$$- X(\sigma^2(X'X)^{-1} + \beta\beta')X'$$

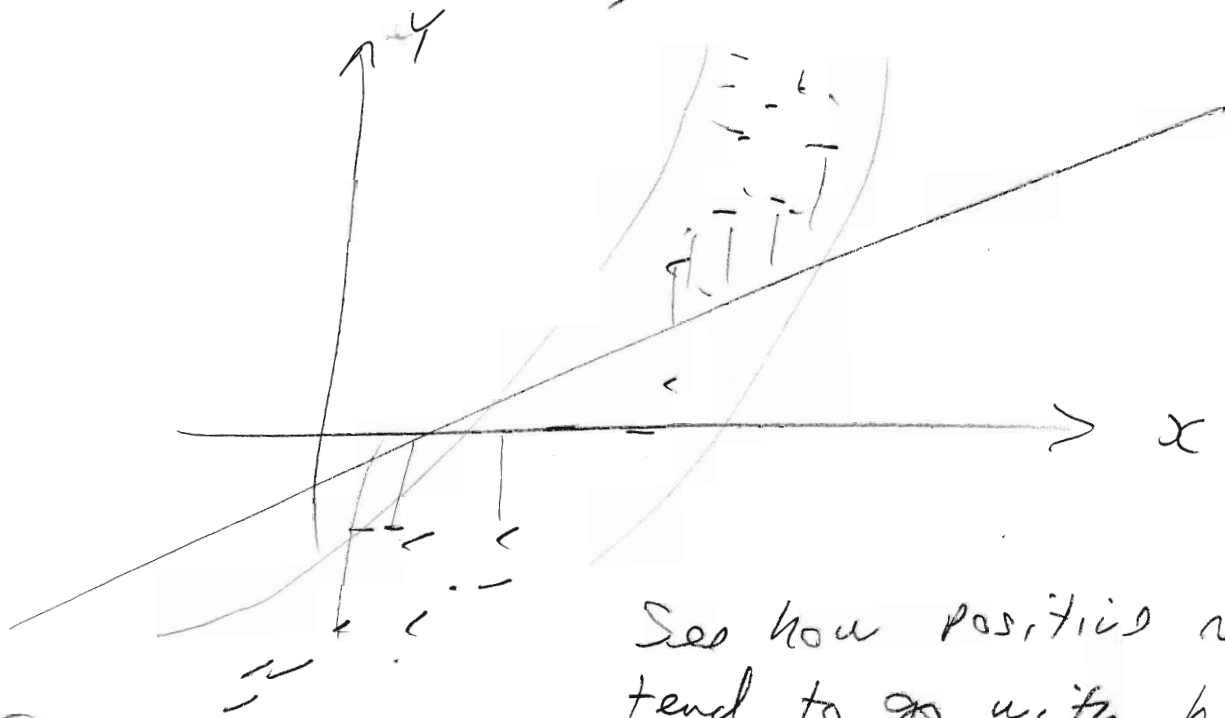
$$= \sigma^2 X(X'X)^{-1}X' + X(X'X)^{-1}X'X\beta\beta'X'$$

$$- \sigma^2 X(X'X)^{-1}X' - X\beta\beta'X' = 0$$

What might you see if model is wrong?

15.10

Well, outliers, or



See how positive residuals tend to go with higher predicted y ?

Not nec linear pattern

So we might see evidence of

- OUTLIERS
- CURVILINEAR PATTERN
- Systematic non-constant variance, like more error in the prediction of large values

(Reasonable: predict restaurant tip.)

To see systematic non-constant variance, better plot standardized or (deleted) studentized Resid.

Plot of x variables vs Residuals. What should you see? Recalling

$$R = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

Look at the sample correlation between x_i & $\hat{\varepsilon}_i$ values. Assuming the model has an intercept so $\sum \hat{\varepsilon}_i = 0$, numerator is

$$\begin{aligned} \sum_{i=1}^n (x_i - \bar{x})(\hat{\varepsilon}_i - 0) &= \sum_{i=1}^n (x_i - \bar{x})\hat{\varepsilon}_i = \sum_{i=1}^n x_i \hat{\varepsilon}_i - \bar{x} \sum_{i=1}^n \hat{\varepsilon}_i \\ &= \sum_{i=1}^n x_i \hat{\varepsilon}_i, \text{ Element } i \text{ of } \begin{bmatrix} 1 & \dots & 1 \\ \bar{x}_1 & \dots & \bar{x}_n \end{bmatrix} \begin{bmatrix} \hat{\varepsilon}_1 \\ \vdots \\ \hat{\varepsilon}_n \end{bmatrix} = X' \hat{\varepsilon} \end{aligned}$$

$$X' \hat{\varepsilon} = X'(Y - X\hat{\beta}) = X'Y - X'X(X'X)^{-1}X'Y = X'Y - X'Y = 0$$

So the sample correlation between IV & Residual is exactly zero, whether the model is correct or not.

There will be no linear trend, but

- Possibly curvilinear
- If curvilinear, symmetric



Regardless of the shape of the curve relating x to y

Plots of x vs $\hat{\epsilon}$

15.12

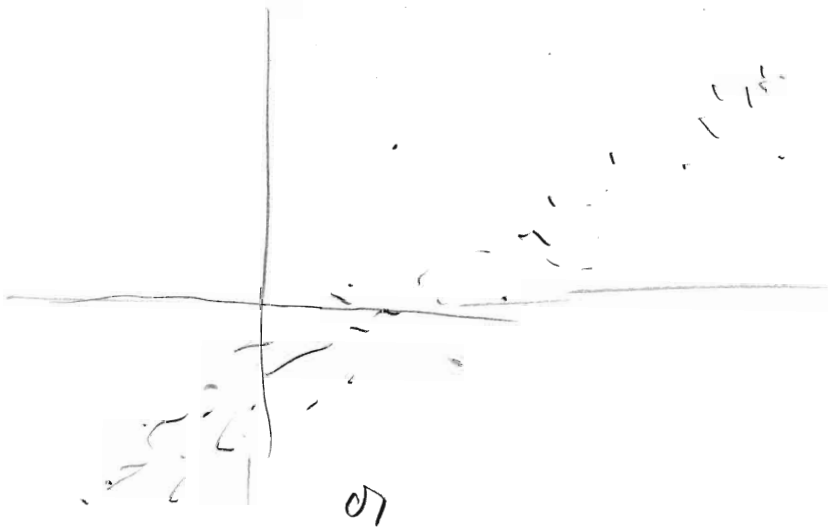
could also show

- Systematically non-constant variance



If data are collected in time sequence,

Plot time vs $\hat{\epsilon}$ or π or t



Tendency to under-predict earlier, over-predict later.

Betrays sloppy data collection



Maybe prediction gets better after some practice

Standardized (or deleted
Studentized) residuals are

15.13

- Easier to interpret (What's "BIG") and
- Have constant variance, so it's better to plot them rather than the raw residuals.

1
E

Zero correlation, Guaranteed
If model is correct, should see NO

1
1

~~POS or NEG~~
POS or NEG
Slope acc to
 β_1

$R^2_{\hat{y}} = R^2$
Nice scatterplot
Showing overall
strength of rel.

Independent of model is correct. Should
see NOTHING

X
X

ORDINARY
SCATTERPLOT

X

X

X

X

1

1

1
E