

Goodness of Fit

P.1

→ consider the regression model:

$$y = \alpha + \beta x + \varepsilon$$

→ we want to know how well the model explains the variance of  $y$

→ the predictions

$$\hat{y} = \hat{\alpha} + \hat{\beta} x \quad \text{where: } \hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}$$

$$\hat{\beta} = \frac{\text{cov}(y, x)}{\text{var}(x)}$$

must be compared ~~stated~~ with the observed values of  $y$

→ so compute the observed residuals

$$\hat{\varepsilon} = y - \hat{y}$$

→ note that

$$(y - \bar{y}) = (\hat{y} - \bar{y}) + \hat{\varepsilon}$$

→ also note that we use  $(y - \bar{y})$  to compute the variance of  $y$

→ so we're essentially going to look at variance today

$$\sum (y - \bar{y})^2 = \sum [(\hat{y} - \bar{y}) + \hat{\varepsilon}]^2$$

$$= \sum [(\hat{y} - \bar{y})^2 + 2 \cdot (\hat{y} - \bar{y})\hat{\varepsilon} + \hat{\varepsilon}^2]$$

$$= \sum (\hat{y} - \bar{y})^2 + 2 \cdot \sum (\hat{y} - \bar{y})\hat{\varepsilon} + \sum \hat{\varepsilon}^2$$

$\underbrace{\phantom{0}}$   
Focus on this term

$$\sum (\hat{y} - \bar{y})\hat{\varepsilon} = \sum (\hat{\alpha} + \hat{\beta}x - \bar{y})\hat{\varepsilon}$$

$$= \hat{\alpha} \underbrace{\sum_{\text{ZERO}} \hat{\varepsilon}}_{\text{ASSUMED}} + \hat{\beta} \underbrace{\sum x \hat{\varepsilon}}_{\text{ZERO}} - \bar{y} \underbrace{\sum_{\text{ZERO}} \hat{\varepsilon}}_{\text{ZERO}}$$

(by Gauss-Markov assumption)

→ so if Gauss Markov assumption  
of  $\text{cov}(\varepsilon, x) = 0$  holds then:

$$\sum (y - \bar{y})^2 = \sum (\hat{y} - \bar{y})^2 + \sum \hat{\varepsilon}^2$$

$$\begin{array}{lcl} \text{TOTAL} & = & \text{EXPLAINED} \\ \text{sum of} & & \text{sum of} \\ \text{SQUARES} & & \text{SQUARES} \\ (\text{TSS}) & & (\text{ESS}) \end{array} + \begin{array}{l} \text{RESIDUAL} \\ \text{sum of} \\ \text{SQUARES} \\ (\text{RSS}) \end{array}$$

→ if we have a good regression model, then RSS will be low and ESS will be high relative to TSS

$$\text{TSS} = \text{ESS} + \text{RSS}$$

$$1 = \frac{\text{ESS}}{\text{TSS}} + \frac{\text{RSS}}{\text{TSS}}$$

$$\boxed{R^2 = \frac{\text{ESS}}{\text{TSS}}} = \boxed{1 - \frac{\text{RSS}}{\text{TSS}}}$$

→ In theory, we could use either  $ESS/TSS$  or  $1 - \frac{RSS}{TSS}$  to compute  $R^2$ , but in practice we use the latter

$$R^2 = 1 - \frac{\sum \hat{\epsilon}^2}{\sum (y - \bar{y})^2}$$

→ Notice that when residuals are smaller  $R^2$  is higher

→ This measure essentially compares the variance of the residuals (remember that expected value of residual is zero) to the variance of  $y$

→ A higher  $R^2$  implies a better model fit

→ So  $R^2$  is great, but if you just keep adding variables to the model,  $R^2$  will rise

→ "Adjusted  $R^2$ " penalizes you for adding terms, but rewards you for better model fit

$$\text{Adj } R^2 = 1 - \frac{\sum \hat{\epsilon}^2 / (n-k)}{\sum (y - \bar{y})^2 / (n-1)}$$

where  $k$  is number of estimated parameters

~~W~~

→ degrees of freedom

- $(n-1)$  is df when computing the variance of  $y$

- $(n-k)$  is df of regression model

→ degrees of freedom

- suppose you have 100 observations
- to compute variance of  $\bar{y}$  you first must compute  $\bar{y}_j$ , so subtract one ~~term~~ to get df

$$df = 100 - 1 = 99$$

\*\*\*

- in the case of a regression model you have estimated  $k$  parameters

$$y = \alpha + \beta x$$

here  $k=2$  (one for  $\alpha$ , one for  $\beta$ )

$$\text{so } df = 100 - 2 = 98$$

\*\*\*

- the reverse occurs when working with explained sum of squares

et

- when working w/ ESS

$$\hat{y} = \hat{\alpha} + \hat{\beta}x$$

so only one variable determines  $\hat{y}$

$$df = 1$$

- Note however that if:

$$\hat{y} = \hat{\alpha} + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2$$

then there would be  $df = 2$

~~W~~

### F-test

→ in previous lectures we focused on estimating coefficients & comparing their values to the std error (testing for statistical significance of coefficient)

- But we may also want to test the overall model or compare one model to another
- For that purpose, we use F-test
- Simple example: test  $H_0: \beta = 0$

- Regression model:  $y = \alpha + \beta x + \varepsilon$
- if  $\beta = 0$ , then  $\bar{y} = \hat{\alpha}$   
because  $\bar{y} = \hat{\alpha} + \hat{\beta} \bar{x}$
- so we want to test  
 $y = \alpha + \varepsilon$  against  $y = \alpha + \beta x + \varepsilon$
- the regression model has  $(n-2)$  degrees of freedom
- the test imposes ONE restriction  
i.e.  $H_0: \beta = 0$ , so  $m = 1$

F-test

$$F = \frac{\frac{(RSS_R - RSS_{UR})}{m}}{RSS_{UR}/(n-k)} = \frac{\frac{(R^2_{UR} - R^2_R)}{m}}{\frac{(1-R^2_{UR})}{(n-k)}}$$

$RSS_R$  is restricted sum of squares  
in the "restricted case"  
(ie when  $\beta=0$ )

$RSS_{UR}$  is unrestricted RSS

$m$  is "numerator df" or in this case it's the number of ~~the~~ restrictions (in this case: one)

$(n-k)$  is "denominator df" or in this case it's df of regression model

$\approx$

Note that when  $\beta=0$  the ~~predicted values of~~ ~~regression~~ ~~equation~~ regression model becomes:

$$\text{or } \sum \hat{\epsilon}^2 = \sum (y - \bar{y})^2$$

$$RSS = TSS$$

and

$$\boxed{\begin{aligned} \hat{y} &= \alpha + \epsilon \\ \hat{\alpha} &= \bar{y} \end{aligned}}$$

$$R^2_R = 0$$