

- by dividing SSR + SSE by their df we obtain mean sums for regression called (MSR) + for error (MSE)

$$\text{ie. } \text{MSR} = \frac{\text{SSR}}{\text{df/SSR}} \quad \text{MSR} = \frac{\text{SSR}}{1} = \text{SSR}$$

$$\text{MSE} = \frac{\text{SSE}}{\text{df/SSE}} \quad \text{MSE} = \frac{\text{SSE}}{n-2} (=s^2)$$

identical g.l.s

ie. since  $s^2$  is an unbiased estimator of  $\sigma^2$  (ie.  $E(s^2) = \sigma^2$ ) and since  $s^2 = \text{MSE} \Rightarrow \text{MSE}$  is an unbiased estimator of  $\sigma^2$  (ie.  $E(\text{MSE}) = \sigma^2$ )

- it can be shown that under

$$H_0: \beta_1 = 0$$

MSR is also an unbiased estimator of  $\sigma^2$  where expected value of  $\text{MSR} = \sigma^2 + \beta_1^2 (\sum x_i - \bar{x})^2$

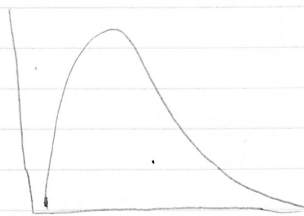
$\Rightarrow$  test-statistic

R.R.

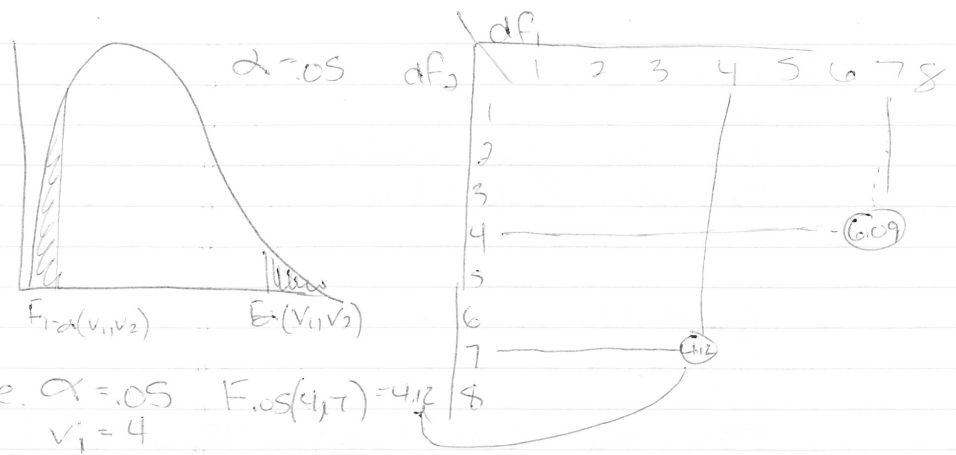
$$F = \frac{\text{MSR}}{\text{MSE}}$$

Reject  $H_0$  if  $F > F_{(1, n-2)}$   
accept  $H_0$  if  $F = 1$

### \* F-Dist



- skewed to the right (ie. positively skewed)
- not symmetric + it <sup>(only has)</sup> pos. values
- has 2 different df ( $v_1$  = numerator df and  $v_2$  = denominator df)



ie.  $\alpha = .05$   $F_{.05}(4, 7) = 4.12$   
 $v_1 = 4$   
 $v_2 = 7$

Not in this course

$$F_{1-\alpha}(v_1, v_2) \Rightarrow \frac{1}{F_{\alpha}(v_2, v_1)}$$

$$F_{.05}(4, 7) = \frac{1}{F_{.05}(7, 4)} = \frac{1}{6.09} = 0.16$$

ANOVA Table

$$TSS = SSR + SSE$$

$$df_{TSS} = df_{SSR} + df_{SSE}$$

source of variation	df	SS	MS	F
Regression	1	SSR	MSR	F
Error	n-2	SSE	MSE	
Total	n-1	TSS		

$H_0: \beta_1 = 0; \alpha$

$H_a: \beta_1 \neq 0$

test-stat. R.R.

$F = \frac{MSR}{MSE}$  We reject  $H_0$  if  $F > F_{\alpha}(1, n-2)$

Excont'd

$$n=10$$

$$\sum x_i = 460 \quad \sum x_i^2 = 23634$$

$$\sum y_i = 760 \quad \sum y_i^2 = 59816$$

$$\sum x_i y_i = 36854$$

$$\hat{\beta}_0 = 40.784$$

$$\hat{\beta}_1 = 0.176$$

$$TSS = S_{yy} = \sum y_i^2 - \frac{(\sum y_i)^2}{n} = 59816 - \frac{760^2}{10} = \underline{2,056}$$

$$SSR = \frac{(S_{xy})^2}{S_{xx}} = \frac{\left[ \sum x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n} \right]^2}{S_{xx} \left[ \sum x_i^2 - \frac{(\sum x_i)^2}{n} \right]} = \frac{(36854 - \frac{460 \cdot 760}{10})^2}{23634 - \frac{460^2}{10}} =$$

$$\frac{(1,894)^2}{2,474} = \underline{1,449.9741}$$

$$SSE = TSS - SSR = 2,056 - 1,449.9741 = 606.02587$$

$$MSR = \frac{SSR}{1} = \underline{1,449.9741}$$

$$MSE = \frac{SSE}{n-2} = \frac{606.02587}{8} = 75.7532 \quad (FS^2)$$

F-test-stat

$$F = \frac{MSR}{MSE} = \underline{19.14}$$

Source of variation      df.

Variation	df	SS	MS	F
Regression	1	1449.9741	1449.9741	19.14
Error	n-2=8	606.02587	75.7532	
Total	n-1=9	2056		

always better to add 1 p. even if using rounded values.

$$H_0: \beta_1 = 0 \quad \text{-not linear}$$

$$H_a: \beta_1 \neq 0 \quad \text{-linear} \quad \alpha = .05$$

test stat.

$$F = \frac{MSR}{MSE} = 19.14$$

R.R. we reject if

$$F > F_{\alpha}(1, n-2) = F_{.05}(1, 8) = 5.32$$

Since  $F = 19.14 > 5.32$ , we reject  $H_0$  & conclude that @ 5% level of significance there is evidence that math test results & final stats grade are linearly related.

F & t test must give same ~~answers~~ <sup>results</sup>

NOTE: t-test & F-test must give you identical answers

-in fact t-test & F-test are equivalent only if we are testing 1 parameter &  $H_0: \beta_1 = 0$

otherwise we can only use t-test for testing  $H_0: \beta_1 < 0$  or  $H_0: \beta_1 \leq 0$   
 $H_a: \beta_1 < 0$        $H_a: \beta_1 > 0$

i.e.  $t^2 = F$

$$\left( \frac{t_{\alpha/2; n-2}}{1} \right)^2 = F_{\alpha}(1, n-2) \quad \text{eg. } t^2 = (4.375)^2 \Rightarrow 19.14 = F$$

$$t_{.025; 8}^2 = (2.306)^2 = 5.3176 = 5.32 = F_{.05}(1, 8)$$

$$* \left( \frac{t_{(1-\alpha/2)} \hat{\beta}_1}{\sqrt{S^2}} \right)^2 = \frac{\hat{\beta}_1^2}{S^2} = \frac{\hat{\beta}_1^2}{MSE} = \frac{\hat{\beta}_1^2 \cdot S_{xx}}{MSE} = \frac{S_{xy}^2}{S_{xx} \cdot MSE}$$

$$= \frac{(S_{xy})^2}{S_{xx} \cdot MSE} = \frac{SSR/1}{MSE} = \frac{MSR}{MSE} = F(1; n-2)$$

Review:  $X$  vs  $y$  (scatter plot)

② SLR model  $y = \beta_0 + \beta_1 x + \epsilon$

or  $E(y) = \beta_0 + \beta_1 x$ ;  $\epsilon \sim N(0, \sigma^2)$ ,  $\forall x$

$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$  - least squares line

$H_0: \beta_1 = 0$  (no linear)

$H_a: \beta_1 \neq 0$  (linear)

$S^2 = SSE$

$(n-2)$

①  $t = \frac{\hat{\beta}_1}{\sqrt{\frac{S^2}{S_{xx}}}} = \frac{\hat{\beta}_1}{\frac{s}{\sqrt{S_{xx}}}}$  R.R. - Reject  $H_0$  if  $t > t_{\alpha/2, n-2}$  or  $t < -t_{\alpha/2, n-2}$

②  $F = \frac{MSR}{MSE}$  R.R. - we reject  $H_0$  if  $F > F_{\alpha}(1, n-2)$

$MSR = \frac{SSR}{1}$

$MSE = \frac{SSE}{n-2} = S^2$

$t^2 = F$  only if we are testing one parameter  
 $H_0: \beta_1 = 0$  v  $H_a: \beta_1 \neq 0$

### Inferences Concerning $E(y)$

- assuming  $H_0: \beta_1 = 0$  has been rejected (i.e. we have a linear relationship btwn  $x$  &  $y$ )
- $\Rightarrow$  We might be interested in following/estimating the average  $\overset{E(y)}{}$  sales per month for a given level of expenditure
- mean  $\overset{E(y)}{}$  response for a given drug dose
- average  $\overset{E(y)}{}$  grade in stats for a given math test score  $x$

$\hat{y} = \hat{E}(y) \rightarrow$  average mean of expected value

- here model is  $y = \beta_0 + \beta_1 x + \epsilon$

or  $E(y) = \beta_0 + \beta_1 x$

$\Rightarrow \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \leftarrow$  fitted line (regression) can be used to estimate the unknown pop. line  $t(y) = \beta_0 + \beta_1 x$

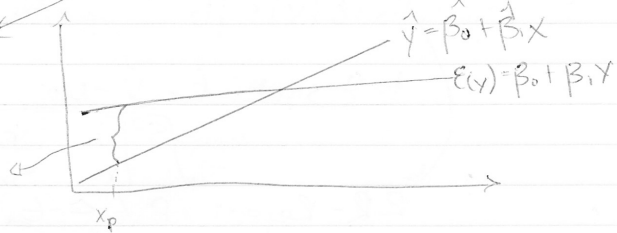
$E(\text{estimator}) = \text{parameter}$

$$E(\hat{y}) = E(\hat{\beta}_0 + \hat{\beta}_1 x) = E(\hat{\beta}_0) + E(\hat{\beta}_1) = \beta_0 + \beta_1 x = E(y)$$

unbiased estimator  $\rightarrow$   $\sigma^2$  unknown

$$V(\hat{y}) = \sigma^2 \left[ \frac{1}{n} + \frac{(x_p - \bar{x})^2}{S_{xx}} \right]$$

$\hat{y} \sim N(\underbrace{E(y)}_{\text{mean}}, \underbrace{V(\hat{y})}_{\text{variance}})$   
Error of estimating  $E(y)$  (when using regression line)



$$V(\hat{y}) = s^2 \left[ \frac{1}{n} + \frac{(x_p - \bar{x})^2}{S_{xx}} \right] = \text{MSE} \left[ \frac{1}{n} + \frac{(x_p - \bar{x})^2}{S_{xx}} \right]$$

- We might be interested, for example, whether the mean can yield per plot is 16 when we apply some amount of fertilizer (say 6 lbs)

$\sigma^2$  unknown  
confidence interval

$H_0: E(y) = 16$  ;  $\alpha$   $x_p = 6 \Rightarrow y_{x_p=6} = \beta_0 + \beta_1(6)$

$H_1: E(y) \neq 16$

test stat

$$t = \frac{\text{estimate} - E(\text{estimator})}{\sqrt{\text{var}(\text{estimator})}} = \frac{\hat{y} - E(y)}{\sqrt{s^2 \left[ \frac{1}{n} + \frac{(x_p - \bar{x})^2}{S_{xx}} \right]}}$$

$$\frac{\hat{y} - \mu_0}{\sqrt{\frac{\text{MSE}}{n} \left(1 + \frac{(b-x)^2}{S_{xx}}\right)}} \rightarrow \text{plug + solve}$$

- In general  $H_0: E(y) = \mu_0$ ;  $\alpha$   
 $H_a: E(y) \neq \mu_0$

test-stat:  $t = \frac{\hat{y} - \mu_0}{\sqrt{\frac{s^2}{n} \left(1 + \frac{(x_p - \bar{x})^2}{S_{xx}}\right)}}$

R.R. - Reject  $H_0$  if  $t > t_{\frac{\alpha}{2}, n-2}$  or  $t < -t_{\frac{\alpha}{2}, n-2}$

ii)  $H_0: E(y) > \mu_0$ ;  $\alpha$     iii)  $H_0: E(y) \leq \mu_0$ ;  $\alpha$   
 $H_a: E(y) < \mu_0$                        $H_a: E(y) > \mu_0$

test-stat  $t = \frac{\hat{y} - \mu_0}{\sqrt{\frac{s^2}{n} \left(1 + \frac{(x_p - \bar{x})^2}{S_{xx}}\right)}}$

R.R. - Reject if  $t < -t_{\alpha, n-2}$

R.R. - Reject if  $t > t_{\alpha, n-2}$

$(1-\alpha) 100\%$  C.I for  $E(y)$

$$E(y) \in \left( \hat{y} \pm t_{\frac{\alpha}{2}, n-2} \sqrt{\frac{s^2}{n} \left(1 + \frac{(x_p - \bar{x})^2}{S_{xx}}\right)} \right) = \left( t_{\frac{\alpha}{2}, n-2} \sqrt{\frac{\text{MSE}}{n} \left(1 + \frac{(x_p - \bar{x})^2}{S_{xx}}\right)} \right)$$

$$\hat{y}|x_p = \hat{\beta}_0 + \hat{\beta}_1 x_p$$

Ex) cont

- Find a 95% C.I. for the  $E(y)$  (Average) expected final grade in stats course given that the math test score was 50 (i.e.  $x_p = 50$ )

$$1 - \alpha = .95 \Rightarrow \alpha = .05 \Rightarrow \alpha/2 = .025$$

$$E(y) \in \left( \hat{y} \pm t_{\alpha/2, n-2} \sqrt{\text{MSE} \left[ \frac{1 + (x_p - \bar{x})^2}{n S_{xx}} \right]} \right)$$

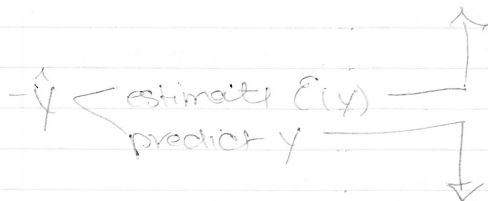
$$\begin{aligned} \hat{y} &= \beta_0 + \beta_1 x_p = \\ &= 40.784 + (0.766)(50) \\ &= 79.08 \end{aligned}$$

$$\begin{aligned} \text{MSE} &= 75.7532 \\ S_{xx} &= 2,474 \\ \bar{x} &= 46 \end{aligned}$$

$$= 79.08 \pm \underbrace{t_{0.025, 8}}_{2.306} \sqrt{75.7532 \left[ \frac{1 + (50 - 46)^2}{10 \cdot 2,474} \right]}$$

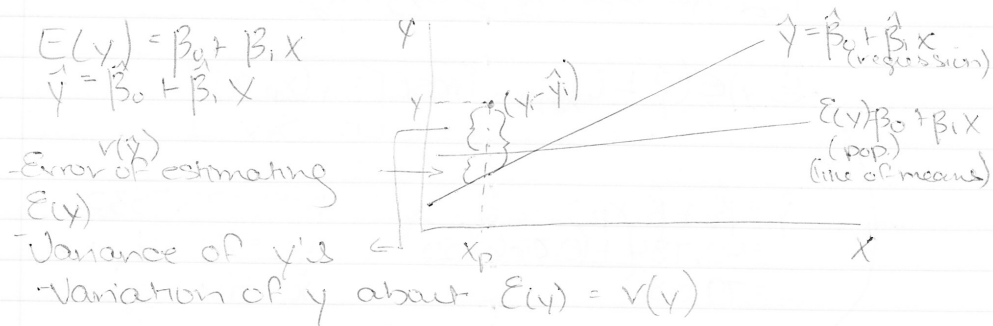
$$= 79.08 \pm 6.55 = \underline{(72.53, 85.63)} \text{ "on average" the grade will be then}$$

i.e. we are 95% confident (that in repeated response sampling) that for students who scored 50 on a math test, the average final grade in stats course will fall in the interval (72.53, 85.63)



Predicting  $y$  for a given value of  $x$  (say  $x_p$ )

Suppose instead of estimating all stats course grad, we want to predict the final grade in stats for a new student, selected from pop. of interest, who scored 50 on a math test ( $X_p = 50$ )



$$E(\hat{y}) = E(y)$$

$$v(X - \hat{y}_i) = v(y) - v(\hat{y}) = \sigma^2 + \sigma^2 \left( \frac{1}{n} + \frac{(X_p - \bar{X})^2}{S_{xx}} \right)$$

$$= \sigma^2 \left( 1 + \frac{1}{n} + \frac{(X_p - \bar{X})^2}{S_{xx}} \right)$$

(for prediction error is greater than for estimation)  
 - CI will be wider  
 $\hat{y} \sim N(E(y), v(y - \hat{y}))$

- but  $\sigma^2$  unknown, we use  $S^2$

$$\Rightarrow v(y - \hat{y}) = S^2 \left( 1 + \frac{1}{n} + \frac{(X_p - \bar{X})^2}{S_{xx}} \right) = \text{MSE} \left[ 1 + \frac{1}{n} + \frac{(X_p - \bar{X})^2}{S_{xx}} \right]$$

(1 -  $\alpha$ ) 100% Prediction Interval (P.I.) for  $y$

$$Y \in \left( \hat{y} \pm t_{\alpha/2, n-2} \sqrt{S^2 \left[ 1 + \frac{1}{n} + \frac{(X_p - \bar{X})^2}{S_{xx}} \right]} \right) \text{ can sub MSE for } \sigma^2$$

where  $\hat{y}|x_p = \hat{\beta}_0 + \hat{\beta}_1 x_p$

(Ex) cond - Find a 95% PI for  $y$  when  $x_p = 50$

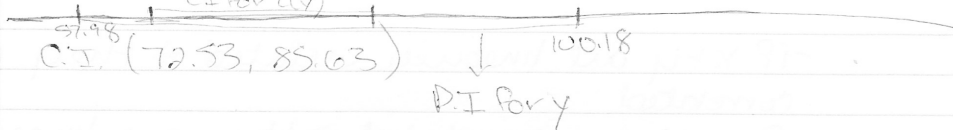
$$y \in \left( \hat{y} \pm t_{\alpha/2, n-2} \sqrt{\text{MSE} \left( 1 + \frac{(x_p - \bar{x})^2}{S_{xx}} \right)} \right) \quad \left| \begin{array}{l} 1 - \alpha = 95 \Rightarrow \alpha = .05 \\ \alpha/2 = .025 \end{array} \right.$$

$$= \left( 79.08 \pm 2.306 \sqrt{75.532 \left( 1 + \frac{(50 - 46)^2}{2474} \right)} \right)$$

$\hat{y} = 40.784 + (7.66)50 = 79.08$

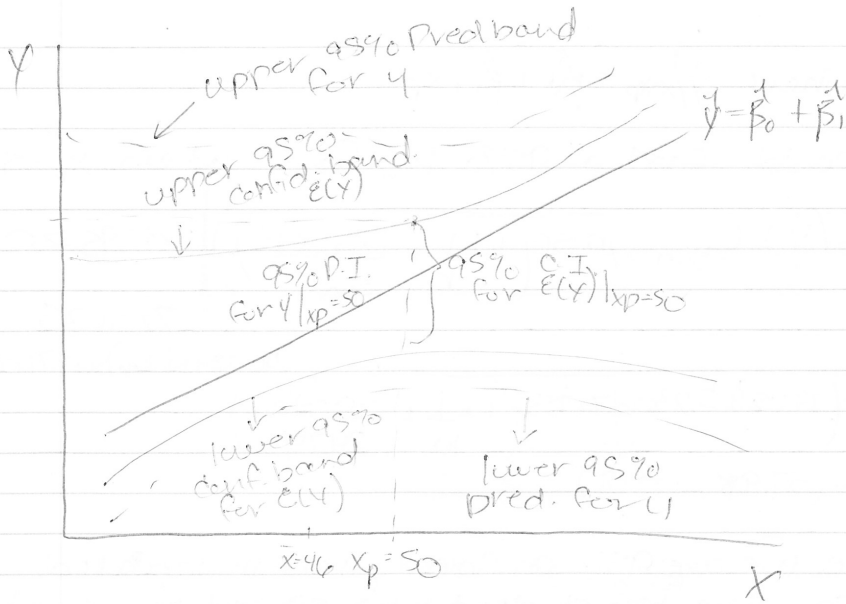
$$= (57.98, 100.18)$$

i.e. We are 95% confident, that in repeated sampling, <sup>for a</sup> ~~each~~ student who scored 50 on the math test, his/her final grade in stats will lie in  $(57.98, 100.18)$ .



NOTES • P.I. for an individual is always greater than C.I. for a group

- Both P.I. & C.I. have their min. value at the pt where  $X = \bar{X}$
- By plotting the end pts of the C.I.s or P.I.s for different values of  $x_p$  & connecting their end pts we obtain a general  $100\%(1-\alpha)$  confidence bands for  $E(y)$  or Prediction bands for  $y$



Correlation - measures the strength & direction of a linear relationship btwn  $X$  &  $Y$

- if  $X$  &  $Y$  are linearly related  $\Rightarrow$  they are correlated

- if  $X$  &  $Y$  are correlated  $\Rightarrow$  they are linearly related

$r$  - the sample "coefficient of correlation", also known as "Pearson product-moment coefficient of correlation"

$$r = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} \quad \text{or} \quad r = \beta_1 \sqrt{\frac{S_{xx}}{S_{yy}}} \quad \text{or}$$

- Properties of  $r$  "coefficient of correlation"

①  $-1 \leq r \leq 1$

②  $r = 0 \Rightarrow X$  &  $Y$  are not linearly related (ie they are not correlated)

③  $r > 0 \Rightarrow$  pos. linear relation.

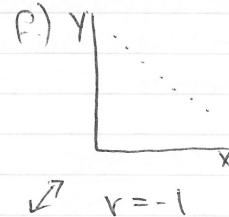
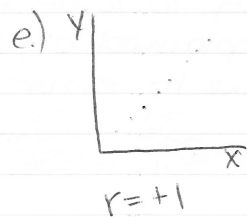
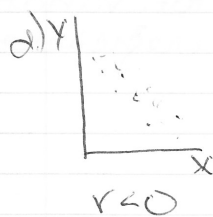
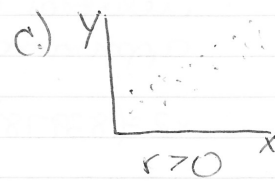
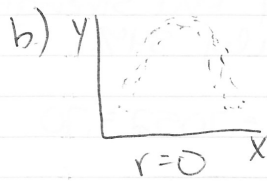
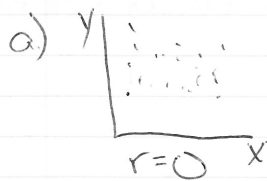
④  $r < 0 \Rightarrow$  neg. linear relation.

$\rightarrow$  cont

⑤  $r^2$  "correlation of determination"

$= \frac{SSR}{TSS}$  - measures how much of the tot. variation in the data can be explained by the regression line (hence how much is due to error.)

i.e.  $0 \leq r^2 \leq 1$  @ 70 model is iff



a) + b) - no lin. rel.  
 $\Rightarrow SSR = 0$   
 $\Rightarrow TSS = SSE$   
 $r^2 = \frac{SSR}{TSS} = \frac{0}{TSS} = 0$

perfect fit  $\Rightarrow$  no error.  
 $SSE = 0$   
 i.e.  $TSS = SSR$   
 i.e.  $r^2 = \frac{SSR}{TSS} = 1$

Ex) Covid

- we wish to calculate the coeff. of correlation btwn math test results + the final grade in Stats course + also calculate the coefficient of determination,  $r^2$

- scatter plot suggests a pos. linear relation. btwn  $x$  +  $y$



$$S_{xx} = 2474 \quad r = \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} = \frac{1894}{\sqrt{(2474)(2056)}} = 0.83978 = .84$$

$$S_{yy} = 2056$$

$$S_{xy} = 1894$$

i.e. math test results & final grade in stats are positively correlated (or positively linearly) related w/ the strength of their relationship approximately 84%.

$$r^2 = (.83978)^2 = .7052 = .70$$

i.e. approximately 70% of the tot. variation in the data is explained by the regression line, hence 30% is due to error

$$r^2 = \frac{SSR}{TSS} \quad (TSS = \overset{70}{SSR} + \overset{30}{SSE})$$

NOTES •  $r$  = sample coefficient of correlation and is

① an estimate of a unknown pop. coefficient of correlation

•  $\rho$  = pop. coeff. of correlation

②  $H_0: \rho = 0$  (no correlation)  $\Leftrightarrow H_0: \beta_1 = 0$  (no lin. relat.)  
 $H_a: \rho \neq 0$  (correlation)  $H_a: \beta_1 \neq 0$  (lin. relat.)

IF one test is done in a prob, you don't have to do the other, they are the same

(for correlation) test-stat.  $\Leftrightarrow t = \frac{\hat{\beta}_1}{\frac{s}{\sqrt{S_{xx}}}} = \frac{\hat{\beta}_1}{\frac{s}{\sqrt{S_{xx}}}}$

$\downarrow$   $t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$   $\Leftrightarrow$  equivalent  $\frac{s}{\sqrt{S_{xx}}}$   $\sqrt{\frac{s^2}{S_{xx}}}$

## Examining Lack of Fit in Linear Regression (Residual Analysis)

After finding

- So far we've been working w/ the model  $y = \beta_0 + \beta_1 x + \epsilon$  without checking of its validity + the validity of its assumptions
- ∴ We need to perform residual analysis to check whether the model + its assumptions have no violations.
  - should be performed immediately after finding the fitted regression line + before any hypothesis testing, estimation, +/or prediction
  - 2 ways of performing residual analysis:
    - ① graphical
    - ② numerical

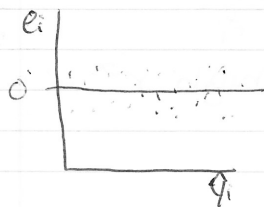
### ① graphical residual analysis

(i) scatter plot (plot  $x$  vs  $y$ ) ← to check for violations of linearity

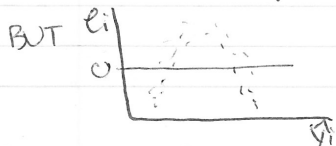
- to check for any outliers ect.

(ii) Plot of predicted values vs. residuals

(i.e.  $\hat{y}_i$  vs  $e_i$ )  $e_i = (y_i - \hat{y}_i)$  ← to check for any violations of assumptions of  $y$ 's or  $\epsilon$ 's being independently distributed + for violations of linearity.

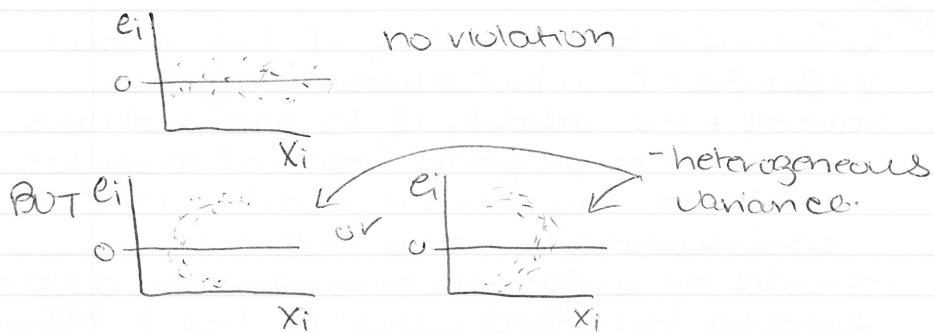


- IF no pattern  $\Rightarrow$  No violations

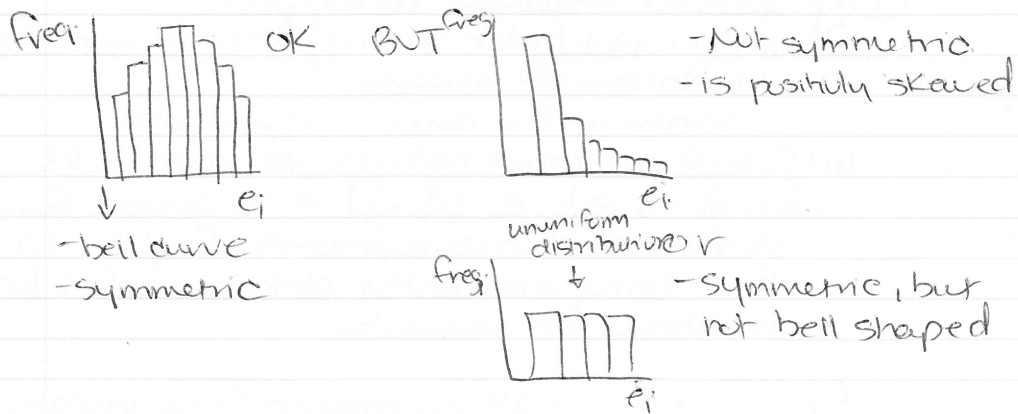


- curved linear pattern  $\Rightarrow$  violation  
 $\Rightarrow$  Higher order terms need to be added to the model

(iii) Plot  $X$  vs  $e_i$  ← To check for validity of the assumption of constant variance



(iv) Plot histogram of errors → To check validity of the assumption of errors being normally distributed



- Good fit = no violations + high  $r^2$

— missing Lecture