

# 3Y03 - PROBABILITY AND STATISTICS FOR ENGINEERING

WS19 Lecture 32

Last Time Even more on Linear Regression  $Y = \beta_0 + \beta_1 X + \varepsilon$

- Tests
  - ↳ on  $\beta_1$  slope (t-test for linearity)
- Confidence Intervals
  - ↳ for  $\beta_1$  slope and mean of  $Y$  at  $x_0$
- Prediction Intervals
  - ↳ for future value of  $Y$  at  $x_0$

All assume  $\varepsilon \sim N(0, \sigma^2)$  and use estimate  $\hat{\sigma}^2 = \frac{1}{n-2} SS_E$  &  $t_{n-2}$  - distribution

where

$$SS_E = SS_T - SS_R = \hat{\beta}_1 S_{xy}$$

"Error Sum of Squares"  $\sum e_i^2$       "Total Sum of Squares"  $SS_T = \sum y_i^2 - n\bar{y}^2$       "Regression Sum of Squares"  $SS_R$

$$SS_T = SS_E + SS_R$$

"Total variability of  $Y$ "  $\sum (y_i - \bar{y})^2$       "Variability of  $Y$  unexplained by regression"  $\sum (y_i - \hat{y}_i)^2$       "Amount of variability of  $Y$  accounted for by regression line"  $\sum (\hat{y}_i - \bar{y})^2$

ANOVA = Analysis of Variance Approach

$$SS_T = SS_E + SS_R$$

Also can show:

$$E\left(\frac{1}{n-2} SS_E\right) = \sigma^2$$

(from before)

$$E(SS_R) = \sigma^2 + S_{xy}$$

So we say  $SS_E$  has  $n-2$  d.o.f. &  $SS_R$  has 1 d.o.f. d.o.f. must balance so  $SS_T$  has  $n-1$  d.o.f.

It can be shown that, under  $H_0: \beta_1 = 0$ ,  
 ( $H_1: \beta_1 \neq 0$ )

$$F_0 = \frac{MS_R \leftarrow \begin{matrix} \text{mean} \\ \text{squares} \\ \text{(regression)} \end{matrix}}{MS_E \leftarrow \begin{matrix} \text{mean} \\ \text{squares} \\ \text{(error)} \end{matrix}} := \frac{SS_R / 1}{SS_E / (n-2)}$$

$\sim$  "F-distribution with numerator d.o.f. = 1,  
 & denominator d.o.f. = n-2"

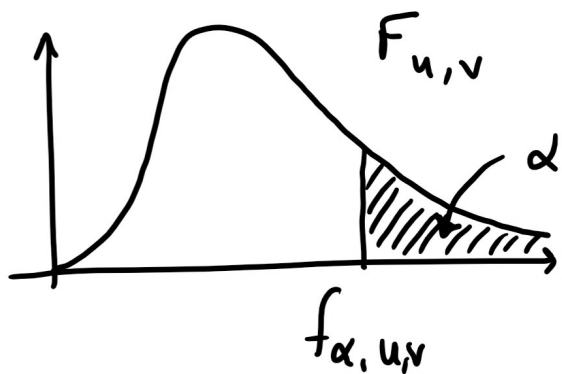
written  $F_{1, n-2}$

The bigger  $F_0$  is the more of variability is explained  
 by the regression so more evidence for linear rel.

i.e. when testing  $H_0: \beta_1 = 0$

$H_1: \beta_1 \neq 0$  (2-sided)

reject  $H_0$  at sig. level  $\alpha$  if  $f_0 > \underbrace{f_{\alpha, 1, n-2}}$



Look up  $\alpha$  f-table,  
 column 1, row n-2

Example - see sheet

Notice:  $T_0 = \frac{\hat{\beta}_1}{\sqrt{\hat{\sigma}^2 / S_{xx}}}$

so  $T_0^2 = \frac{\hat{\beta}_1^2 S_{xx}}{\hat{\sigma}^2} = \frac{\hat{\beta}_1^2 S_{xy}}{SS_E / (n-2)} = \frac{SS_R}{SS_E / (n-2)} = F_0$

So for 2-sided test on  $H_0: (\beta_1 = 0)$ , t-test & f-test are equivalent (but for 1-sided, need to use t-test).

## 11.7 Adequacy of the Regression Model

Assumptions made ① there is a linear rel.

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

where  $E(\varepsilon_i) = 0$  ( $V(\varepsilon_i) = \sigma^2$ )

②  $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$  for  $i \neq j$

For tests on  $\beta_1$  : ③  $\varepsilon_i \sim N(0, \sigma^2)$   
(& C.I., P.I.)

Goal check assumptions

③  $\varepsilon_i$  has Normal distr.: make a prob. plot of observed values  
i.e.  $e_1, \dots, e_n$

Recall: Reorganize  $e_{(1)} < \dots < e_{(n)}$

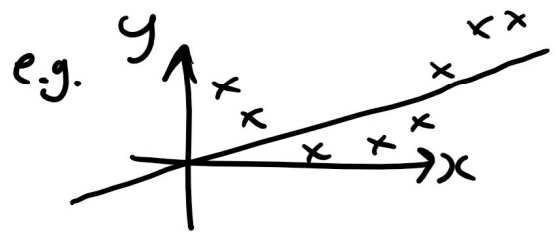
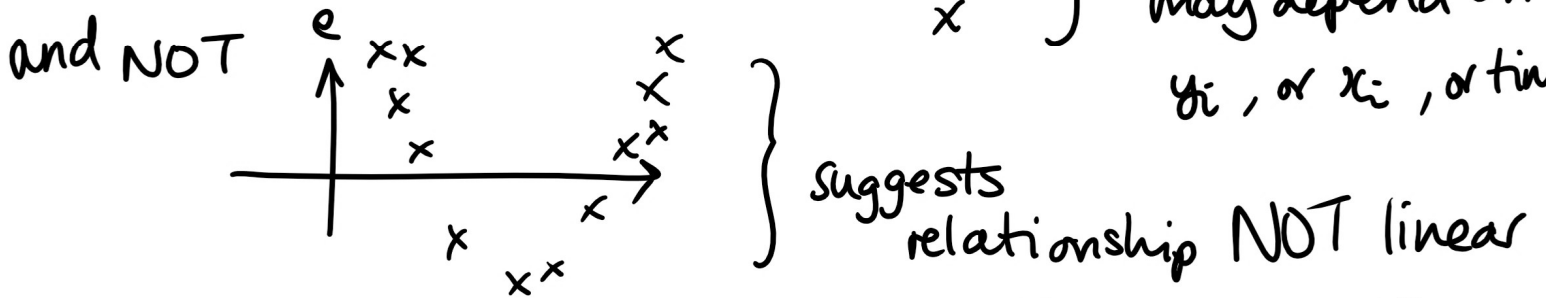
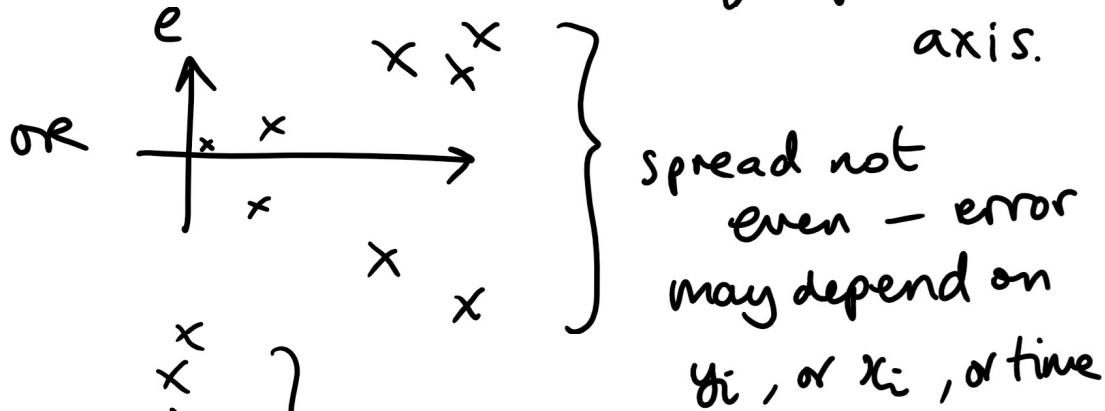
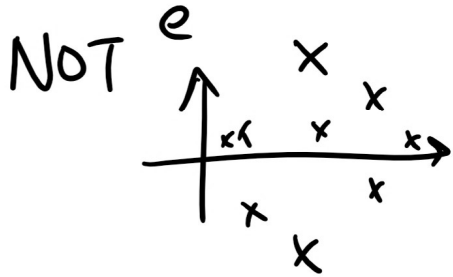
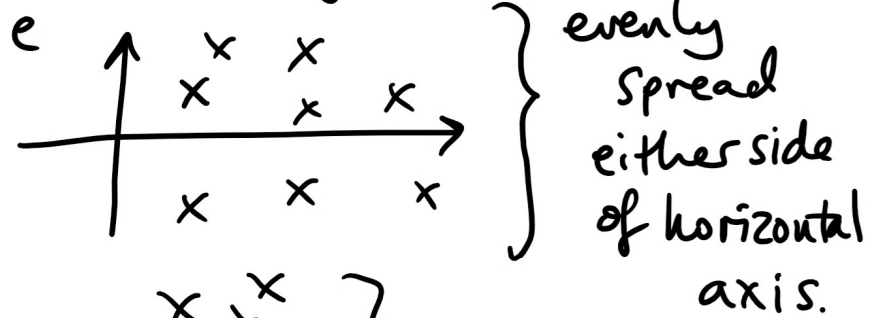
& plot  $e_{(j)}$  against  $z_j$  where

$$\Phi(z_j) = P(Z < z_j) = \frac{j - 0.5}{n}$$

If correct i.e.  $\varepsilon_i$  Normal  $\rightarrow$  get a straight line.

②  $Cov(\epsilon_i, \epsilon_j) = 0$  : plot  $e_i$  against  $y_i$ , or  $x_i$ , or time  
 $i \neq j$  (none of these should influence error)

So should get:



③ Linear rel. between X & Y

(a) t-test on  $\beta_1$

(b) ANOVA f-test on  $\beta_1$

(c) ANOVA:

The coefficient of determination  $R^2 = \frac{SS_R}{SS_T}$



- proportion of total variation in  $Y$  explained by regression  $\in [0, 1]$

- bigger : more likely to be a linear relationship

Example - see sheet.

## 11.8 Correlation

↳ assume  $X$  &  $Y$  jointly distributed

Recall underlying correlation coefficient:  $\rho = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$

Under our assumptions *it turns out that*

$\beta_1 = \frac{\sigma_y}{\sigma_x} \rho$  i.e. NO linear rel. exactly when  $\rho = 0$ .

Estimator for  $\rho$  : 
$$R = \frac{\sum (Y_i - \bar{Y})(X_i - \bar{X})}{\sqrt{\sum (X_i - \bar{X})^2 \sum (Y_i - \bar{Y})^2}}$$
$$= \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}}$$

Note :  $R^2 = \dots = (\text{coeff. of det.})^2$  \* see below

Example :  $R = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} = \frac{-16}{\sqrt{(32.8)(10)}} = \boxed{-0.88}$  *This says X&Y are very negatively correlated*

↳ Now added to separate sheet.

$$\textcircled{*} \underset{\substack{\uparrow \\ \text{sample correlation} \\ \text{coefficient}}}{(R)^2} = \frac{S_{xy}^2}{S_{xx} S_{yy}} = \left( \frac{S_{xy}}{S_{xx}} \right) \frac{S_{xy}}{S_{yy}} = \frac{\hat{\beta}_1 S_{xy}}{SS_T} = \frac{SS_R}{SS_T} = R^2$$

coefficient  
of determination

Phew!

$R^2 = R^2!$

Observe, though, that we can't get  $R$  from  $R^2$  without knowing something more:

in our running example on the separate sheet,  $R^2 = 0.784$ .

So we know  $R = \sqrt{0.784}$  or  $-\sqrt{0.784}$   
i.e. 0.88 or -0.88

In order to conclude  $R = -0.88$  this way we would need to know something more e.g. that the slope of the line  $\hat{\beta}_1$  is negative (and here we do know that;  $\hat{\beta}_1 = -0.49$ ).

### Least Squares Regression Example

$n=5$

i	x	y	$x^2$	$y^2$	xy
1	1	4	1	16	4
2	2	3	4	9	6
3	4	1	16	1	4
4	6	2	36	4	12
5	8	0	64	0	0
SUM	21	10	121	30	26



$$\bar{x} = \frac{21}{5} = 4.2 \quad \bar{y} = \frac{10}{5} = 2$$

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{-16}{32.8} = -0.49 \text{ slope}$$

$$S_{xy} = \sum xy - n\bar{x}\bar{y} = 26 - 5(4.2)(2) = -16$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1\bar{x} = 2 - (-0.49)(4.2) = 4.05 \text{ intercept}$$

$$S_{xx} = \sum x^2 - n\bar{x}^2 = 121 - 5(4.2)^2 = 32.8$$

$$\hat{y} = 4.05 - 0.49x$$

same

$$S_{yy} = \sum y^2 - n\bar{y}^2 = 30 - 5(2)^2 = 10$$

$$SS_T = S_{yy} = 10$$

$$\sigma^2 = \frac{1}{n-2} SS_E = \frac{1}{5-2} (10 - (-0.49)(-16)) = \frac{1}{3} (2.16) = 0.72$$

$$SS_R = \hat{\beta}_1 S_{xy} = (-0.49)(-16) = 7.84$$

$$R^2 = \frac{SS_R}{SS_T} = \frac{7.84}{10} = 0.784$$

$$f_0 = \frac{SS_R/1}{SS_E/n-2} = \frac{7.84}{2.16/3} = 10.8$$

$$SS_E = SS_T - SS_R = 10 - 7.84 = 2.16$$

$$R = \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} = \frac{-16}{\sqrt{(32.8)(10)}} = -0.88$$

$$f_{0.05, 1, 3} = 10.13 \quad \checkmark \text{ so reject } H_0$$

Tests for  $\beta_1$ :

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

$$\alpha = 0.05$$

95% CI for  $\beta_1$ :

$$t_0 = \frac{\hat{\beta}_1 - (\beta_1)_0}{\sqrt{\hat{\sigma}^2 / S_{xx}}} = \frac{-0.49 - 0}{\sqrt{0.72 / 32.8}} = -3.31$$

$$t_{\alpha/2, n-2} = t_{0.025, 3} = -3.182$$

more extreme so reject  $H_0$   
-yes, a linear rel.

$$\hat{\beta}_1 \pm t_{0.025, 3} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} = -0.49 \pm 3.182(0.14) = -0.49 \pm 0.47 = (-0.96, -0.02)$$

95% CI for Y at  $x_0=3$ :

$$(\hat{\beta}_0 + \hat{\beta}_1 x_0) \pm t_{0.025, 3} \sqrt{\hat{\sigma}^2 \left( \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)} = (4.05 - 0.49(3)) \pm 3.182 \sqrt{0.72 \left( \frac{1}{5} + \frac{(3-4.2)^2}{32.8} \right)}$$

95% PI for  $Y_0$  at  $x_0=3$ :

$$\hat{y} \pm t_{0.025, 3} \sqrt{\hat{\sigma}^2 \left( 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)} = (1.25, 3.91)$$

$$= (4.05 - 0.49(3)) \pm 3.182 \sqrt{0.72 \left( 1 + \frac{1}{5} + \frac{(3-4.2)^2}{32.8} \right)} = (-0.43, 5.59)$$

Notice 0 not in here:  
another way to reject

$$H_0: \beta_1 = 0 \text{ for } H_1: \beta_1 \neq 0$$

$$\sqrt{\frac{(3-4.2)^2}{32.8}}$$