## Last time   LINEAR REGRESSION

Model linear relationship between $X$ and $Y$ using

$$Y = \beta_0 + \beta_1 X + \varepsilon \qquad \leftarrow \text{random error}$$

with $E(\varepsilon) = 0$  $\boxed{V(\varepsilon) = \sigma^2}$

least squares estimates

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\sum_{i=1}^{n} x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^{n} x_i^2 - n\bar{x}^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Estimate $\sigma^2$ with

$$\hat{\sigma}^2 = \frac{1}{n-2} SS_E = \frac{1}{n-2}(SS_T - SS_R)$$

where $SS_T = \sum_{i=1}^{n} y_i^2 - n\bar{y}^2$, $SS_R = \hat{\beta}_1 S_{xy}$

$$= \sum_{i=1}^{n}(y_i - \bar{y})^2 \text{ (see below)}$$

## 11.3  Properties of Least Squares Estimators $(\hat{\beta}_0, \hat{\beta}_1)$.

For fixed $x$,  $Y = \beta_0 + \beta_1 x + \varepsilon$  is a r.v.

$\nwarrow E(\varepsilon) = 0$

So  $E(Y) = \beta_0 + \beta_1 x$  $\nwarrow V(\varepsilon) = \sigma^2$

$V(Y) = \sigma^2$.

Recall  $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\sum x_i \overset{Y_i}{y_i} - n\bar{x}\overset{\bar{Y}}{\bar{y}}}{\sum x_i^2 - n\bar{x}^2}$

$\sum = \sum_{i=1}^{n}$  everywhere

$$E(\hat{\beta}_1) = \frac{\sum x_i E(Y_i) - n\bar{x} E(\bar{Y})}{\sum x_i^2 - n\bar{x}^2}$$

$\bar{Y} = \frac{1}{n} \sum_{i=1}^{n} Y_i$

$E(\bar{Y}) = \frac{1}{n} \sum_{i=1}^{n} E(Y_i) = \frac{1}{n} \sum_{i=1}^{n} (\beta_0 + \beta_1 x_i) = \frac{1}{n} \cdot n\beta_0 + \beta_1 \frac{1}{n} \sum_{i=1}^{n} x_i = \beta_0 + \beta_1 \bar{x}$.

$$= \frac{\sum x_i (\beta_0 + \beta_1 x_i) - n\bar{x}(\beta_0 + \beta_1 \bar{x})}{\sum x_i^2 - n\bar{x}^2}$$

(Now, here, $\hat{\beta}_1$ is a function of r.v.s — an estimator for $\beta_1$, while in the last lecture it stood for the estimate for $\beta_1$ from the data)

$$= \frac{\beta_0 \left( \sum x_i - n\bar{x} \right) + \beta_1 \left( \sum x_i^2 - n\bar{x}^2 \right)}{\sum x_i^2 - n\bar{x}^2}$$

$$= \beta_1$$

$\hat{\beta_1}$ is an unbiased estimator for $\beta_1$.

And $V(\hat{\beta_1}) = \underset{\text{Check!}}{\ldots} \quad \sigma^2/S_{xx}$ $\quad \left( \begin{array}{l} \text{similarly !} \\ \text{Use 5.4 !} \end{array} \right)$

<span style="color:cyan">( Just to remind us that today we are talking estimators — functions of r.v.s $Y_i$, not of data points $y_i$ )</span>

& $E(\hat{\beta_0}) = E(\underset{Y}{\bar{y}} - \hat{\beta_1}\bar{x})$

$= E(\bar{Y}) - \bar{x} E(\hat{\beta_1})$

$= \beta_0 + \beta_1\bar{x} - \bar{x}\beta_1 = \beta_0 \quad$ so $\hat{\beta_0}$ <span style="color:cyan">is an</span>

<span style="color:cyan">unbiased estimator</span> for $\beta_0$.

And $V(\hat{\beta_0}) = V(\bar{Y} - \hat{\beta_1}\bar{x}) \qquad$ <span style="color:cyan">(See 5.4!)</span>

$= V(\bar{Y}) + \bar{x}^2 V(\hat{\beta_1})$

$= \frac{\sigma^2}{n} + \frac{\bar{x}^2 \sigma^2}{S_{xx}} = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right).$

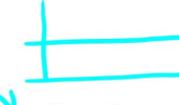Both variances depend on $\sigma^2$, so we can estimate with

$$\hat{\sigma}^2 = \frac{1}{n-2} SS_E. \qquad \text{<span style="color:cyan">(See last lecture.)</span>}$$

The estimated standard errors of $\hat{\beta_1}$ and $\hat{\beta_0}$ are their estimated standard deviations i.e.

$$se(\hat{\beta_1}) = \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \qquad \text{and} \qquad se(\hat{\beta_0}) = \sqrt{\hat{\sigma}^2 \left( \frac{1}{n} - \frac{\bar{x}^2}{S_{xx}} \right)}.$$

# 11.4 Hypothesis Tests in Simple Linear Regression

We want to test $H_0 : \beta_1 = (\beta_1)_0$ ← some #

$H_1 : $ e.g. $\beta_1 \neq (\beta_1)_0$

Slope = 0

e.g. $(\beta_1)_0 = 0$, which would mean NO linear relationship between X & Y.

We can't get anywhere without assuming something, so we assume $\varepsilon \sim N(0, \sigma^2)$ & so

i.e. this section will only apply in situations where this is a valid assumption

We found these above.

$$Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$

(linear comb. of Normal r.v.s is Normal)

We found these above.

So then $\hat{\beta_1} \sim N\left(\beta_1, \dfrac{\sigma^2}{S_{xx}}\right)$ ( $\parallel$ )

Under $H_0$ ($: \beta_1 = (\beta_1)_0$), $\dfrac{\hat{\beta_1} - (\beta_1)_0}{\sigma / \sqrt{S_{xx}}} \sim N(0,1)$

(Standardizing.)

Since we don't know $\sigma^2$ we're in similar territory to "Tests on Mean of Normal pop. Variance Unknown" setting.

We replace $\sigma^2$ with $\hat{\sigma}^2$ & get as test statistic:

$$T_0 = \dfrac{\hat{\beta_1} - (\beta_1)_0}{\hat{\sigma} / \sqrt{S_{xx}}} \sim t\text{-distribution with } n-2 \text{ degrees of freedom}$$

As usual, with 2-sided $H_1$, reject $H_0$ if

$$|t_0| > t_{\frac{\alpha}{2}, n-2} \quad \text{or} \quad -|t_0| < -t_{\frac{\alpha}{2}, n-2}.$$

Example  In the example from last time, we found

$$\hat{\beta_1} = -0.49, \quad S_{xx} = 32.8, \quad \hat{\sigma}^2 = 0.72.$$

Test the claim that there is a linear relationship between X and Y, at significance level $\alpha = 0.05$.

$$(n = 5).$$

Solution   $H_0 : \beta_1 = 0$   (no lin. relation.)

$H_1 : \beta_1 \neq 0$   ↳ We must put the "equality statement" as $H_0$.

$$t_0 = \frac{\hat{\beta_1} - 0}{\hat{\sigma}/\sqrt{S_{xx}}} = \frac{-0.49}{\sqrt{0.72}/\sqrt{32.8}} = -3.31.$$

Compare against $-t_{\frac{0.05}{2}, 5-2} = -t_{0.025, 3} = -3.182.$

$-3.31 < -3.182$  so reject $H_0$ for $H_1$. (Yes, there is a linear relationship.)

Tests for $\beta_0$   Exactly same idea. With assumption $\varepsilon \sim N(0, \sigma^2)$ we have

$$\hat{\beta_0} \sim N\left(\beta_0, \sigma^2\left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)\right).$$

$t$-test :  $H_0 : \beta_0 = (\beta_0)_0$     Use $T_0 = \dfrac{\hat{\beta}_0 - (\beta_0)_0}{\sqrt{\hat{\sigma}^2\left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)}}$

$H_1$ :e.g. $\beta_0 \neq (\beta_0)_0$

$\sim t$ – distribution

also with $n-2$ degrees of freedom.

## Analysis of Variance Approach (ANOVA)

Recall :  $SS_E = SS_T - SS_R$

$$\underset{\substack{\text{Sum of Squares}\\\text{Error}}}{\overset{n}{\underset{i=1}{\sum}} e_i^2} \quad " \qquad \Big| \qquad \underset{\substack{\text{Total sum}\\\text{of squares}}}{\overset{i}{\sum} y_i^2 - n\bar{y}^2} \qquad \overset{"}{\underset{}{\hat{\beta}_1 S_{xy}}}$$

$$= \overset{\hat{n}}{\underset{i=1}{\sum}} (y_i - \hat{y}_i)^2 \qquad = \sum (y_i - \bar{y})^2 \qquad = \sum (\hat{y}_i - \bar{y})^2$$

<span style="color:cyan">Regression sum of squares; residual amount only explained by regression.</span>

So  $SS_T = SS_E + SS_R$

$\qquad\qquad\qquad \uparrow \qquad\quad \uparrow$

$\qquad\qquad\quad n-2 \qquad$ 1 degree

$\qquad\qquad$ degrees of $\quad$ of freedom

$\qquad\qquad$ freedom

We have $E(SS_E / n-2) = \sigma^2$ $\qquad$ $E(SS_R / 1) = \overset{\text{Check!}}{\ldots} =$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \sigma^2 + \hat{\beta}_1 S_{xx}$

$\qquad\qquad\qquad\qquad\qquad\qquad$ <span style="color:cyan">$E(\hat{\beta}_1 S_{xy}) = \ldots$</span>

And  so

$$F_0 = \frac{SS_R / 1}{SS_E / (n-2)} \sim F\text{-distribution}$$

with 1 d.o.f. in numerator & $n-2$ d.o.f. in denominator.

T. B. C.

. . .

(Next time we'll explain something more about this — for now, the point is that we've come up with a test statistic — $F_0$ — whose underlying distribution is known & understood.)