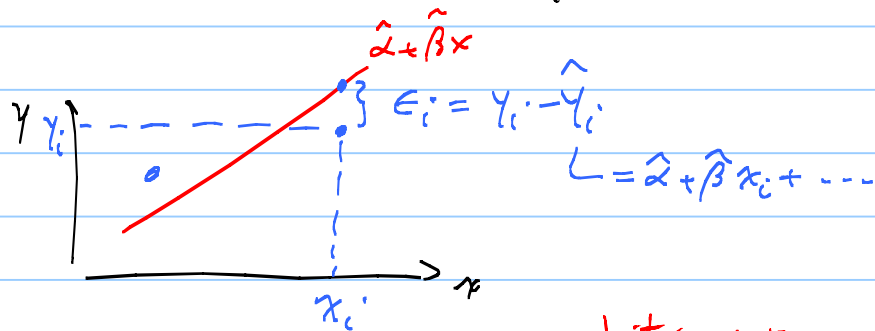# Lecture 26 (Ch. 11)

We did regression $y_i = \alpha + \beta x_i + \cdots + \epsilon_i$     Ch. 3.

We did inference on $\mu, \pi, \mu_2 - \mu_1, \pi_2 - \pi_1, \ldots$   Ch 7, 8.

Now we do inference on $\beta$ (and $\alpha$), $y, \ldots$     Ch. 11.

Review:



For a sample we write $y_i = \alpha + \beta x_i + \epsilon_i$ ← arbitrary params to be estimated by OLS.

$a, b$ in book

and   $\hat{y}_i = \hat{\alpha} + \hat{\beta} x_i$

where $\hat{\alpha}, \hat{\beta}$ are the OLS estimates of $\alpha, \beta$, ie.

$$\hat{\beta} = \frac{\overline{xy} - \overline{x}\,\overline{y}}{\overline{x^2} - \overline{x}} = \frac{S_{xy}}{S_{xx}} \quad , \quad \hat{\alpha} = \overline{y} - \hat{\beta}\,\overline{x} \quad .$$

where $S_{xx} = \sum\limits_{i=1}^{n} (x_i - \overline{x})^2$

$S_{xy} = \sum\limits_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y})$

Recall That

Sample Var. $= S_x^2 = \frac{1}{n-1} \sum (x_i - \overline{x})^2 = \frac{S_{xx}}{n-1}$

For a population, There exists an OLS fit as well!
The math/notation for obtaining That fit is exactly same as above.

How can we distinguish between The sample and The pop? E.g. $\overline{x}, \mu_x$
I will use The following notation for The predictions:

$\hat{y}(x) = \hat{\alpha} + \hat{\beta} x$ (for sample)     $y(x) = \alpha + \beta x$ (for population)

But you have to keep in mind The $\alpha, \beta$ in here are **not** arbitrary params to be estimated; They are OLS "estimates" obtained from The population.

Then There is The Analysis of Variance:

$$SST = \sum (y_i - \bar{y})^2 = \underbrace{SS_{explained}}_{\hat{\beta} S_{xy}} + \underbrace{SS_{unexpl.}}_{SSE}$$

df:    $n-1$    $=$    $k$    $+$    $n-(k+1)$

$$R^2 = \frac{SS_{expl.}}{SST}$$
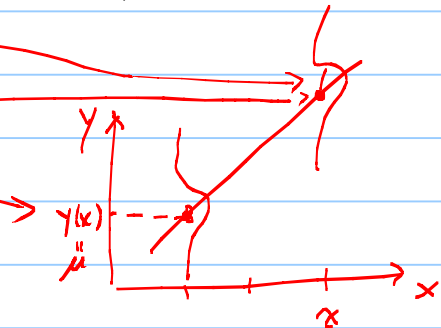
percent of var. in y
explained by x ...
(Goodness of fit)

$\# \text{ of } \beta's$
(excluding $\alpha$)
$\# \text{ of predictors}$
$y = \alpha + \beta_1 x_1 + \cdots \beta_k x_k$

$$S_e = \sqrt{\frac{SSE}{n-(k+1)}} \quad \sim RMSE$$

std. dev. of errors
$\sim$ Typical error
or spread about fit.

---

Now, to do inference we need a **Probability model (for regression)**:

Assume That The y's are **Normally** distr. **at each x**, with params
$\mu = Y(x)$ , $\sigma = \sigma_E$ (denoted $\sigma$ in book)

e.g. $Y(x) = \alpha + \beta x + \cdots$

This allows us to say things like:

1) $\hat{Y}(x) = \hat{\alpha} + \hat{\beta}x =$ estimates <u>mean</u> of y, <u>given</u> x
(You saw This in qz4)

2) In about 95% of The cases, we expect to have
y-values within $Y(x) \pm 1.96 \, \sigma_E$, for a given x

like 95% of cases are
within $\mu \pm 1.96\sigma$ (Ch.1)

3) other probs. e.g. $prob(a < y < b \mid x) =$

like $pr(a < x < b) =$ (Ch.1)

— True prediction = True mean | x.

$$prob\left( \frac{a - Y(x)}{\sigma_E} < \boxed{\frac{Y - Y(x)}{\sigma_E}} < \frac{b - Y(x)}{\sigma_E} \right) = Table \, T$$

$$Z \sim N(0,1)$$

$$pr\left( \frac{a-\mu}{\sigma} < \boxed{\frac{x-\mu}{\sigma}} < \frac{b-\mu}{\sigma} \right)$$

$$Z \sim N(0,1)$$

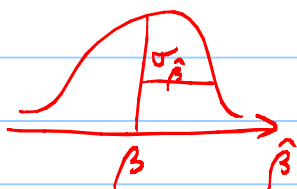Note: $Y \sim N(Y(x), \sigma_E^2) \implies E = Y - Y(x) \sim N(0, \sigma_E^2)$

Let's build a CI (and hyp. test) for ONE $\beta$ : $\quad Y_i = \alpha + \beta x_i + \epsilon_i$

**Theorem:** If $\epsilon \sim N(0, \sigma_\epsilon^2)$, Then $\hat{\beta}$ is normal with params:

Expected value (or mean) of The
sampling dist. of $\hat{\beta}$

$E[\hat{\beta}] \equiv \mu_{\hat{\beta}} = \beta \leftarrow$ pop. slope

$\sqrt{V[\hat{\beta}]} \equiv \sigma_{\hat{\beta}} = \dfrac{\sigma_\epsilon}{\sqrt{S_{xx}}} = \dfrac{\sigma_\epsilon}{\sqrt{n-1}\,(S_x)} \leftarrow$

Sample std. dev. of $x$.



Ch. 7
If $x \sim N(\mu_x, \sigma_x^2)$, Then
$\bar{x}$ is Normal with params

$E[\bar{x}] = \mu_{\bar{x}} = \mu_x$

$\sqrt{V[\bar{x}]} = \sigma_{\bar{x}} = \sigma_x/\sqrt{n}$

**Q1: What is The quantity That has a std. Normal dist?**

A) $\hat{\beta}$  B) $\dfrac{\hat{\beta} - \mu_Y}{\sigma_Y/\sqrt{n}}$  c) $\dfrac{\hat{\beta} - \beta}{\sigma_\beta/\sqrt{n}}$  D) $\dfrac{\hat{\beta} - \beta}{\sigma_{\hat{\beta}}}$

If $w \sim N(\mu_w, \sigma_w)$, Then $z = \dfrac{w - \mu_w}{\sigma_w} \sim N(0,1)$.

$z = \dfrac{\hat{\beta} - \beta}{\sigma_\epsilon/\sqrt{S_{xx}}} \sim N(0,1)$

$t = \dfrac{\hat{\beta} - \beta}{s_e/\sqrt{S_{xx}}} \sim t\text{-dist.} \quad df = n-2$ $\overset{k+1}{\downarrow}$

Ch. 7
$z = \dfrac{\bar{x} - \mu}{\sigma/\sqrt{n}} \sim N(0,1)$

$t = \dfrac{\bar{x} - \mu}{s/\sqrt{n}} \sim t\text{-dist.} \quad df = n-1$

Then, self-evident fact gives:

C.I. for $\beta$ : $\hat{\beta} \pm t^* \dfrac{s_e}{\sqrt{S_{xx}}} \qquad df = n-2 \text{ (Table VI)}$
$\qquad\qquad\qquad\qquad\qquad\qquad \underset{k+1}{\uparrow} \qquad\quad \text{or IV}$

$H_0: \beta \,\square\, \beta_0$

$H_1: \beta \,\square\, \beta_0$

$t_{obs} = \dfrac{\hat{\beta}_{obs} - \beta_0}{s_e/\sqrt{S_{xx}}}$

$df = n - 2 \overset{k+1}{\swarrow}$

p-value $= (1,2) \cdot pr(\hat{\beta} \,\square\, \hat{\beta}_{obs}) = pr(t \,\square\, t_{obs}) = $ Table VI
$\qquad\qquad \underset{1 \text{ or } 2 - \text{sided.}}{\curvearrowright}$

$\boxed{\text{problem 11.16}}$ [Revised; remove the word "positive", ie. do 2-sided]

$n = 13$   $x =$ nickel content,   $y =$ percentage austentite.

Data:   $\Sigma(x_i - \bar{x})^2 = 1.183$        $= S_{xx}$

$\Sigma(y_i - \bar{y})^2 = 0.0508$        $= S_{yy}$

$\Sigma(x_i - \bar{x})(y_i - \bar{y}) = 0.2073$   $= S_{xy}$

Question: Is There a statistically significant $(\alpha = 0.05)$
relationship between $x$ and $y$?

1) C.I. $\beta$ :   $\hat{\beta} \pm t^* \, S_e / \sqrt{S_{xx}}$

$\hat{\beta} = \dfrac{S_{xy}}{S_{xx}} = \dfrac{.2073}{1.183} = .1752 \longrightarrow$   SSE $= SST - \hat{\beta} S_{xy}$

$= .0508 - (.1752)(.2073)$

$S_e = \sqrt{\dfrac{SSE}{n-2}} = \sqrt{\dfrac{.014}{13-2}} = 0.0357$        $= .014$

$\therefore$ 95% CI for $\beta$:   $.1752 \pm 2.201 \left( \dfrac{.0357}{\sqrt{1.183}} \right) = 0.0328$        $= (0.10, 0.24)$
                                             $df = 13-2$

We are 95% Confident That The pop. $\beta$ is in here.

Also, Zero is not included $\Rightarrow$ Relationship is statistically significant

2) $H_0 : \beta = 0$      $t_{obs} = \dfrac{.1752 - 0}{.0328} = 5.31$ ,

$H_1 : \beta \neq 0$

p-value $= 2 \, pr(\hat{\beta} > \hat{\beta}_{obs}) = 2 \, pr(t > t_{obs})$

$= 2 \, pr(t > 5.31) < 0.001$

p-value $< \alpha$                                                         $\uparrow$

$\therefore$ Evidence That $\beta \neq 0$. (same conclusion        Table VI
                                as above).                                  $df = 13-2$

$$\boxed{FYI}$$

Note that the test of $\beta = 0$ is equivalent to testing
if there is a linear relationship between $x$ and $y$.
But if a linear relationship is all that you are testing,
then we can test the population correlation coeff

$$H_0: \varphi = 0$$

$$H_1: \varphi \neq 0$$

Typo, p. 506, blue box
$$H_a: \varphi \neq 0 \; (\text{no linear} \ldots$$

The test statistic for this test is a bit weird:

$$\Rightarrow \quad t = \frac{r - 0}{\sqrt{\frac{1-r^2}{n-2}}} \quad \text{has a } t \text{ distr, with } df = n-2.$$

Recall $r = S_{xy} / \sqrt{S_{xx} S_{yy}}$

This way, you take your data $(x_i, y_i)$, compute the sample
corel. coeff $(r)$, then $t_{obs}$, and then p-value,
all without any fitting.

3) For the above example:

$$H_0: \varphi = 0$$
$$H_1: \varphi \neq 0$$

$$r = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} = \cdots = .8456$$

$$t_{obs} = \frac{r - 0}{\sqrt{\frac{1-r^2}{n-2}}} = \cdots = 5.3 \quad \leftarrow \text{Same value as } t_{obs} \text{ we got above when testing } \beta.$$

p-value $= 2 \, \text{prob}(t > t_{obs}) = $ same as above.

∴ Same conclusion.

In summary: We have 3 ways of testing if
there is a useful relation between $x$ & $y$:

1) C.I. for $\beta$     2) Testing $H_0: \beta = 0$     3) $H_0: \varphi = 0$

The very beginning of section 3.3 in lab4 shows how to make/simulate data on x and y that are linearly associated. The x data consists of 100 cases from a uniform distribution, and the TRUE/population relationship between x and y is given by y = 10 + 2x.

a) What is the value of sigma_epsilon in that simulation?

b) Using the same settings used insection 3.3, write code to build the (empirical) sampling distribution of beta_hat based on 5000 trials. This code should produce a histogram.

c) According to the lecture, the mean of the histogram is supposed to be equal (or close) to what quantity? Is it?

d) According to the lecture, the standard deviation of that histogram is supposed to be equal (or close) to what quantity? Is it?

e) According the lecture, the distribution of the beta_hat is supposed to be normal with certain parameters. Use qqnorm() and abline() to confirm that.