

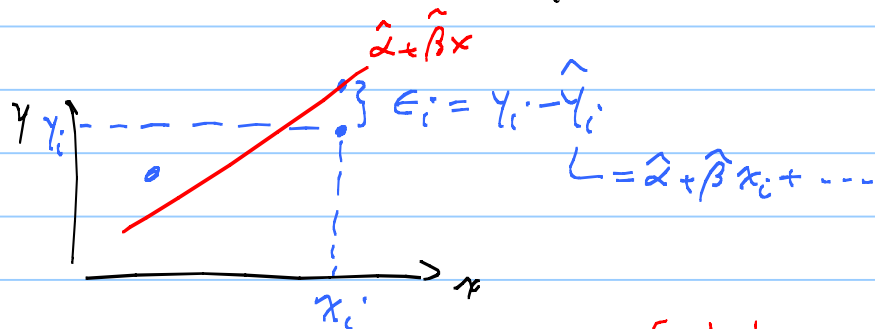
## Lecture 24 (Ch. 11)

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

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

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

Review:



For a sample we write  $y_i = \alpha + \beta x_i + \epsilon_i$  and  $\hat{y}_i = \hat{\alpha} + \hat{\beta} x_i$  [arbitrary params to be estimated by OLS, i.e.  $\frac{\partial}{\partial \alpha} SSE$ , etc.]

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

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

where  $S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2$

Recall That

$$\text{Sample var.} = s_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{S_{xx}}{n-1}$$

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

For a population, There exists an OLS fit as well!

Just because we can fit a line to the whole pop., it does not follow that there are no errors when predicting  $y$  from  $x$ !

So, now, in regression we need notation that can distinguish between pop. and sample quantities; like  $\bar{x}$  and  $\mu_x$ , but for regression.

I will use the following notation for the predictions:

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

But, in Ch. 11, you have to keep in mind that this  $\alpha, \beta$  are NOT free params that we can do things like  $\frac{\partial}{\partial \alpha} SSE$ , etc.; They are fixed quantities obtained by "fitting" a line to the whole population.

Then There is The Analysis of Variance:

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

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

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

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

$S_e = \sqrt{\frac{SSE}{n - (k+1)}}$   $\sim \text{RMSE}$   
 std. dev. of errors  
 $\sim$  Typical error or spread about fit.

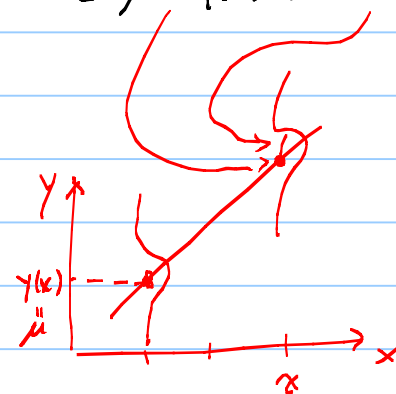
Now, to do inference we need a probability model (for regression):

Assume  $y$ 's are Normally distr. at each  $x$ , with params  $\mu = y(x)$ ,  $\sigma = \sigma_e$

e.g.  $\mu = y(x) = \alpha + \beta x + \dots$   $\sigma = \sigma_e = \text{fixed}$   
 estimate  $\alpha, \beta$  with  $\hat{\alpha}, \hat{\beta}$  estimate, with  $S_e$

Note:  $y \sim N(y(x), \sigma_e)$

$$e = y - y(x) \sim N(0, \sigma_e^2)$$



This allows us to say things like:

- 1)  $\hat{y}(x) = \hat{\alpha} + \hat{\beta}x = \text{estimates mean of } y, \text{ given } x$
- 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) o the v probs. e.g.  $\text{prob}(a < y < b | x) =$

True prediction = True mean  $| x$ .

$$\text{prob}\left(\frac{a - y(x)}{\sigma_e} < \frac{y - y(x)}{\sigma_e} < \frac{b - y(x)}{\sigma_e}\right) = \text{Table I}$$

$z \sim N(0, 1)$

like  $\text{pr}(a < x < b) =$  (Ch. 1)

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

$z \sim N(0, 1)$

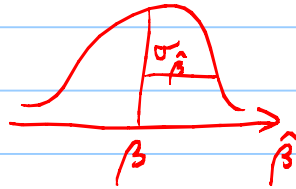
Let's build a CI (and hyp. test) for ONE  $\beta$  :  $y_i = \alpha + \beta x_i + \epsilon_i$

**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  $\hat{\beta}$

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

$$\sqrt{V[\hat{\beta}]} = \sigma_{\hat{\beta}} = \frac{\sigma_\epsilon}{\sqrt{S_{xx}}} = \frac{\sigma_\epsilon}{\sqrt{n-1} S_x}$$



$$S_{xx} = \sum_i^n (x_i - \bar{x})^2 = (n-1) S_x^2$$

Defn. of sample var.

Recall  $\sigma_\epsilon$  is const. and  $S_x$  does not vary as  $\frac{1}{n-1}$  because of  $\sum$  in the numerator of  $S_x$ .  
So,  $\sigma_{\hat{\beta}}$  falls off as  $1/\sqrt{n}$

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)  $\frac{\hat{\beta} - \mu_y}{\sigma_y / \sqrt{n}}$

C)  $\frac{\hat{\beta} - \beta}{\sigma_\beta / \sqrt{n}}$

D)  $\frac{\hat{\beta} - \beta}{\sigma_{\hat{\beta}}}$

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

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

$$t = \frac{\hat{\beta} - \beta}{s_\epsilon / \sqrt{S_{xx}}} \sim t\text{-dist. } df = n-2$$

Ch. 7

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

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

Then, self-evident fact gives:

C.I. for  $\beta$  :  $\hat{\beta} \pm t^* \frac{se}{\sqrt{S_{xx}}}$   $df = n-2$  (Table VI or IV)

$H_0: \beta \square \beta_0$

$H_1: \beta \square \beta_0$

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

p-value = (1,2)  $\cdot$   $pr(\hat{\beta} \square \hat{\beta}_{obs}) = pr(t \square t_{obs})$  = Table VI

$\uparrow$  1 or 2-sided.

problem 11.17 [Revised; remove the word "positive", i.e. do 2-sided]

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

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

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

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

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

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

$$\hat{\beta} = \frac{S_{xy}}{S_{xx}} = \frac{.2073}{1.183} = .1752$$
$$SSE = SST - \hat{\beta} S_{xy} = .0508 - (.1752)(.2073) = .014$$

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

$$\therefore 95\% \text{ CI for } \beta: .1752 \pm 2.201 \left( \frac{.0357}{\sqrt{1.183}} \right) = 0.0328 = (0.10, 0.24)$$

$df=13-2$

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

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

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

$H_1: \beta \neq 0$

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

$$= 2 \Pr(t > 5.31) < 0.001$$

$$p\text{-value} < \alpha$$

$\therefore$  Evidence That  $\beta \neq 0$ . (same conclusion as above).

Table VI  
 $df = 13-2$



In Summary, we have 2 ways of testing an association between  $xy$ .  
(A Third way, next FYI)

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

$$H_0: \rho = 0$$

$$H_1: \rho \neq 0$$

The test statistic for this test is a bit weird:

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

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

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

3) For the above example:

$$H_0: \rho = 0$$

$$H_1: \rho \neq 0$$

$$r = \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} = \dots = 0.8456$$

$$t_{obs} = \frac{r - 0}{\sqrt{\frac{1-r^2}{n-2}}} = \dots = 5.3$$

← same value as  $t_{obs}$  we got above when testing  $\beta$ .

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

∴ some 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)  $H_0: \rho = 0$

hw-lect 24-1

by R

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  $y = 10 + 2x$ .

- What is the value of  $\sigma_{\epsilon}$  in that simulation?
- Using the same settings used in section 3.3, write code to build the (empirical) sampling distribution of  $\beta_{\text{hat}}$  based on 5000 trials. This code should produce a histogram.
- According to the lecture, the mean of the histogram is supposed to be equal (or close) to what quantity? Is it?
- According to the lecture, the standard deviation of that histogram is supposed to be equal (or close) to what quantity? Is it?
- According the lecture, the distribution of the  $\beta_{\text{hat}}$  is supposed to be normal with certain parameters. Use `qqnorm()` and `abline()` to confirm that.

hw-lect 24-2

In a problem dealing with flow rate ( $y$ ) and pressure drop ( $x$ ) across filters, it is known that  $y = -0.12 + 0.095x$ . I.e. This is the "fit" to the population. Suppose it is also known that  $\sigma_{\epsilon} = 0.025$ . Now, IF we were to make repeated observations of  $y$  when  $x = 10$ , what's the prob. of a flow rate exceeding 0.835?

Hint:  $\frac{y - (\text{true mean of } y \text{ at some } x)}{\sigma_{\epsilon}} \sim N(0, 1)$ .

This document was created with Win2PDF available at <http://www.win2pdf.com>.  
The unregistered version of Win2PDF is for evaluation or non-commercial use only.  
This page will not be added after purchasing Win2PDF.