# Lecture 11 (Ch.3)

We did regression (fitting) by assuming a model for data $(x_i, y_i)$:

$$y_i = \alpha + \beta x_i + \epsilon_i \leftarrow \text{error/residual.}$$

obs. y at $x_i$    y of line $y(x) = \alpha + \beta x$   at $x = x_i$

To find the "best" $\alpha, \beta$ (ie. line), we minimized SSE:

$$SSE = \sum_i^n \epsilon_i^2 = \sum_{i=1}^n \left[ y_i - (\alpha + \beta x_i) \right]^2$$

obs       pred.

Compare!       Compare!

and got $\quad \hat{\beta} = \dfrac{\overline{xy} - \overline{x}\,\overline{y}}{\overline{x^2} - \overline{x}} \quad, \quad \hat{\alpha} = \overline{y} - \hat{\beta}\,\overline{x}.$

Then, the eqn of the "best fit" is $\quad \hat{y}(x) = \hat{\alpha} + \hat{\beta} x$

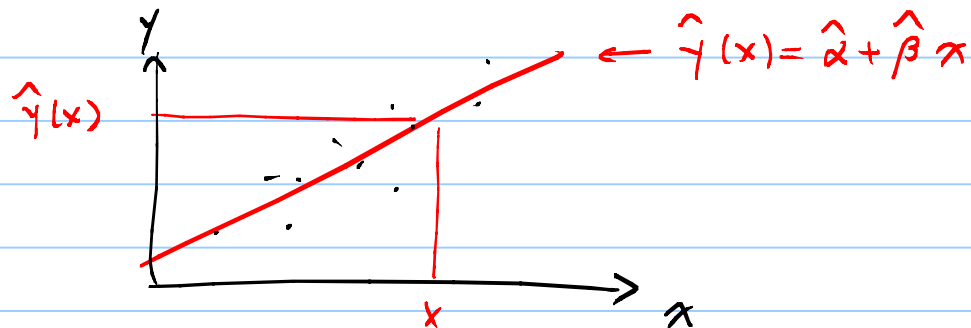Note: $\hat{y}(x_i) = \hat{\alpha} + \hat{\beta} x_i \quad (\text{no } \epsilon\,!)$

$\hat{y}(x_i)$ is sometimes written as $\hat{y}_i$.

---

I forgot to mention that regression is most useful when $x$ is easy to measure, but $y$ is hard to measure.

E.g.   $x$ = Blood flow velocity (FV) with ultra sound.
      $y$ = Intracranial Pressure (ICP).

When $\hat{\alpha}, \hat{\beta}$ are obtained from regression, then, given $\underset{\text{FV}}{\textcircled{$x$}}$, we can predict $\underset{\text{ICP}}{\textcircled{$y$}}$ from $\hat{y}(x) = \hat{\alpha} + \hat{\beta} x$.   (No $\epsilon_i\,!$)
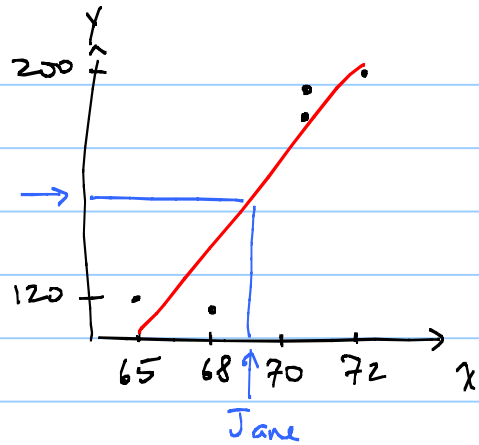


---

Also $\hat{\beta} = \dfrac{S_{xy}}{S_{xx}}$ where $S_{xx} = \sum_i (x_i - \overline{x})^2$   $S_{xy} = \sum_i (x_i - \overline{x})(y_i - \overline{y})$

Example

or FV ⏞ Data ⏞ and ICP

| height (x) | weight (y) | xy | x² |
|---|---|---|---|
| 72 | 200 | · | · |
| Joe: 70 | 180 | | |
| 65 | 120 | | |
| 68 | 118 | | |
| 70 | 190 | | |
| $\overline{x}$ | $\overline{y}$ | $\overline{xy}$ | $\overline{x^2}$ |

$$\hat{\beta} = \frac{\overline{xy} - \overline{x}\,\overline{y}}{\overline{x^2} - \overline{x}^2} = \frac{11224.8 - 69(161.6)}{4766.6 - 69(69)} = 13.28$$

**Interpret:** A change of 1 in is associated with an avg. change of 13.28 pounds.

$$\hat{\alpha} = \overline{y} - \hat{\beta}\overline{x} = 161.6 - 13.28(69) = -755$$

$$\text{lm}(y \sim x) \implies \hat{\beta} = 13.3, \ \hat{\alpha} = -755.11 \implies \hat{y}(x) = -755 + 13.28x$$

⟹ E.g. Joe's predicted weight, according to his height, is

$$\hat{y} = 13.28(70") - 755.11 \approx 174.9 \text{ pounds.}$$

⟹ We can now predict everyone's (weight) or ICP, from their (height) or FV.

| Height (x) | Weight (y) | ... | $\hat{y}$ | $(y - \hat{y})$ |
|---|---|---|---|---|
| 72 | 200 | ... | 201.5 | -1.5 |
| Joe= 70 | 180 | | 174.9 | 5.1 |
| 65 | 120 | | 108.5 | 11.5 |
| 68 | 118 | | 148.3 | -30.3 |
| 70 | 190 | | 174.9 | 15.1 |

$\hat{y} = \hat{\alpha} + \hat{\beta}x$ predicted y

any other fit will have a larger SSE.

⟹ For the people in the data set, we can also find their error/residual

⟹ For people outside the data set (e.g. Jane) we can predict their y from their x, but we cannot compute error, because we don't know their true y. In Ch.11, we'll address this issue.

However, be WARNED if you extrapolate

$$x = 0 \implies y = -755 \text{ pounds!}$$

<u>Shifting gears again.</u>

There is a <u>different</u> (more useful) way of looking at regression, via variance. This way, we will arrive at quantities called $R^2$ an $s_e$. which together assess how good the fit is.
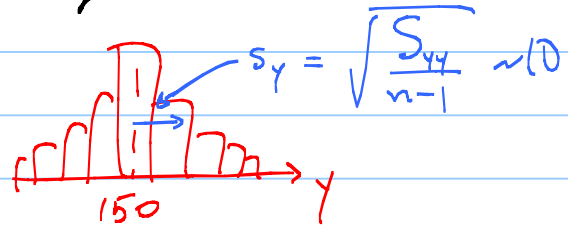
Let me motivate it:

→ Suppose we measure my Tablet's Length.

→ Repeat, and histogram:

→ One may report :

<span style="color:red">True length = 150 ± 10 cm</span>

$s_y = \sqrt{\dfrac{S_{yy}}{n-1}} \sim 10$



150

→ Now, suppose you are unhappy with the <u>large</u> $s_y$. <span style="color:red"><u>low precision.</u></span>

→ You may wonder, could some of that variability be due to something else that is varying everytime you make a measurement of $y$. <span style="color:red"><u>x</u> = temperature ? humidity ?</span>

If so, then by measuring <u>y</u> and <u>x</u>, we may be able to reduce the ± of our report, by specifying <u>y at a given x.</u>

# Analysis of variance (ANOVA) approach to regression:

**Q** How much of the variation in y is due to the (linear) relationship between y and x? ← Table length ... temperature.

**A** variance of y $= \frac{1}{n-1} \sum_{i=1}^{n} (y_i - \bar{y})^2$  [OLS]

[OLS]  $+ \hat{y}_i - \hat{y}_i$ , $\hat{y}_i = \hat{\alpha} + \hat{\beta} x_i$

*E for Error*   Error  *E for Error, not for Explained*

$$\underbrace{S_{yy} = \sum_i (y_i - \bar{y})^2}_{SS_{total}} = \underbrace{\sum_i (\hat{y}_i - \bar{y})^2}_{SS_{explained}} + \underbrace{\sum_i (y_i - \hat{y}_i)^2}_{SS_{unexplained}}$$

total variation in y.

variation in y explained by (or due to) x

variation in y unexplained by x

$$\underset{\sim (10)^2}{SST} = \underset{}{SS_{explained}} + \underset{\sim (3)^2}{SS\textcircled{E}}$$

Variability is reduced from $\pm(10)^2$ to something smaller, say $\pm(3)^2$.

Therefore, $\frac{SS_{exp}}{SST}$ , called $\boxed{R^2}$ , measures how good the fit is. ($\times 100$)

percent variation in y, explained by x.

(Bad Model/fit)  $0 < R^2 < 1$  (Good Model/fit)

The other piece, $SS_{unexpl.} = SSE$, is a sum (of squares), and so can be can be "Averaged" to provide a measure of typical error

$$\sqrt{\frac{SSE}{n-2}} = \underbrace{\sqrt{\frac{1}{n-2} \sum_i^n \underbrace{(y_i - \hat{y}_i)^2}_{error}}}_{\text{"funny Avg."}} \equiv S_e \sim \text{std. dev. of errors}$$
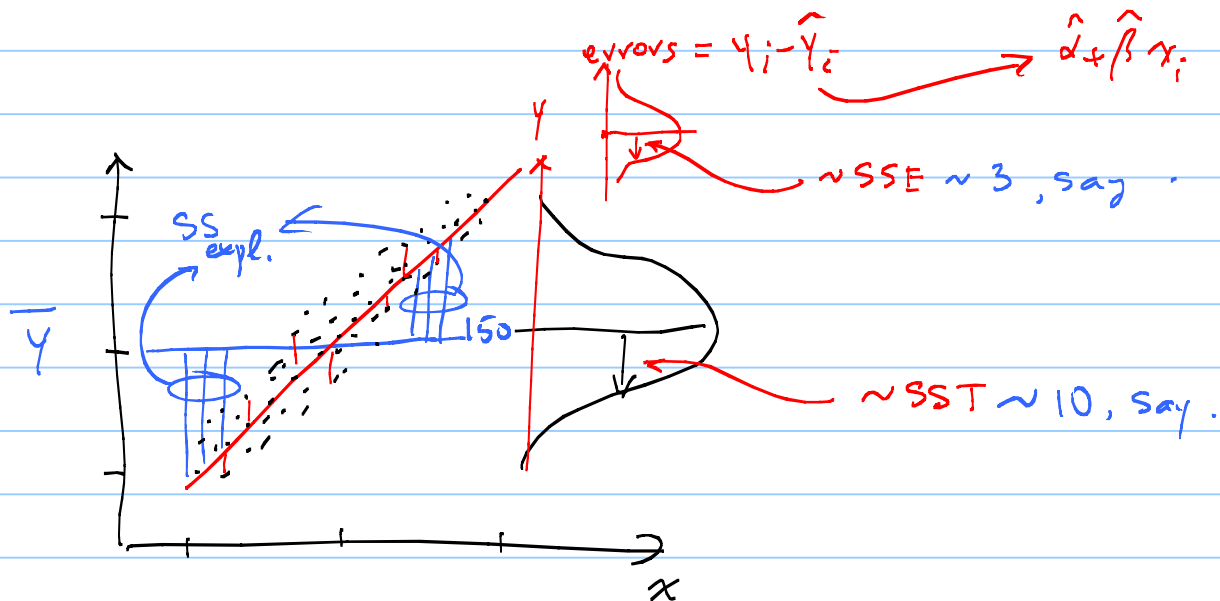
$\sim$ typical error.

Compare with $S_y^2 = \frac{1}{n-1} \sum_i (y_i - \bar{y})^2$
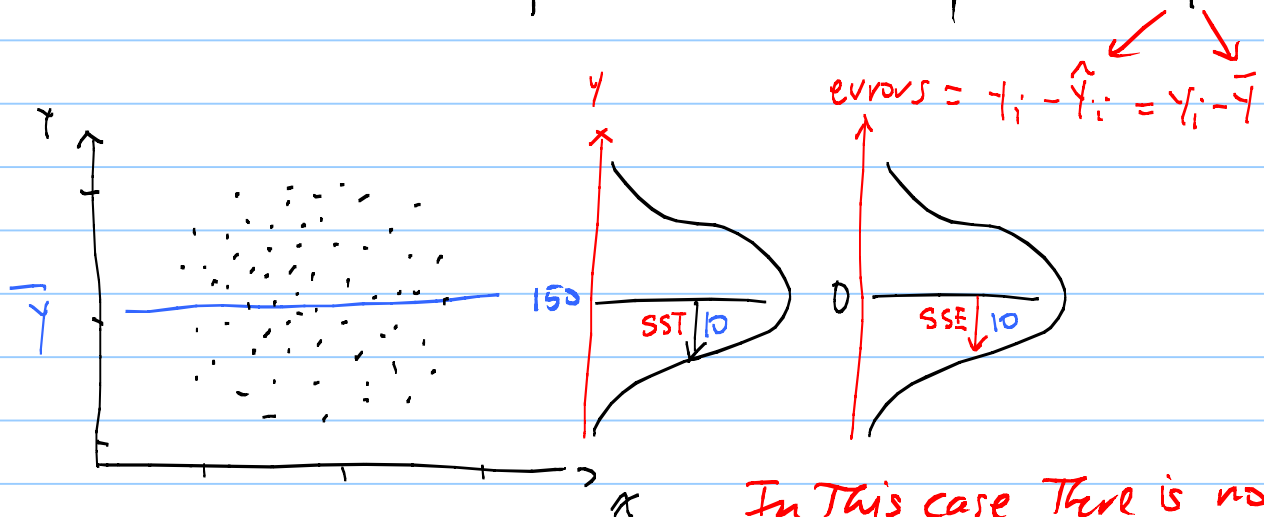
Report $\hat{y}(x) \pm S_e$ .

Picture for The ANOVA decomposition:

errors $= y_i - \hat{Y}_i$ $\longrightarrow \hat{\alpha} + \hat{\beta} x_i$

$\sim SSE \sim 3$, say.

$SS_{expl.}$

150

$\bar{Y}$

$\sim SST \sim 10$, say.

$x$

So, when there is a (linear) relationship between $x$ & $y$, Then some portion of the variation in $y$ can be attributed to (or explained by) $x$. That portion is $SS_{exp.}$, and the (unexplained) rest is $SS_{unexp} = SSE$.

$(10)^2$ $(3)^2$

So The variability in $y$, $SST$, is reduced to $SSE$.

When There is no relationship between $x$ and $y$, Then The fig looks like below. Note That This situation is equivalent to the situation where we have data only on $y$, and NOT on $x$ at all. In That case The best prediction for every case is $\bar{Y}$ (see hw):

$y$

errors $= y_i - \hat{Y}_i = y_i - \bar{Y}$

$\bar{Y}$

150

$SST$ 10

0

$SSE$ 10

$x$

In This case There is no reduction in $SST$ at all, as expected.

**Example** (same as in last few lectures):

$$SST = \sum_i (y_i - \bar{y})^2 = \cdots = 6251.2$$

$$SSE = \sum_i (y_i - \hat{y}_i)^2 =$$ *last column in table in prev. lecture.*

$$= (-1.5)^2 + (5.1)^2 + (11.5)^2 + (-30.3)^2 + (15.1)^2 = 1307$$

$$\Rightarrow \quad R^2 = \text{Coef. of det.} = \frac{SST - SSE}{SST} = \frac{6251.2 - 1307}{6251.2} = 0.79.$$

**Conclusion:** 79% of the variability (or variation)
**(Meaning)** in y (weight, or Tablet length) is due to (can be explained by) the linear relation with x (height, or temperature).

The other piece of the decomposition:

$$\Rightarrow \quad S_e = \sqrt{\frac{1307}{5 - 2}} = 20.9 \text{ pounds}$$

**Conclusion:** The typical deviation of the y values (weight / Tablet length)
**(Meaning)** (ie. error or residual) about the fit is about 21 pounds.

Report weight (or Tablet length):  $\hat{y} \pm 20.9$    with $R^2 = 0.79$

or ICP

or ...     $\overset{\text{``}}{-755 + 13.3\,\boxed{x}}$ ← height or FV or ...

**Q1.** In the prev. clicker qz we found that if $y(x) = \beta$, then the OLS estimate of $\beta$ is $\hat{\beta} = \bar{y}$. Then $S_e$ is (proportional to)

A) 0      B) $S_x$      C) $S_y$      D) None of the above.

*Make Sense?*

$$S_e = \sqrt{\frac{SSE}{n-2}} = \sqrt{\frac{\sum(y_i - \hat{y}_i)^2}{n-2}} = \sqrt{\frac{\sum(y_i - \bar{y})^2}{n-2}} = \sqrt{\frac{n-2}{n-1}}\, S_y \quad \Rightarrow \quad \hat{y}(x) \pm S_e = \bar{y} \pm S_y$$

**BTW:** for this $\hat{y}(x) = \hat{\beta}$ example, $R^2 = 1 - \frac{SSE}{SST} = 1 - \frac{SST}{SST} = 0$

We learned from the last 2 clicker questions, That
If there is no $x$ data, Then the OLS prediction $\hat{y}$ is just $\bar{y}$.

I.e.



What are the $R^2$ and $s_e$ ?

(Y) No x.

$\hat{y}_i = \bar{y}$    defn. of $s_y^2$.

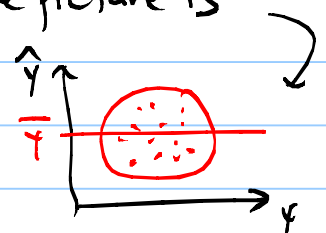$$s_e^2 = \frac{SSE}{n-2} = \frac{\sum_i (y_i - \hat{y}_i)^2}{n-2} = \frac{\sum_i (y_i - \bar{y})^2}{n-2} = \left(\frac{n-1}{n-2}\right) s_y^2 \implies s_e \sim s_y$$

$$R^2 = 1 - \frac{SSE}{SST} = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2} = 1 - \frac{\sum_i (y_i - \bar{y})^2}{\sum_i (y_i - \bar{y})^2} = 0 \implies R^2 = 0.$$
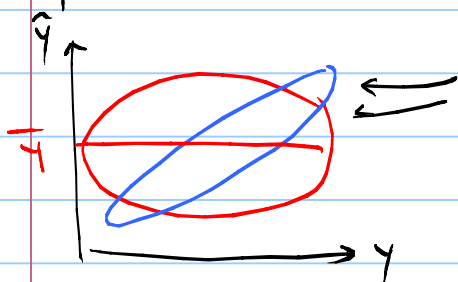
∴ So, if we use $\bar{y}$ as our prediction, Then $R^2 = 0$ (Bad), and
$s_e \sim s_y$, ie. The typical error $\sim$ typical dev. in $y$, ie. nothing gained.

Another situation when nothing is gained is if we make random
predictions, e.g. $\hat{y}_i =$ random. Suppose The mean and The var.
of These random predictions are The same as Those of observations,
ie. $\hat{y}_i =$ random with $\bar{\hat{y}} = \bar{y}$, $s_{\hat{y}} = s_y$. The picture is

But now, something strange happens:

Although one can use The formula for $R^2$ to
arrive at a number, That number does not
have the usual interpretation (ie. percentage of var. in $y$, explained by $x$),
because $\hat{y}_i =$ random are not OLS predictions. So, we don't have
The ANOVA decomposition at all. Same objection applies to $s_e$.
Again, The ANOVA decomposition is correct only for OLS $\hat{y}$;
$\hat{y} =$ random are not OLS predictions.



These both have equal/comparable $s_{\hat{y}}$ (ie. $R^2$).
But The blue one has lower $s_e$.
This doesn't contradict ANOVA, because
The red $\hat{y}$ is not OLS.

In short, both have equal precision, but blue is more accurate

For the data shown here:
x= 45, 58, 71, 71, 85, 98, 108
y = 3.20, 3.40, 3.47, 3.55, 3.60, 3.70, 3.80
a) Compute the eq. of the OLS fit.
b) Compute the total variation, SST.
c) Decompose it into explained and unexplained.
d) Compute R2 and interpret it (in English),
e) Compute the std. dev of errors, s_e, and interpret it (in English).

All by hand. You may use R to compute sums, means, std. deviations, but not a function that does regression or analysis of variance.

**hw-lect11-2** Consider the following decomposition:

$$\sum_i (y_i - \bar{y})^2 = \sum_i [(\hat{y}_i - \bar{y}) + (y_i - \hat{y}_i)]^2 \qquad \boxed{ByR}$$

$$= \sum_i (\hat{y}_i - \bar{y})^2 + \sum_i (y_i - \hat{y}_i)^2 + 2 \sum_i (\hat{y}_i - \bar{y})(y_i - \hat{y}_i)$$

In past hws I have asked students to prove that the last term is zero if $\hat{y}_i = \hat{\alpha} + \hat{\beta} x_i$, with $\hat{\alpha}, \hat{\beta}$ being the OLS estimates (ie. $\hat{\alpha}, \hat{\beta}$ given in lects, book). Unfortunately, it's a long calculation; so this time we'll try to show that it's zero using simulation in R. Write code to

a) generate a sample of size 100 from the unif dist. between -1 and +1. Call it $x$.

b) generate y such that $y_i = 2 + 3 x_i + \epsilon_i$ with $\epsilon_i$ having a normal distr. with $\mu = 0, \sigma = 0.5$.

c) Do regression on $x, y$, and call the predictions $\hat{y}$.

d) compute $\sum_i (\hat{y}_i - \bar{y})(y_i - \hat{y}_i)$. It should be (very) zero!

$SS_{exp}$ can be computed from its defining relation $\left(\sum_i (\hat{y}_i - \bar{y})^2\right)$ Or from $(SST - SSE)$, or from $\hat{\beta}$ and $S_{xx}$, as follows.

Explain what has happened at every step.

$$SS_{exp} = \sum_i (\hat{y}_i - \bar{y})^2$$

?

$$= \sum_i (\hat{\alpha} + \hat{\beta} x_i - \bar{y})^2$$

?

$$= \sum_i (\bar{y} - \hat{\beta}\bar{x} + \hat{\beta} x_i - \bar{y})^2$$

?

$$= \sum_i (\hat{\beta})^2 (x_i - \bar{x})^2$$

?

$$= (\hat{\beta})^2 \sum_i (x_i - \bar{x})^2$$

?

$$= (\hat{\beta})^2 S_{xx}$$

⇑

If you would like, print out This page, and write your answers in The space here.