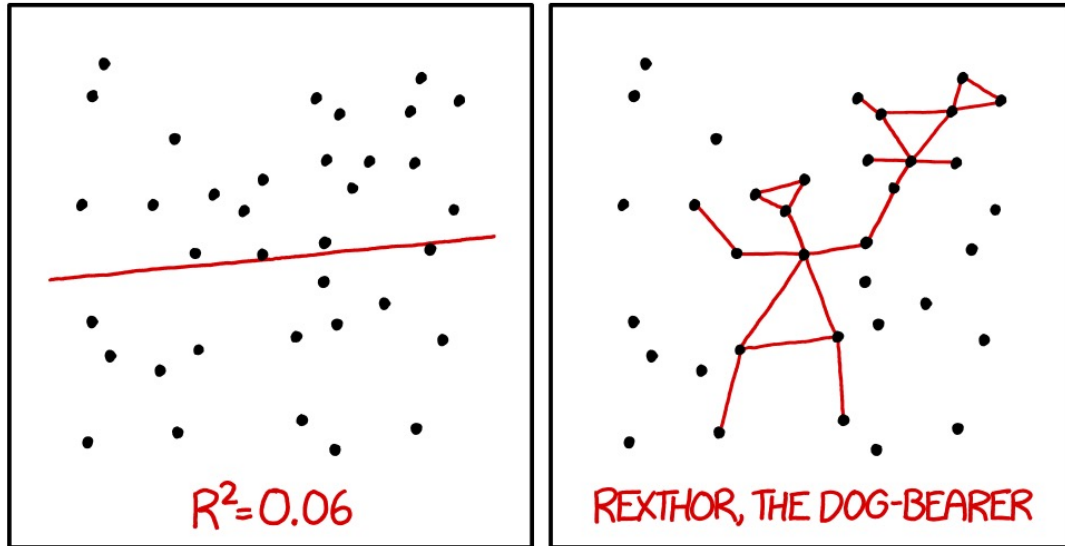


# Stat 88: Probability & Mathematical Statistics in Data Science



I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER  
TO GUESS THE DIRECTION OF THE CORRELATION FROM THE  
SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

Lecture 40 : 4/30/2021

Chapter 12

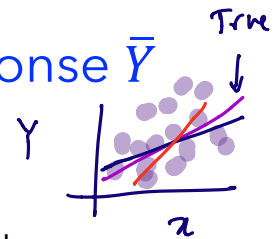
Finishing up regression

<https://xkcd.com/1725/>

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \text{ (True reg. line)} : \hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

The individual response  $Y_i$  and the average response  $\bar{Y}$

$$\varepsilon_i \sim N(0, \sigma^2) \quad Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$



- $Y_i$  are normal with expectation  $\beta_0 + \beta_1 x_i$  and variance  $\sigma^2$
- Note that the individual responses are independent of each other.

- Let  $\bar{Y}$  be the average response. 
$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$
- $E(\bar{Y}) = \beta_0 + \beta_1 \bar{x}$  (the expected average response is the signal at the average value of the predictor variable)

- $Var(\bar{Y}) = \frac{\sigma^2}{n}$  (only involves the error variance since the randomness in the  $Y_i$ 's comes only from the errors or noise)

- Since  $\bar{Y}$  is a linear combination of independent normally distributed random variables, it is also normal.

$$\bar{Y} \sim N(\beta_0 + \beta_1 \bar{x}, \frac{\sigma^2}{n})$$

from Ch 11,  $\hat{\alpha} = \frac{E(D_x D_y)}{\sigma_x^2}$

The estimated slope  $\hat{\beta}_1$  (estimates the true slope  $\beta_1$ , which is an unobservable parameter)

• The least squares estimate of the true slope  $\beta_1$  is  $\hat{\beta}_1 = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$  } sample cov of  $x, Y$

Let's define  $a_i = (x_i - \bar{x})$  }  
 $\sum_{i=1}^n (x_i - \bar{x}) = \sum_{i=1}^n a_i \hookrightarrow A$  }  
 $\hat{\beta}_1 = \frac{\sum_{i=1}^n a_i (Y_i - \bar{Y})}{A} = \sum_{i=1}^n \frac{a_i}{A} (Y_i - \bar{Y})$  } sample variance

- Notice that  $\hat{\beta}_1$  is random (because of the  $Y_i$ ).
- Also, since  $Y_i$  is normal, and  $\bar{Y}$  is normal, so is  $Y_i - \bar{Y}$ , therefore  $\hat{\beta}_1$  is also normally distributed

$$E(\hat{\beta}_1) = \beta_1$$

- $E(Y_i - \bar{Y}) = \beta_1(x_i - \bar{x})$
- $E(\hat{\beta}_1) = \beta_1$ , so  $\hat{\beta}_1$  is an unbiased estimator of  $\beta_1$
- $Var(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$  (to be taken as fact, proof beyond the scope of this class)

## SD of the estimated slope $\hat{\beta}_1$

- $SD(\hat{\beta}_1) = \frac{\sigma}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$ 

$\sigma$  ← constant unknown

→ we know these

- Need to estimate  $\sigma$ , which we will do by using the SD of the residuals. The larger the  $n$ , the better our estimate of  $\sigma$

$$\hat{\sigma} = SD(D_1, D_2, \dots, D_n) = \sqrt{\frac{1}{n} \sum_{i=1}^n (D_i - \bar{D})^2}, \quad \leftarrow$$

- $D_i = Y_i - \hat{Y}_i = Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$  (The  $D_i$  are the residuals and estimate the errors)  
residuals
- Since we are estimating the SD from the data, we will call it the **standard error** of the estimator.
- That is, we will denote this estimated  $SD(\hat{\beta}_1)$  by  **$SE(\hat{\beta}_1)$** .

$$\begin{array}{l} x \quad Y \\ (1, 3) \\ (0, 2) \\ (-1, 2) \end{array} \left\{ \begin{array}{l} \hat{\beta}_1 \\ \hat{\beta}_0 \end{array} \right\} \begin{array}{l} \text{compute } \hat{Y}_i \\ Y_i - \hat{Y}_i \end{array}$$

$$SE(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

$\hat{\sigma}$  ← estimate of  $\sigma$  which estimated by the sample SD of residuals

$D_i: Y_i - \hat{Y}_i = \text{residual}$

## Confidence intervals for $\beta_1$

- $SE(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$ , for large  $n$ ,  $SE(\hat{\beta}_1) \rightarrow SD(\hat{\beta}_1)$

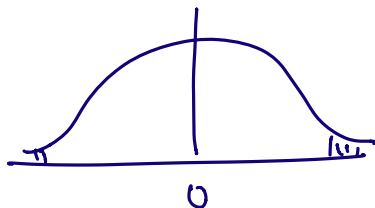
Therefore, for large  $n$ , the distribution of  $\hat{\beta}_1$ , standardized, is approximately standard normal.

$$T = \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} \sim N(0,1) \quad \leftarrow \text{if } n \text{ not large, then } T \not\sim N(0,1)$$

- A 95% CI for  $\beta_1$  is given by  $\hat{\beta}_1 \pm 2SE(\hat{\beta}_1)$  ( , )

If 95% CI does not contain 0, then can reject  $H_0: \beta_1 = 0$  (vs.  $\beta_1 \neq 0$ ) at 5% significance level.

- Note that if the sample size is not large enough, the distribution of  $T$  is not necessarily normal, since the assumption that  $SE(\hat{\beta}_1) \approx SD(\hat{\beta}_1)$  may not hold.
- In this situation, we model the distribution of  $T$  using a family of bell-shaped distributions, called the *t-distributions*.



## Hypothesis tests to test $\beta_1 = 0$

- $\beta_1 = 0$  is a very important question: is there any linear relationship at all?
- A 95% CI for  $\beta_1$  is given by  $\hat{\beta}_1 \pm 2SE(\hat{\beta}_1)$ : we can use this CI. If 0 is not in this interval, then we reject the null hypothesis of the slope being 0 at the 5% significance level.
- We can set up a test:  $H_0: \beta_1 = 0$  vs  $H_1: \beta_1 \neq 0$  and use the fact that under the null hypothesis,

$$T = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)} \sim N(0,1)$$

$$H_0: \beta_1 = 0$$
$$H_1: \beta_1 \neq 0$$

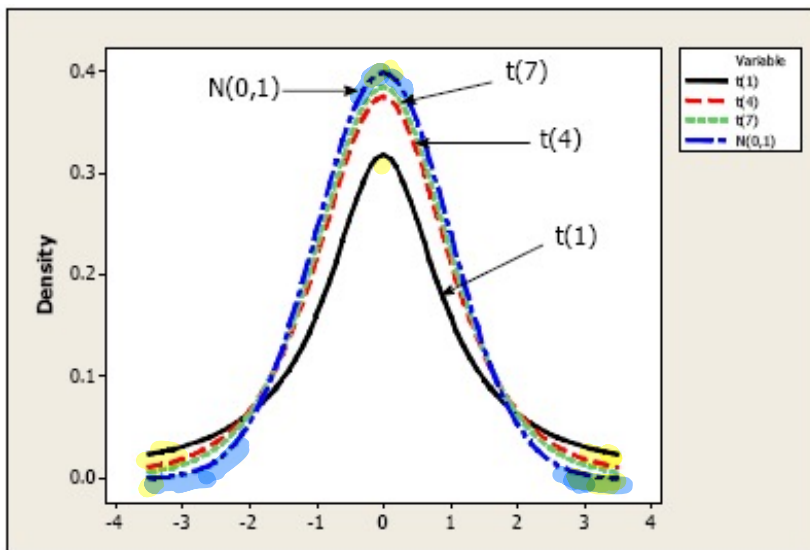
- Let's look at the example from the text on pulse rates after looking at the t-distribution

# The $t$ -distribution

$t$ -dsn has parameter "degrees of freedom"

Rather than a normal curve, a  $t$ -curve is used. For regression, "degrees of freedom" for  $T$  equals  $n - 2$ . For large enough  $n$ , use the normal curve.

(When the sample size  $n$  is large, so is  $n - 2$ , so we might as well use the normal curve. When the sample size is small, using the appropriate  $t$  curve gives more accurate answers.)



$$T = \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)}$$

## Example (12.4.3)

slope, intercept, r, p, se\_slope=

$\hat{\beta}_1$  (1.142879681904831,  
 $\hat{\beta}_0$  13.182572776013345,  
 $r$  0.6041870881060092,  
 $p\text{-value}$  1.7861044071652305e-24, ( $H_0: \beta_1 = 0$  vs  $H_1: \beta_1 \neq 0$ )  
 $SE(\hat{\beta}_1)$  0.09938884436389145)

mean\_active, sd\_active

(91.29741379310344, 18.779629284683832)  
 $\bar{y}$ ,  $SD(y_1, \dots, y_n)$

c) Find the SD of the residuals.

$\hat{\sigma}^{??}$

$$SD(\text{residuals}) \approx \sqrt{1-r^2} SD(y)$$

$$\approx 14.96$$

mean\_resting, sd\_resting

(68.34913793103448, 9.927912546587986)  
 $\bar{x}$ ,  $SD(x_1, \dots, x_n)$   
 $SD(x_i)$

$$SE(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

$$\hat{\sigma} = SE(\hat{\beta}_1) \cdot \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$= SE(\hat{\beta}_1) \cdot \sqrt{n \cdot (SD(\text{resting}))^2}$$



$$232 \cdot (\hat{\sigma}_x)^2$$

Quick look back