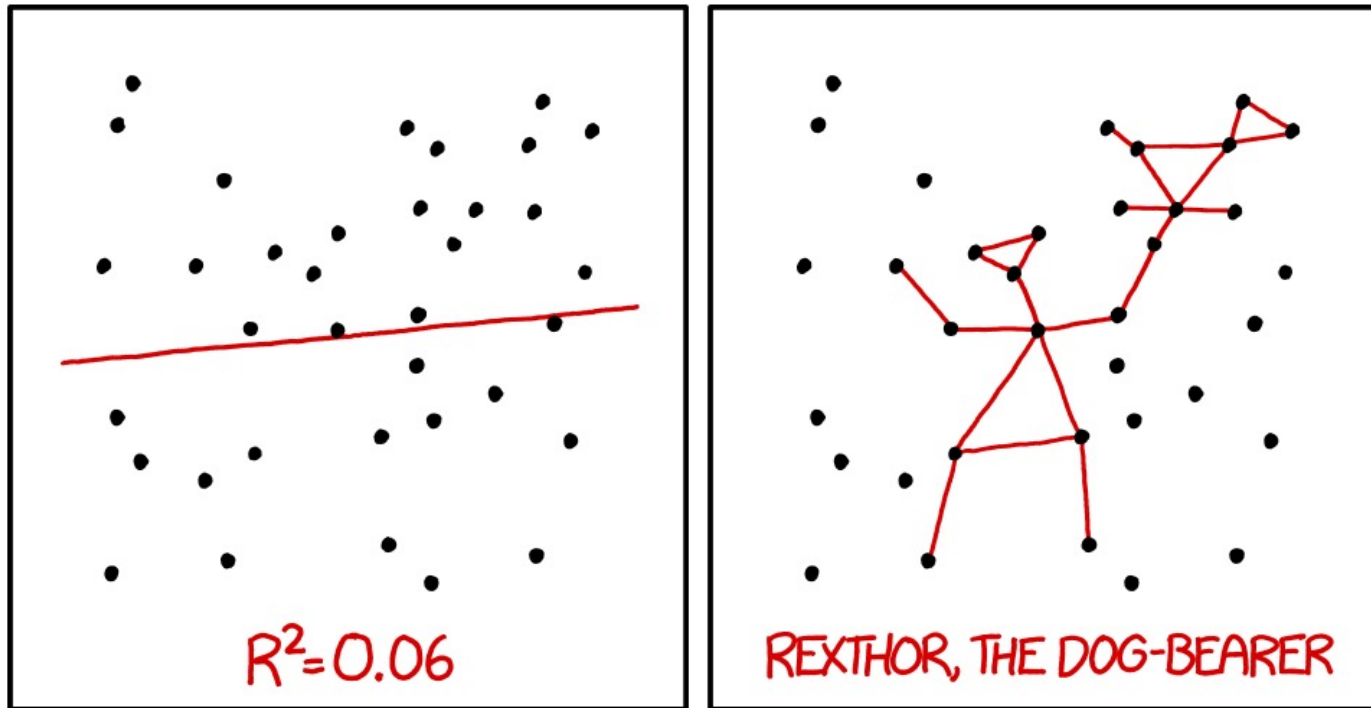


# Stat 88: Probability & Mathematical Statistics in Data Science



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 39 : 4/28/2021

Chapter 12

More about regression

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

So far:

- $\hat{Y} = \hat{a}X + \hat{b}$  - line of "best fit" : minimizes mean squared error =  $E[(Y - \hat{Y})^2]$
- $\hat{Y}$  is called the fitted value of  $Y$ , where:  $\hat{a} = \frac{r\sigma_Y}{\sigma_X}$ ,  $\hat{b} = \mu_Y - \hat{a}\mu_X$
- $\hat{Y} = \hat{a}X + \hat{b} = \hat{a}X + \mu_Y - \hat{a}\mu_X = \hat{a}(X - \mu_X) + \mu_Y = \hat{a}D_X + \mu_Y$
- Correlation:  $r = E\left[\left(\frac{X - \mu_X}{\sigma_X}\right)\left(\frac{Y - \mu_Y}{\sigma_Y}\right)\right] = E(Z_X Z_Y)$  and  $-1 \leq r \leq 1$
- Residual  $D = Y - \hat{Y}$ ,  $E(D) = 0$ ,  $Var(D) = (1 - r^2)\sigma_Y^2$
- $r(D, X)$  (the residuals are uncorrelated with the predictor: why?)

# The Simple Linear Regression Model

- Regression model from data 8
- Model has two variables: response ( $Y$ ) & ( $x$ ) predictor/covariate/feature variable
- **Assumptions:** response is a linear function of the predictor (signal) + random error (noise), where the noise has a **normal** distribution, centered at 0. The signal is not random, but the response is, because the noise is random:

$$\text{response} = \text{signal} + \text{noise}$$

- In mathematical language:

# The regression line

- For each  $i$ , we want to get as close as we can to the *signal*  $\beta_0 + \beta_1 x_i$
- There is some “true” regression line  $\beta_0 + \beta_1 x$  that we cannot observe since there is noise. We estimate this line by minimizing the squared observed error.
- Estimate of the line given the data is  $Y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ , where  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are the estimates of the intercept and slope, respectively, given the data.
- We will investigate the distribution of the slope estimate (why is it random?) after looking at the individual and average response.

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

- For any fixed  $i$ ,  $Y_i$  is the sum of the signal and the noise.
- The signal is not random, but the noise is random with  $\epsilon_i \sim N(0, \sigma^2)$
- Therefore what is the distribution of the  $Y_i$  ?
- What can you say about the independence and distribution of each of the  $Y_i$  ? Are they iid?
- Let  $\bar{Y}$  be the average response. What would be its distribution?
- $E(\bar{Y}) =$
- $Var(\bar{Y}) =$

## The estimated slope $\beta_1$

- Recall the slope we derived in the previous chapter

$$\hat{a} = \frac{E(D_X D_Y)}{\sigma_X^2}$$

- Now we have data, so we need to use the empirical distribution
- The least squares estimate of the true slope  $\beta_1$  is:

$$\hat{\beta}_1 =$$

- Notice that  $\hat{\beta}_1$  is random (because of the  $Y_i$ ). How would we find its distribution?
- Note that  $E(Y_i - \bar{Y}) = \beta_1(x_i - \bar{x})$
- $E(\hat{\beta}_1) =$

## Distribution of $\hat{\beta}_1$

- From the formula of  $\hat{\beta}_1$ , we see that it is a linear combination of the independent normal rvs  $Y_1, Y_2, \dots, Y_n$  and therefore  $\hat{\beta}_1$  is also normal.
- $E(\hat{\beta}_1) = \beta_1$  indicating that  $\hat{\beta}_1$  is an \_\_\_\_\_ estimator of  $\beta_1$
- Recall that the common variance of the errors  $\epsilon_i$  is  $\sigma^2$
- FACT:  $Var(\hat{\beta}) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$
- What you want to note is that the numerator is constant, so as we have more terms, the denominator gets larger, and our estimated slope gets closer to the true slope.
- We will need to estimate  $\sigma^2$  since it is an unknown parameter.

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

- $SD(\hat{\beta}_1) =$
- Need to estimate  $\sigma$ , which we will do by using the SD of the residuals. Since we are estimating the SD from the data, we will call it *standard error* of the estimator.
- That is, we will denote this estimated  $SD(\hat{\beta}_1)$  by  **$SE(\hat{\beta}_1)$** .
- The larger the  $n$ , the better our estimate of  $\sigma$

$$\hat{\sigma} = SD(D_1, D_2, \dots, D_n) = \sqrt{\frac{1}{n}}$$

- A 95% CI for  $\beta_1$  is given by  $\hat{\beta}_1 \pm 2SE(\hat{\beta}_1)$
- For large  $n$ , the distribution of  $\hat{\beta}_1$ , standardized, is approximately standard normal.

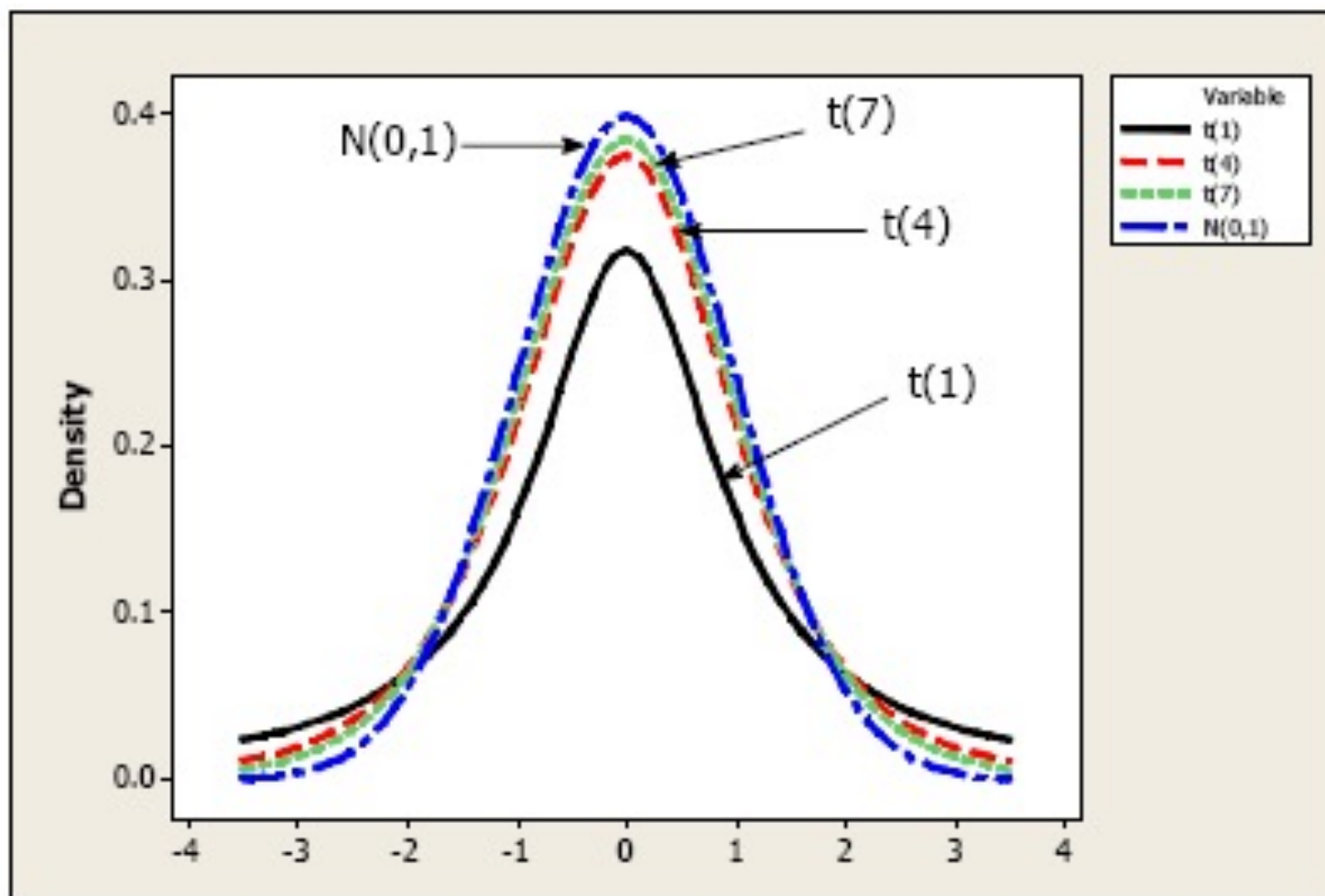
$$T = \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} \sim N(0,1)$$

- Let's look at the example from the text on pulse rates.



# The $t$ -distribution

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.



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