

# CSCI 3022: Intro to Data Science

## Lecture 21: Inference in Regression

---

Rachel Cox

Department of Computer  
Science

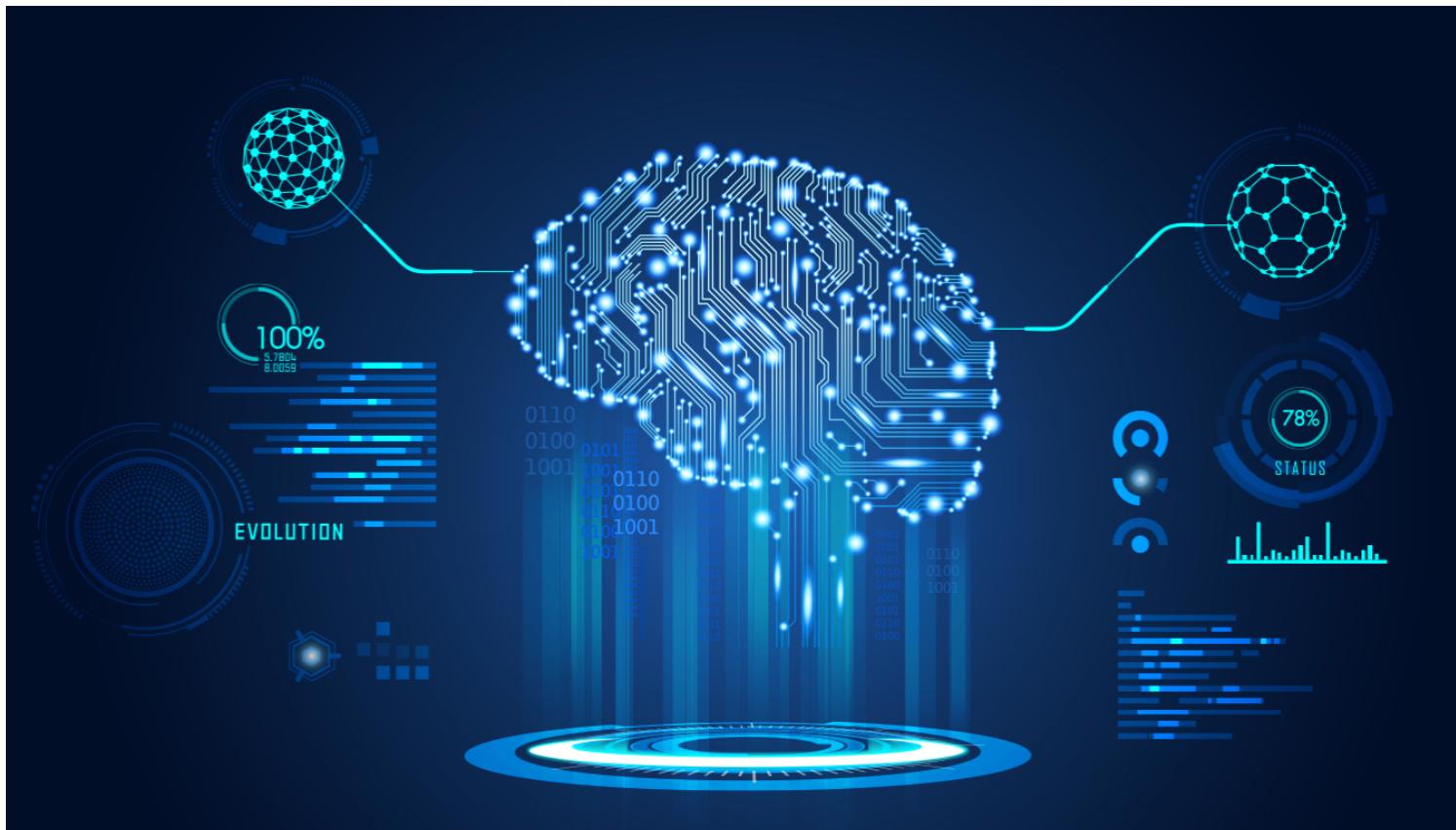
## Announcements & Reminders

---

HWS - Due April 17<sup>th</sup> at 11:59pm

# What will we learn today?

- Estimate the variance in data about the true regression line
- Quantify the goodness-of-fit of our SLR model
- Perform inference on the regression parameters
  
- *A Modern Introduction to Probability and Statistics, Chapter 27.3*



# Simple Linear Regression (SLR)

---

Given data  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ , fit a simple linear regression of the form

$$y_i = \alpha + \beta x_i + \epsilon_i \quad \text{where } \epsilon_i \sim N(0, \sigma^2)$$

Estimates of the intercept and slope parameters are given by minimizing

$$SSE = \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2 \quad \left. \right\} \begin{matrix} SSE - \text{sum of squared} \\ \text{errors} \end{matrix}$$

The least-squares estimates of the parameters are:

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x} \quad \text{and} \quad \hat{\beta} = \frac{\bar{x} \bar{y} - \bar{x} \bar{y}}{(\bar{x})^2 - \bar{x}^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

from our  
book

# Simple Linear Regression (SLR)

---

$$\begin{aligned}\bullet \quad \hat{\beta} &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \leftarrow \text{from book} \\ &= \frac{\sum_{i=1}^n (x_i y_i - x_i \bar{y} - \bar{x} y_i + \bar{x} \bar{y})}{\sum_{i=1}^n (x_i^2 - 2x_i \bar{x} + (\bar{x})^2)} \\ &= \frac{n \bar{xy} - n \bar{x} \bar{y} - n \bar{x} \bar{y} + n \bar{x} \bar{y}}{n \bar{x^2} - 2n(\bar{x})^2 + n (\bar{x})^2} \\ &= \frac{n \bar{xy} - n \bar{x} \bar{y}}{n \bar{x^2} - n (\bar{x})^2} \\ &= \frac{\bar{xy} - \bar{x} \bar{y}}{\bar{x^2} - (\bar{x})^2} \quad \leftarrow \begin{array}{l} \text{we derived} \\ \text{this in the} \\ \text{last lecture} \end{array}\end{aligned}$$

# Estimating the variance

The parameter  $\sigma^2$  determines the spread of the data about the true regression line.

An estimate of  $\sigma^2$  will be used in computing confidence intervals and doing hypothesis testing on the estimated regression parameters.

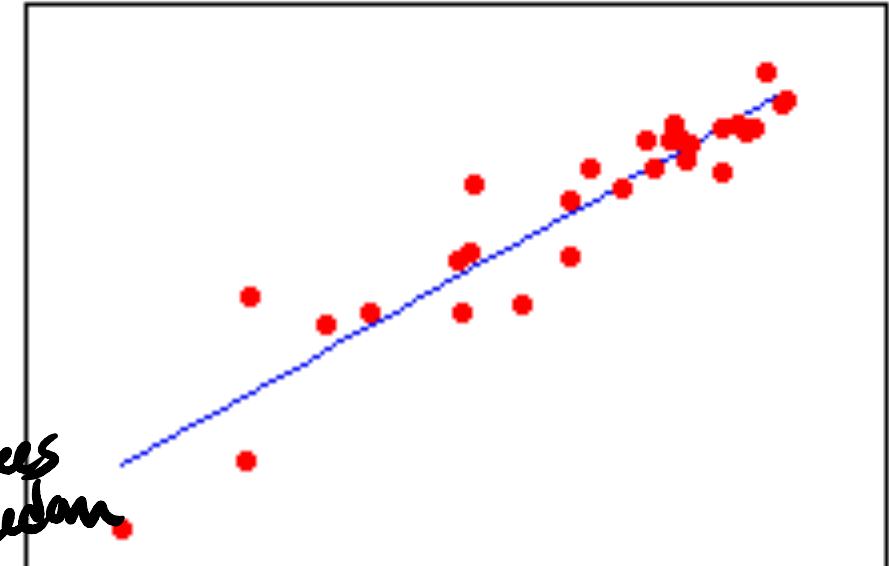
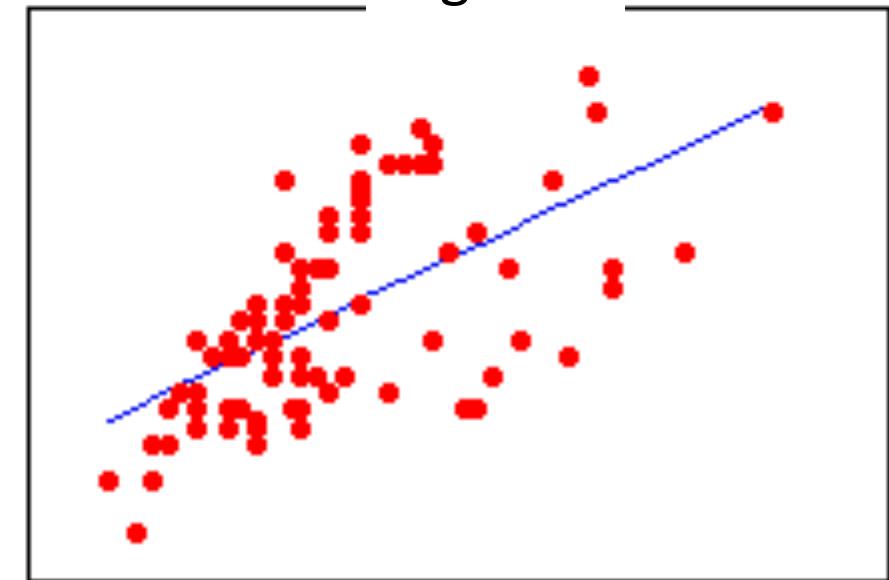
Recall that the sum of squared errors is given by:

$$\cdot SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

*observed data* ↑    ↑ *model*    Because  $\hat{y} = \hat{\alpha} + \hat{\beta} x$

Our estimate or the variance  $\hat{\sigma}^2$  is given by:  $\hat{\sigma}^2 = \frac{SSE}{n-2}$  :

↑  
 $n-2$  degrees  
of freedom



# The Coefficient of Determination - $R^2$

The coefficient of determination,  $R^2$ , quantifies how well the model explains the data.

- $R^2$  is a value between 0 and 1

The **sum of squared errors (SSE)** can be thought of as a measure of how much variation in Y is left unexplained by the model.

$\downarrow$  <sup>model</sup>    $\checkmark$  <sup>mean</sup>

The **regression sum of squares** is given by:  $SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$

It gives a sense of how much variation in Y is explained by our model.

A quantitative measure of the total amount of variation in observed Y values is given by the **total sum of squares**:  $SST = \sum_{i=1}^n (y_i - \bar{y})^2$

$\nwarrow$  <sup>data</sup>    $\swarrow$  <sup>mean</sup>

SST is what we would get for SSE if we used the mean of the data as our model.

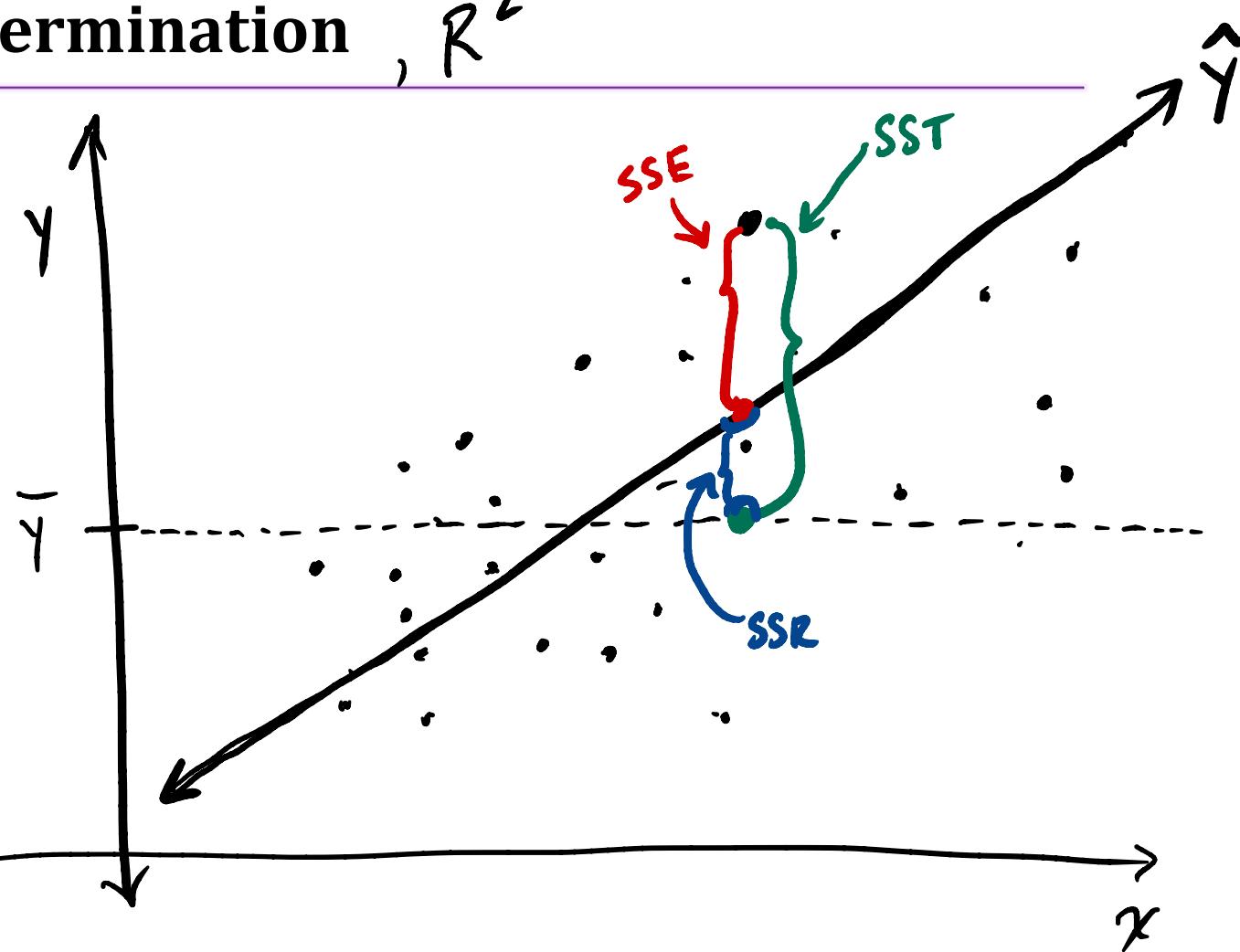
# The Coefficient of Determination , $R^2$

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

$R^2$  is a way to  
quantify "goodness of fit"



$$R^2 = \frac{\text{Variance explained by model}}{\text{Total Variance}} = \frac{SST - SSE}{SST} = 1 - \frac{SSE}{SST}$$

\* In general, the larger  $R^2$  is, the better the regression model fits the data.

# The Coefficient of Determination

---

The sum of squared deviations about the least-squares line is smaller than the sum of squared deviations about any other line.

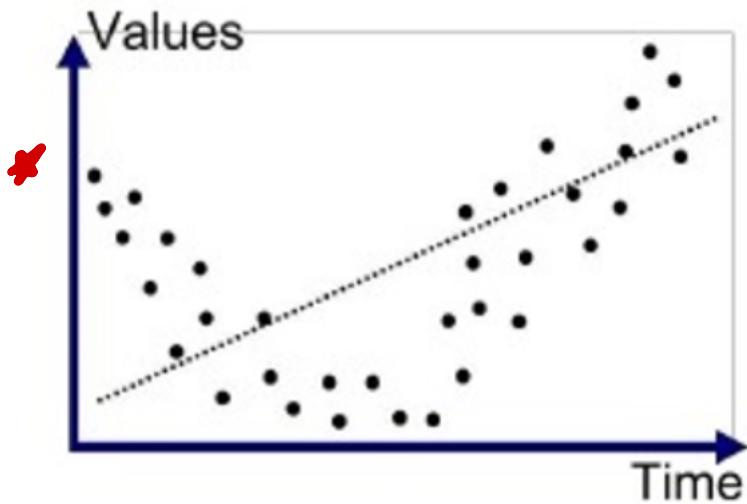
$$SSE \leq SST$$

The ratio  $SSE/SST$  is the proportion of total variation in the data (SST) that cannot be explained by the SLR model (SSE). So we define the coefficient of determination  $R^2$  to be the proportion that can be explained by the model.

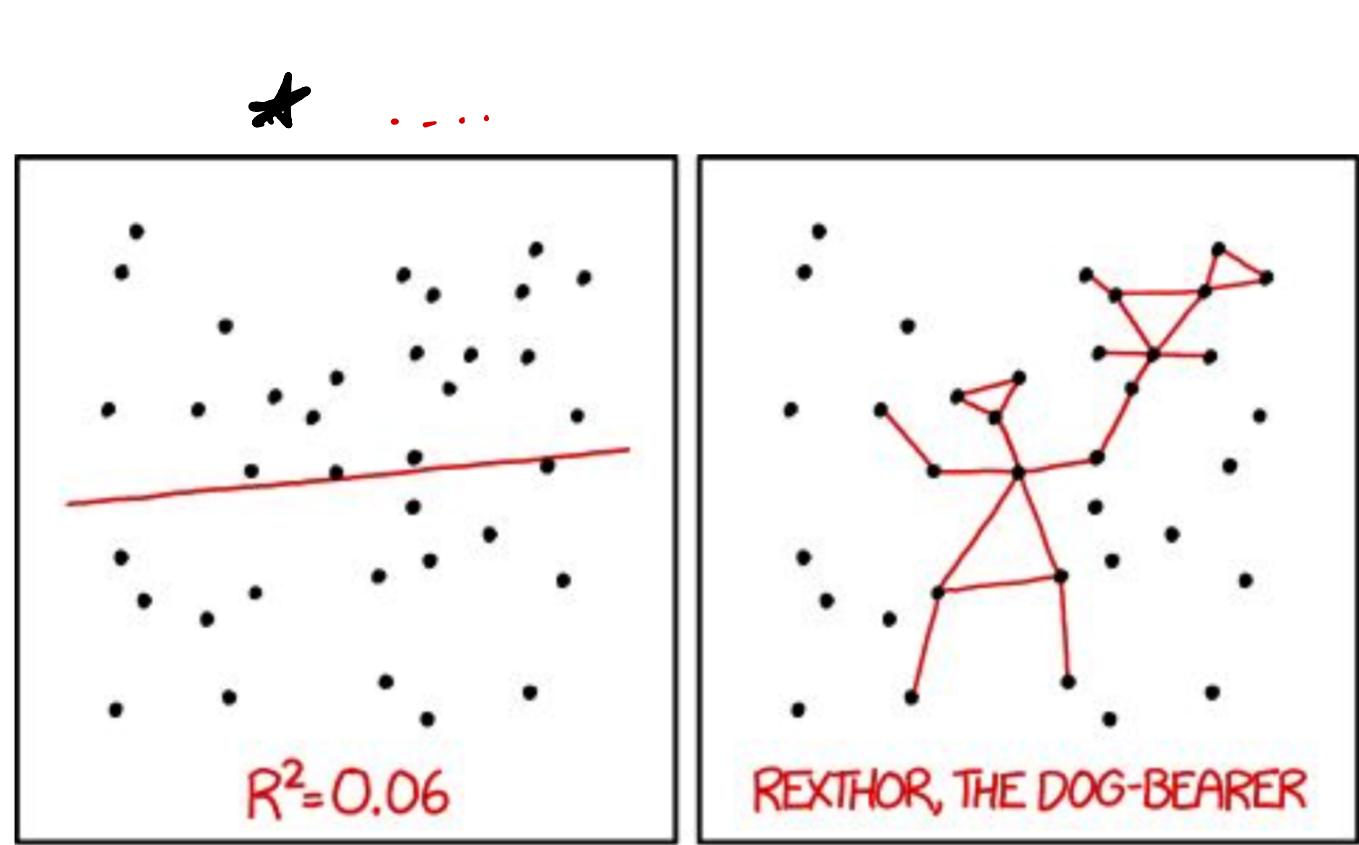
$$R^2 = 1 - \frac{SSE}{SST}$$

# The Coefficient of Determination

$R^2$  is the proportion of total variation in the data that is explained by the model. It does not tell you whether you have the correct model.



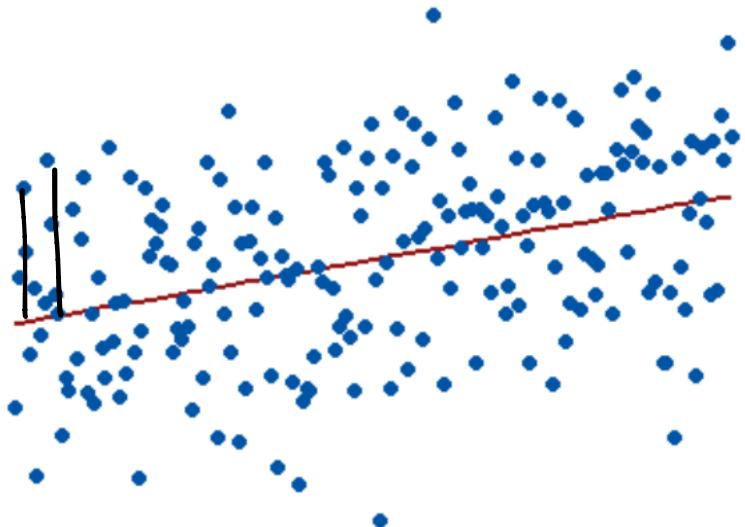
Maybe we should have used a quadratic model for this data



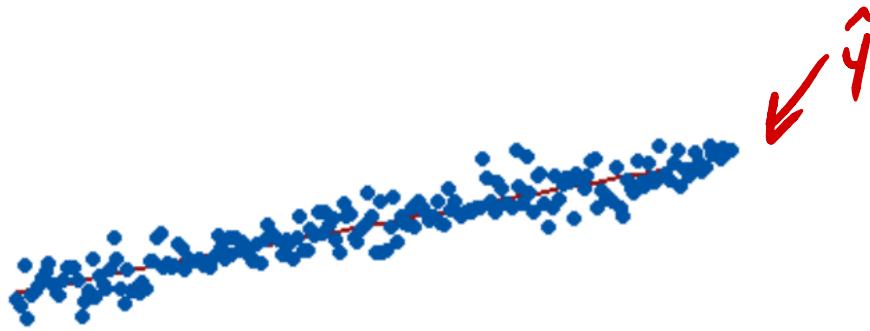
I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

# The Coefficient of Determination

---



$$R^2 = .15$$



$$R^2 = .85$$

$R^2$  is a number between 0 and 1

Another way to examine goodness of fit is to examine a plot of the residuals.

# Inference about SLR parameters

---

$\downarrow \hat{\alpha}, \hat{\beta}$

The parameters in the simple linear regression model have distributions.

From these distributions, we can construct CIs for the parameters, conduct hypothesis tests, etc.

We will focus mainly on the slope parameter  $\beta$  - Is there really a relationship between the feature and the response?

X

Y

The distribution for the estimate of the slope is given by:

$$\hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right) \rightarrow$$


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

# Inference about SLR parameters

- ❖ The confidence intervals for  $\beta$  are given by:

$$\hat{\beta} \pm t_{\alpha/2, df=n-2} \cdot \frac{SE(\hat{\beta})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

$SE(\hat{\beta})$   
 $\hat{\sigma}$

*t-distribution*

- ❖ Hypothesis Testing:

$$H_0: \beta = c$$

$$H_1: \beta \neq c$$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - (\hat{\alpha} + \hat{\beta}x_i))^2}{n-2}$$

See notebook 21  
for examples!

Test statistic:  $t = \frac{\hat{\beta} - c}{SE(\hat{\beta})}$  → Compare to  $t_{\alpha/2, n-2}$  or compute p-value

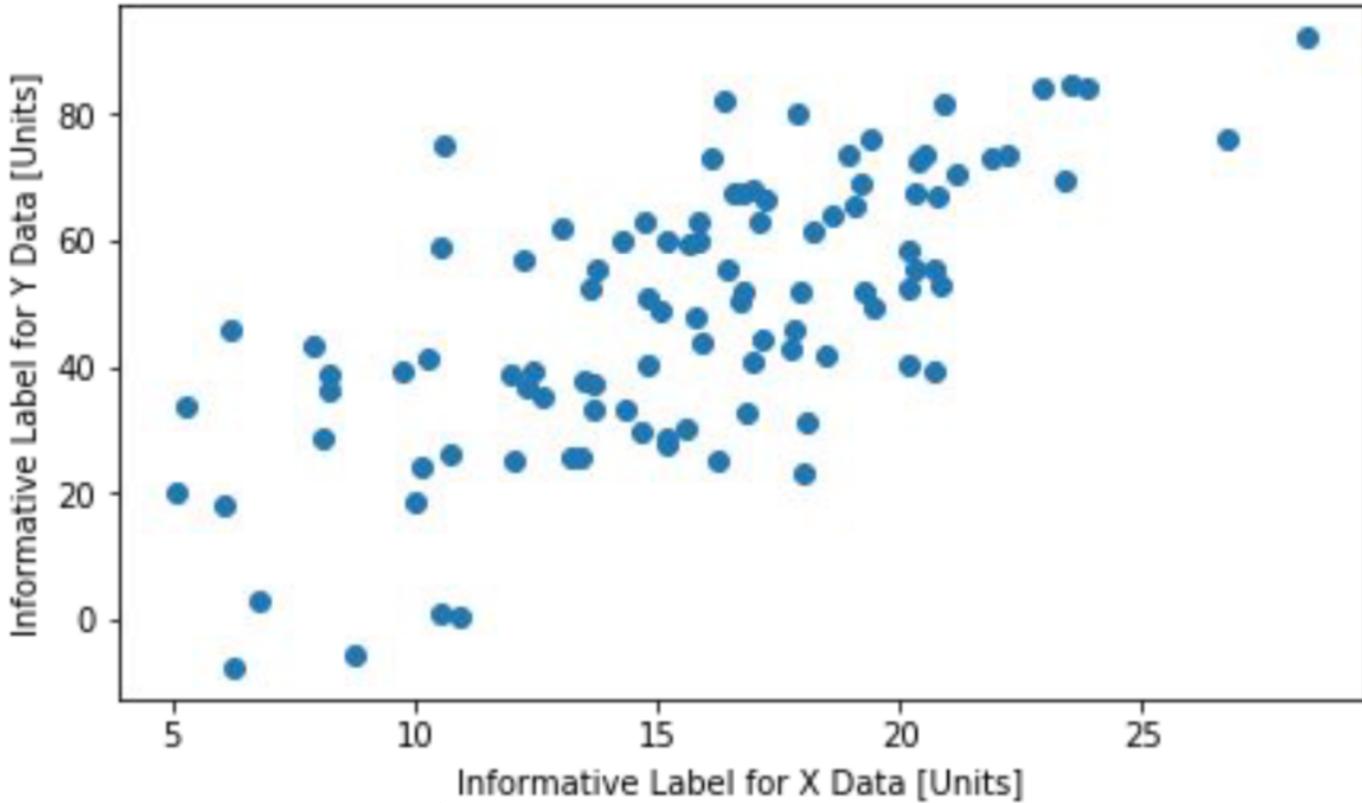
# Inference about SLR parameters

---

Given data ( $x, y$ )

- Explore/plot/histograms
- Formulate null and alternative hypotheses
- Fit a regression line
- Compute CI or p-value/rejection region and test statistic for testing your hypotheses.
- Make conclusions

*Response*



*Feature*

*Next Time:*

❖ Notebook Day