

# Module 3: Linear Regression

Lecture 4  
Jan 19th, 2023



**SPARTANS WILL.**

# Recap

Simple Linear Regression  $Y = f(X) + \epsilon = \beta_0 + \beta_1 X + \epsilon$

Residual sum of squares (RSS)

Confidence interval.

$\hat{\beta}_1 \pm 2 \cdot \text{SE}(\hat{\beta}_1)$

$(X) f(x) = E(Y | X=x)$

$Y = \beta_0 + \beta_1 X + \epsilon$

$E(\beta_0 + \beta_1 x + \epsilon) = \beta_0 + \beta_1 x + E(\epsilon)$

$= \beta_0 + \beta_1 x$

$$\text{SE}(\hat{\beta}_1)^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad \text{SE}(\hat{\beta}_0)^2 = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

where  $\sigma^2 = \text{Var}(\epsilon)$

$$\hat{\sigma}^2 = \frac{\text{RSS}}{n-2}$$

$\beta_1$

$\hat{\beta}_1$

2. Sigma

$\hat{\beta}_0, \hat{\beta}_1, \hat{\sigma}^2$



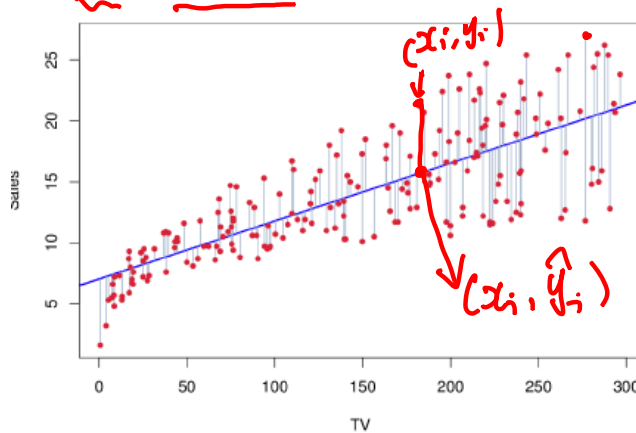
- Confidence interval, hypothesis test, and p-value for coefficient estimates
- Residual standard error (RSE)
- R squared
- Setup for multiple linear regression

# Set up

$$Y = f(X) + \epsilon = \beta_0 + \beta_1 X + \epsilon \quad \Rightarrow$$

- Given  $(x_1, y_1), \dots, (x_n, y_n)$   $\leftarrow$
- Let  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$  be prediction for  $Y$  on  $i$ th value of  $X$ .
- $\underline{e_i} = \underline{y_i} - \underline{\hat{y}_i}$  is the  $i$ th residual

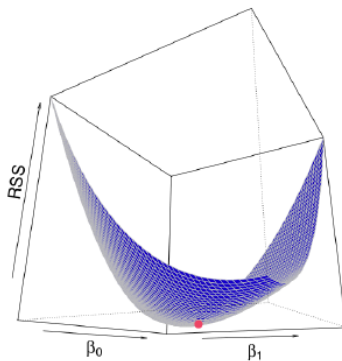
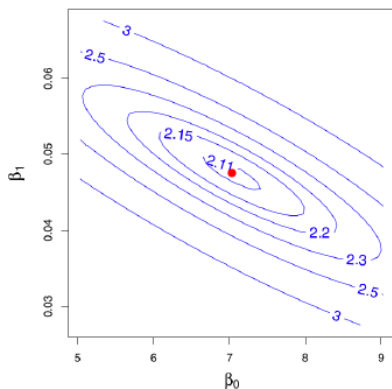
metric ?



# Ordinary Least Squares Regression

MICHIGAN STATE  
UNIVERSITY

OLS



Residual sum of squares RSS is

$$RSS = e_1^2 + \dots + e_n^2$$
$$T = \sum_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

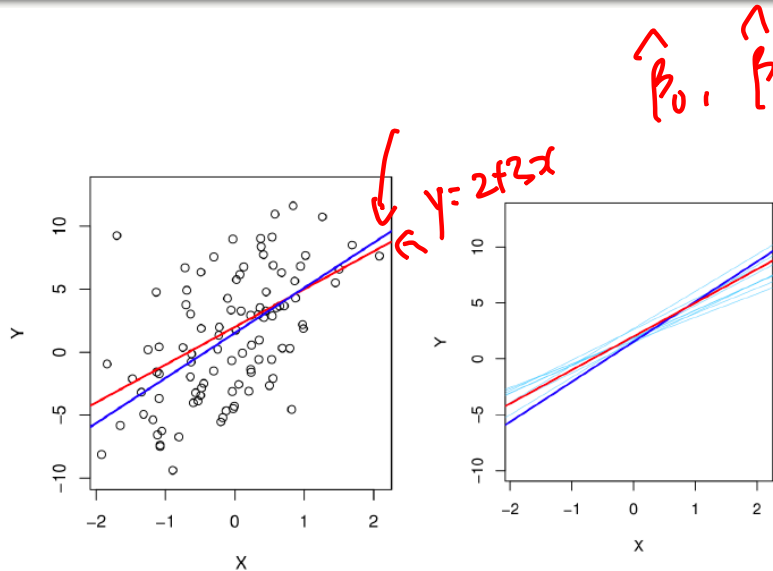
## Least squares criterion

Find  $\beta_0$  and  $\beta_1$  that minimize the RSS.

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

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

# Ordinary Least Squares



- 100 data points drawn from  $Y = 2 + 3X + \varepsilon$
- $\varepsilon$  drawn from normal distribution with mean 0
- Red line is true relationship, blue is least squares estimate
- Repeat this 10 times and plot all the found lines (in variations of blue)
- The resulting models are slightly different but are all around the red true relationship

# Variance of OLS estimates

- Variance of linear regression estimates:

$$SE(\hat{\beta}_0) = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

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

where  $\sigma^2 = \text{Var}(\varepsilon)$

- Residual standard error is an estimate of  $\sigma$

$$RSE = \sqrt{RSS/(n-2)}$$

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

↓

$$\varepsilon = Y - \beta_0 - \beta_1 X$$

$$\begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} - \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} \begin{pmatrix} 1 \\ x_1 \end{pmatrix}$$

$$\text{Var}(\varepsilon) = E((\varepsilon - E\varepsilon)^2)$$

$$\approx \frac{1}{n} \left( \left( \varepsilon_i - \left( \frac{\sum \varepsilon_i}{n} \right) \right)^2 \right)$$

$$= 0$$

$$\hat{\sigma}^2 = \frac{RSS}{n-2} = \frac{1}{n} \sum (\varepsilon_i^2) = \frac{1}{n} \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

That is, there is approximately a 95% chance that the interval

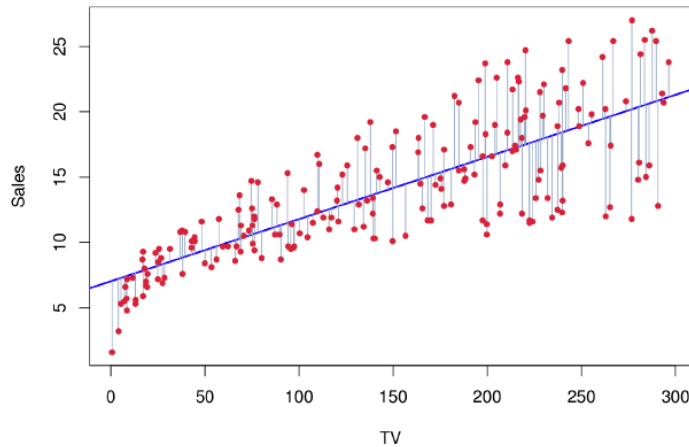
$$\approx \left[ \hat{\beta}_1 - \textcircled{2} \cdot \text{SE}(\hat{\beta}_1), \hat{\beta}_1 + 2 \cdot \text{SE}(\hat{\beta}_1) \right] \quad \text{not ideal.}$$

will contain the true value of  $\beta_1$  (under a scenario where we got repeated samples like the present sample)

For the advertising data, the 95% confidence interval for  $\beta_1$  is  $[0.042, 0.053]$



# Confidence Interval: Ad data



For the advertising data set, the 95% CIs are:

- $\beta_1 :: [0.042, 0.053]$ 
  - ▶ First line through  $[0, 8]$  and  $[300, 20.6]$
  - ▶ Second line through  $[0, 7]$  and  $[300, 21.9]$
- $\beta_0 :: [6.130, 7.935]$

# Hypothesis Testing

$SE(\hat{\beta}_1)$

- Standard errors can also be used to perform **hypothesis tests** on the coefficients. The most common hypothesis test involves testing the **null hypothesis** of

$H_0$  : There is no relationship between  $X$  and  $Y$   $\Leftarrow$   
versus the *alternative hypothesis*

$H_A$  : There is some relationship between  $X$  and  $Y$ .

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

$$H_0: \beta_1 = 0$$

$$Y = \beta_0 + \varepsilon$$

- Standard errors can also be used to perform **hypothesis tests** on the coefficients. The most common hypothesis test involves testing the **null hypothesis** of

$H_0$  : There is no relationship between  $X$  and  $Y$   
versus the *alternative hypothesis*

$H_A$  : There is some relationship between  $X$  and  $Y$ .

- Mathematically, it is

$$H_0 : \beta_1 = 0$$

versus

$$H_A : \beta_1 \neq 0,$$

since if  $\beta_1 = 0$  then the model reduces to  $Y = \beta_0 + \epsilon$ , and  $X$  is not associated with  $Y$ .

- Mathematically, it is

$$H_0 : \underline{\beta_1 = 0}$$

versus

$$H_A : \beta_1 \neq 0,$$

since if  $\beta_1 = 0$  then the model reduces to  $Y = \beta_0 + \epsilon$ , and  $X$  is not associated with  $Y$ .

- We have  $\hat{\beta}_1$  from data and want to test whether it is far from 0. But how far?

$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}$$

# Hypothesis Testing

- Mathematically, it is

$$H_0 : \beta_1 = 0$$

versus

$$H_A : \beta_1 \neq 0,$$

since if  $\beta_1 = 0$  then the model reduces to  $Y = \beta_0 + \epsilon$ , and  $X$  is not associated with  $Y$ .

- We have  $\hat{\beta}_1$  from data and want to test whether it is far from 0. But how far?

$$t = \frac{\hat{\beta}_1 - 0}{\text{SE}(\hat{\beta}_1)}$$

$$\hat{\sigma}^2 = \frac{\text{RSS}}{n-2}$$

- This will have a t-distribution with  $n-2$  degrees of freedom, assuming  $\beta_1 = 0$
- Using statistical software, it is easy to compute the probability of observing any value equal to  $|t|$  or larger. We call this probability the **p-value**.
- We will specify an alpha value (0.05) before Hypothesis testing. If p-value less than the alpha value, we will reject the null hypothesis.

# Test Statistic and p-value

Test statistic:

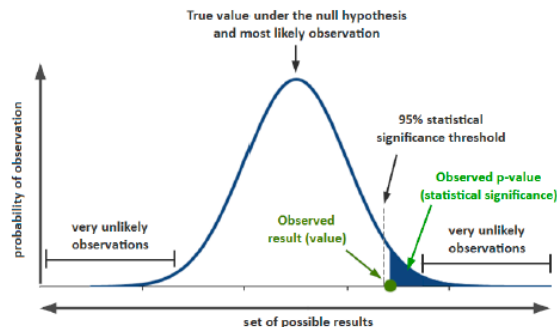
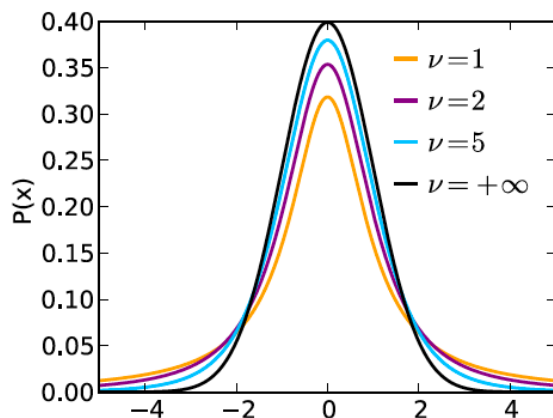
$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}$$

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

$$\varepsilon \sim N(0, \sigma^2)$$

Draw me:

t-distribution with  $n - 2$  degrees of freedom



# Results for the Advertising Data

	Coefficient	Std. Error	t-statistic	p-value
$\beta_0$ Intercept	7.0325	0.4578	15.36	< 0.0001
TV $\beta_1$	<u>0.0475</u>	<u>0.0027</u>	<u>17.67</u>	<u>&lt; 0.0001</u>

0.05

- Since  $p\text{-value} < 0.05$ , we reject the null hypothesis and conclude that TV is related to sale.

# Assessing the Accuracy of the

- Residual Standard Error

$$\hat{\sigma} = \text{RSE} = \sqrt{\frac{1}{n-2} \text{RSS}} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2},$$

- estimate of the standard deviation of  $\varepsilon$
- Avg amount that the response will deviate from the true regression line
- avg amount response will deviate from the true regression line



# Assessing the Accuracy of the

- Residual Standard Error

$$\text{RSE} = \sqrt{\frac{1}{n-2} \text{RSS}} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2},$$

- R-square**: fraction of variance explained by the linear model

$$R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}$$

*- Variance can be explained by your model*

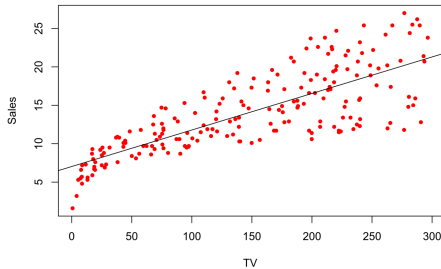
where  $\text{TSS} = \sum_{i=1}^n (y_i - \bar{y})^2$  is the *total sum of squares*.

where  $\text{TSS} = \sum_{i=1}^n (y_i - \bar{y})^2$  is the total sum of squares

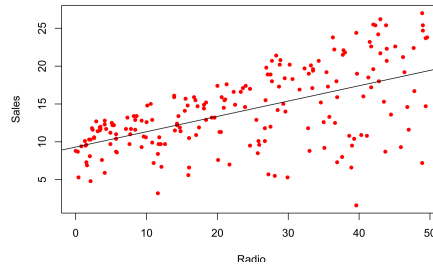
*↑ nothing about X*

$$\begin{aligned} \text{RSS} &= \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \\ &= \sum \varepsilon_i^2 \end{aligned}$$

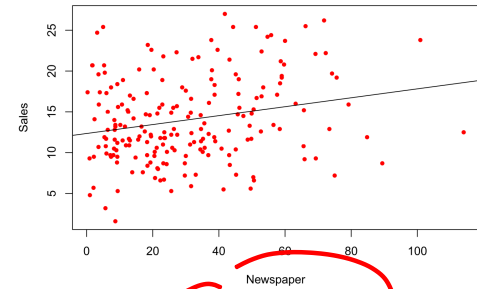
# Advertisement Data



$$R^2 = 0.61$$



$$R^2 = 0.33$$



$$R^2 = 0.05$$



# Correlation Coefficient

- Measures dependence between two random variables  $X$  and  $Y$

$$r = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$$

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

- Correlation coefficient  $r$  is between  $[-1, 1]$ 
  - 0: Variables are not linearly related
  - 1: Variables are perfectly related (same)
  - 1: Variables are negatively related (different)

$$R^2 = r^2$$