

- ]

1

## Module 3: Linear Regression

Lecture 4  
Jan 19th, 2023

$$\text{grade } f(x) = E(Y | X=x)$$



**SPARTANS WILL.**

# Recap

MICHIGAN STATE UNIVERSITY

$$\hat{f}_1(x) = 0.8 + 0.6x$$

$$\hat{f}_2(x) = 0.78 + 0.5x$$

train

$f(x)$

random

$k=hn$

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

Estimating  $f$ : non-parametric vs parametric

Assessing Model Accuracy: Training MSE and Test MSE

$$Var(X) = E((X - EX)^2)$$

train

$\hat{f}$

Test

test MSE

Bias-Variance Trade-off

↑

$$X \sim N(0, 1)$$

$$EX = 0$$

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

$$E(|\hat{f}(x) - f(x)|) = \text{Bias}$$



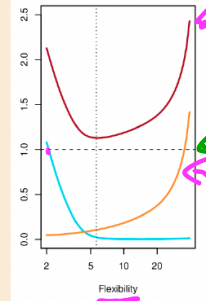
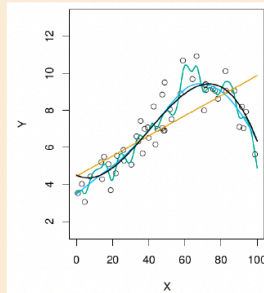
$$Var(\hat{f}(x)) = E((\hat{f}(x) - E(\hat{f}(x)))^2)$$

Estimating  $f$ : non-parametric vs parametric

Assessing Model Accuracy: Training MSE and Test MSE

★

Bias



Label the line corresponding to each of the following:

- MSE
- Bias
- Variance of  $\hat{f}(x_0)$
- Variance of  $\varepsilon$

$$Y_0 = f(x_0) + \varepsilon$$

$$E(y_0 - \hat{f}(x_0))^2 = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\varepsilon)$$

Irreducible error



SPARTANS WILL.

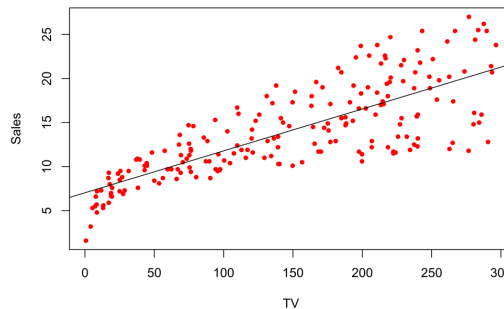
# Simple Linear Regression

- We have only **one** feature

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

Where  $\beta_0$  and  $\beta_1$  are two unknown constants also known as coefficient or parameters.

- Example:



$$\text{Sales} \approx \beta_0 + \beta_1 \times \text{TV}$$

- Least squares coefficient estimates for linear regression
- Residual sum of squares (RSS)
- Confidence interval, hypothesis test, and p-value for coefficient estimates
- Residual standard error (RSE)
- R squared

# Simple Linear Regression

- We have only one feature

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

Where  $\beta_0$  and  $\beta_1$  are two unknown constants also known as coefficient or parameters.

- Given the training data, we can estimate  $\hat{\beta}_0$  and  $\hat{\beta}_1$  for the model coefficients and predict future Y using

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x,$$

Where  $\hat{y}$  indicates a prediction of Y on the basis of  $X = x$ . The hat symbol denotes an estimated value.

# Example

		TV	Radio	Newspaper	Sales
1					
2	1	230.1	37.8	69.2	22.1
3	2	44.5	69.3	45.1	10.4
4	3	17.2	45.9	69.3	9.3
5	4	151.5	41.3	58.5	18.5
6	5	180.8	10.8	58.4	12.9
7	6	8.7	48.9	75	7.2
8	7	57.5	32.8	23.5	11.8
9	8	120.2	19.6	11.6	13.2
10	9	8.6	2.7	1	4.8
11	10	199.8	2.6	21.2	10.6
12	11	66.1	5.8	24.2	8.6

$$e_i = (y_i - \hat{y}_i)$$

$$\text{sales} \approx \beta_0 + \beta_1 \text{TV}$$

- $\beta_0$  intercept;  $\beta_1$  slope
- Coefficients or parameters :  $\{\beta_0, \beta_1\}$
- Once we have good guesses for  $\hat{\beta}_i$ , model is

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$



# How to estimate?

- Let  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$  be the prediction of Y based on the  $i$ th value of X.

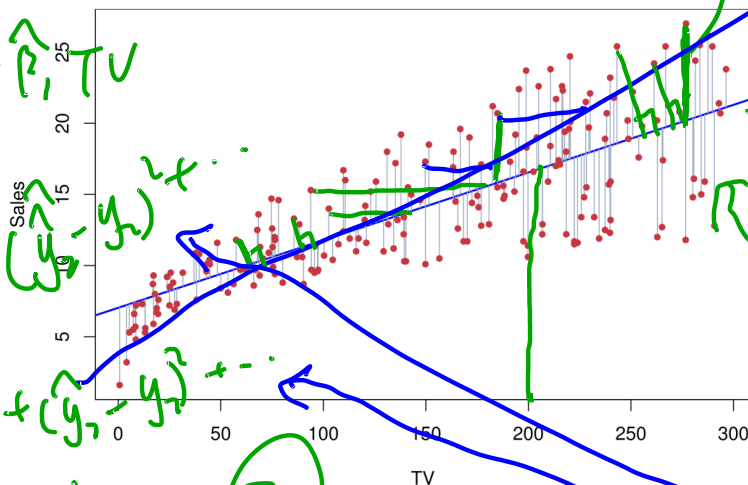
$$TV = \alpha_0 + \alpha_1 \text{Sale}$$

Y. Errors:  $e_i = y_i - \hat{y}_i$

$$\text{Sale} = \hat{\gamma}_0 + \hat{\gamma}_1 TV \leftarrow \min \text{ distance}$$

$$\alpha_0 \neq \beta_0$$

$$\text{Sale} = \hat{\beta}_0 + \hat{\beta}_1 TV$$



$$\hat{\beta}_0 + \hat{\beta}_1 x$$

$$A = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots$$

$$B = (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \dots$$

$$A > B$$

$$Y = \gamma_0 + \gamma_1 x$$

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

# How to estimate?

- Residual Sum of Squares

$$\text{RSS} = e_1^2 + e_2^2 + e_3^2 + \cdots + e_n^2 = \sum_{i=1}^n e_i^2$$

- Equivalence

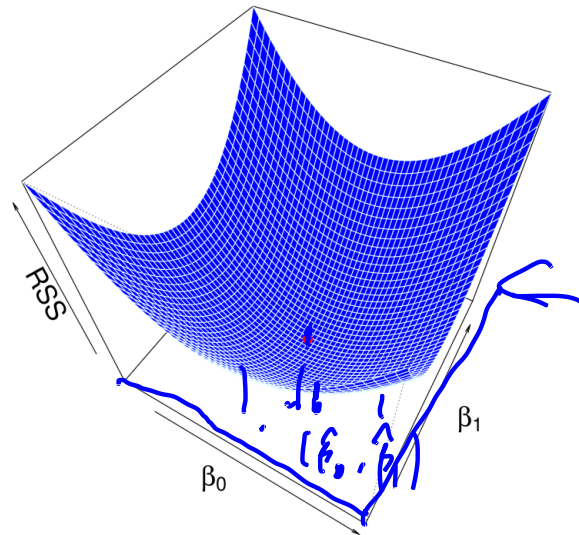
$$\min \text{RSS} = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

$(\hat{\beta}_0, \hat{\beta}_1)$

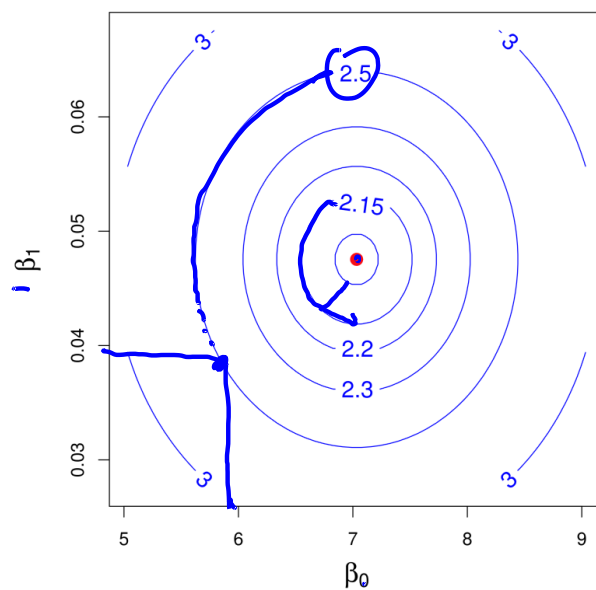
# How to estimate?

- Minimizing RSS to find the best estimates.

$$\min_{\beta_0, \beta_1} \text{RSS} = \min_{\beta_0, \beta_1} \sum_{i=1}^n e_i^2 = \min_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$



$$\min_{\beta_0, \beta_1} \text{RSS} = \min_{\beta_0, \beta_1} \sum_{i=1}^n e_i^2 = \min_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$



# Solving for minimal RSS

$$\min_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

- ▶ RSS is a **convex** function of  $\beta_0, \beta_1$
- ▶ Minimum achieved when (recall the chain rule):

$$\frac{\partial \text{RSS}}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\frac{\partial \text{RSS}}{\partial \beta_1} = -2 \sum_{i=1}^n x_i (y_i - \beta_0 - \beta_1 x_i) = 0$$

# Linear Regression Coefficients

$$\min_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Solution:

$$\begin{aligned} \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n x_i (y_i - \bar{y})}{\sum_{i=1}^n x_i (x_i - \bar{x})} \end{aligned}$$

where

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$



# Linear Regression Coefficients

MICHIGAN STATE  
UNIVERSITY

$$Y = X\beta + \tilde{\epsilon}$$

$$Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$X = \begin{pmatrix} x_{11}, \dots, x_{1p} \\ \vdots \\ x_{n1}, \dots, x_{np} \end{pmatrix}$$

$$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{pmatrix}$$

$$RSS = \underbrace{(Y - X\beta)^T}_{n \times 1} (Y - X\beta) = \underbrace{(Y^T Y)}_{\text{scalar}} - \underbrace{\beta^T X^T X \beta}_{\text{scalar}} - \underbrace{2Y^T \beta}_{\text{scalar}}$$

$$\frac{\partial RSS}{\partial \beta} = 0 + 2X^T X \beta - 2X^T Y = 0$$

$$\Rightarrow (X^T X) \beta = X^T Y \Rightarrow \hat{\beta} = \underbrace{(X^T X)^{-1}}_{p \times p} \underbrace{X^T Y}_{n \times p}$$

$$\hat{\beta} = \beta + (X^T X)^{-1} X^T \tilde{\epsilon}$$

$$\underline{E(\hat{\beta})} = \beta + \underline{E((X^T X)^{-1} X^T \tilde{\epsilon})}$$

$$= \beta + (X^T X)^{-1} X^T \underline{E(\tilde{\epsilon})}$$

$$= \beta$$

$$= (X^T X)^{-1} X^T (X\beta + \tilde{\epsilon})$$

$$= \cancel{(X^T X)^{-1} X^T X} \beta + (X^T X)^{-1} X^T \tilde{\epsilon}$$

# Why Minimizing RSS?

- It has closed form!!

✱ • Maximize likelihood when  $Y = \beta_0 + \beta_1 X + \epsilon$  , when  $\epsilon \sim \mathcal{N}(0, \sigma^2)$

- Best Linear Unbiased Estimator (BLUE): Gauss-Markov Theorem.

OLS  $\hat{\beta}$        $E(\hat{\beta}) = \beta$

↑



# Bias in Estimation

- Assume a true value  $\mu^*$
- An estimate from training data  $\hat{\mu}$
- The estimate is unbiased if  $E(\hat{\mu}) = \mu^*$

$$E \underline{X} = \mu^*$$

$$X_1, X_2, \dots, X_n$$

- Sample mean is unbiased for population mean

$$E(\hat{\mu}) = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \mu$$

- Standard variance estimate is biased

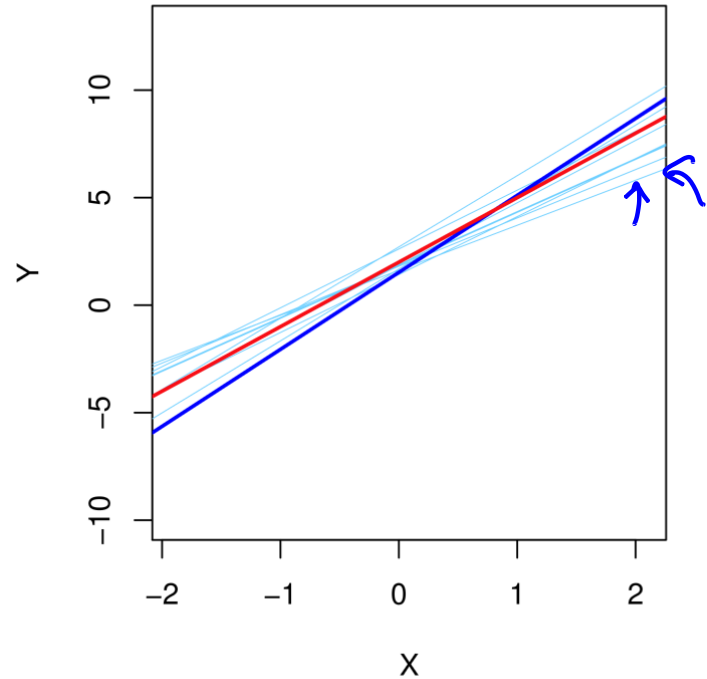
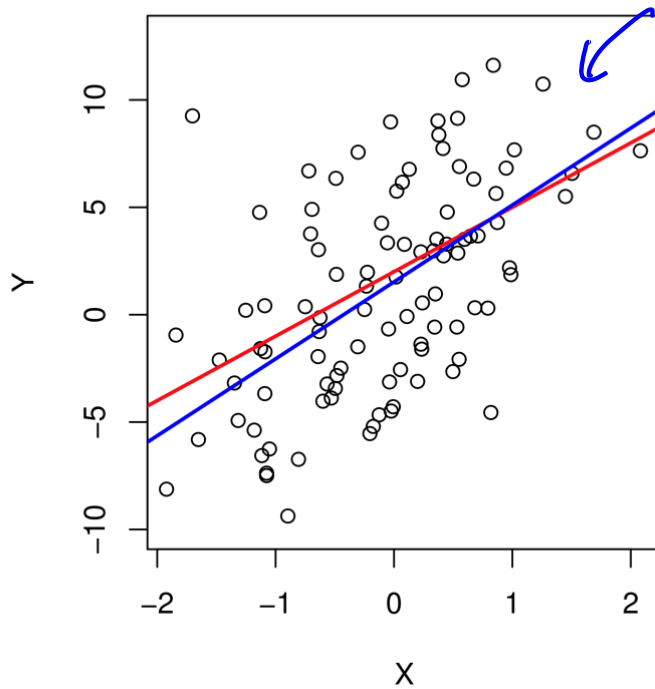
$$E(\hat{\sigma}^2) = E\left[\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right] \neq \sigma^2$$

$$V_n(X) = \sigma^2$$

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2$$

$$\begin{aligned} S &= \frac{X_1 + X_2 + \dots + X_n}{n} \\ E(S) &= E\left(\frac{X_1 + \dots + X_n}{n}\right) \\ &= \frac{EX_1 + EX_2 + \dots + EX_n}{n} \\ &= \frac{n\mu^*}{n} = \mu^* \end{aligned}$$

# Linear Regression is Unbiased



# Variance of the estimates

- Assume a true value  $\mu^*$
- An estimate from training data  $\hat{\mu}$
- The variance of sample mean is:

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

$$\text{Var}(\hat{\mu}) = \frac{\sigma^2}{n}$$

The variances of the linear regression estimates are

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

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

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

$$e_i^2 = (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

$$RSS = \sum_{i=1}^n e_i^2$$

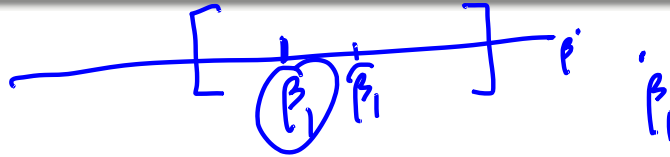
$$\text{Var}(\epsilon) = \sigma^2$$

- We then can calculate the confidence intervals. The 95% confidence interval for  $\beta_1$  approximately takes the form

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

$$\hat{\beta}_1 \pm 2 \hat{\sigma} \leftarrow 95\% \text{ CI}$$

# Confidence Interval



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

$$\left[ \hat{\beta}_1 - 2 \cdot \text{SE}(\hat{\beta}_1), \hat{\beta}_1 + 2 \cdot \text{SE}(\hat{\beta}_1) \right]$$

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  $[0.042, 0.053]$

- 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  $Y$ .

Handwritten blue ink notes showing the linear regression equation  $Y = \beta_0 + \beta_1 X$ . The  $\beta_1$  coefficient is circled, and an arrow points to it from the text  $H_A : \beta_1 \neq 0$ .

# Hypothesis Testing

- Mathematically, it is

$$H_0 : \beta_1 = 0$$

versus

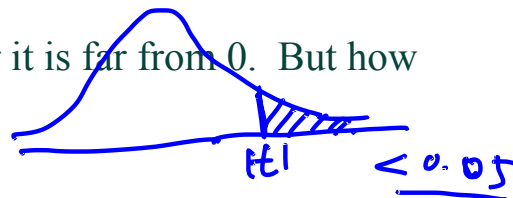
$$H_A : \beta_1 \neq 0,$$

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?

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

- 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.



# Results for the Advertising Data

	Coefficient	Std. Error	t-statistic	p-value
$\hat{\beta}_0$ Intercept	7.0325	0.4578	15.36	< 0.0001
$\hat{\beta}_1$ TV	0.0475	0.0027	17.67	< 0.0001

< 0.05

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