## Lecture 1

Introduction

STAT 8020 Statistical Methods II August 20, 2020



Who is the instructor?

Class Policies Schedule

Tell us about yourself

Simple Linear Regression

SLR Parameter

Residual Analys

Whitney Huang Clemson University



Who is the instructor'

Class Policie Schedule

Tell us about yoursel

Regression

SLR Parameter Estimation

Residual Analysis

# Who is the instructor?

#### Who am I?

- Second year Assistant Professor of Applied Statistics and Data Science
- Born in Laramie, Wyoming, grew up in Taiwan





 With a B.S. in Mechanical Engineering, switched to Statistics in graduate school

Got a Ph.D. (Statistics) in 2017 at Purdue University.





Who is the instructor?

Class Policies Schedule

Tell us about yourself

mple Linear egression

LR Paramete stimation

### How to reach me?

CLEMS N U N I V E R S I T Y

Who is the instructor?

Class Policies Schedule

Tell us about yourself

Regression

SLR Paramete Estimation

Residual Analysis

Email: wkhuang@clemson.edu

Office: O-221 Martin Hall

 Office Hours: TR 11:00am – 12:00pm and by appointment

1.4



Who is the instructor?

#### Class Policies / Schedule

Tell us about yoursel

Regression

SLR Parameter Estimation

Residual Analysis

# Class Policies / Schedule

## Logistics

- We will meet TR 12:30pm 1:45pm via Zoom
- There will be three online exams and a (comprehensive) online final. The (tentative) dates for the three exams are:
  - Exam I: Sept. 24, Thursday
  - Exam II: Oct. 20, Tuesday
  - Exam II: Nov. 12, Tuesday
  - The Final Exam will be given on Wednesday, Dec. 7, 3:00 pm -5:30 pm.
- No classes on Nov. 3 (Fall Break) & 26 (Thanksgiving)



Who is the instructor?

#### Class Policies / Schedule

Tell us about yourself

Regression

estimation

#### **Class Website**

## CANVAS and my teaching website (link:

https://whitneyhuang83.github.io/STAT8020/Fall2020/stat8020\_2020Fall.html)

- Course syllabus [Link] / Announcements
- Lecture slides/notes
- Exam schedule
- Data sets
- R code



Who is the instructor?

Class Policies . Schedule

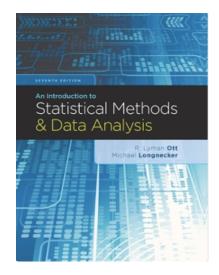
Tell us about yourself

Simple Linear Regression

SLR Paramete Estimation

#### **Recommended Textbook**

An Introduction to Statistical Methods and Data Analysis, 6<sup>th</sup> Edition. Lyman Ott and Micheal T. Longnecker, Duxbury, **2010**; ISBN-13: 978-1305269477





Who is the instructor?

Class Policies Schedule

Tell us about yourself

Simple Linear Regression

LR Parameter stimation

### **Evaluation**

• Grade Distribution:

Exam I:	25%
Exam II	25%
Exam III	25%
Final Exam	25%

<u>>= 90 00</u>

Letter Grade:

/ >= 30.00	_ ^
$88.00 \sim 89.99$	A-
$85.00 \sim 87.99$	B+
$80.00 \sim 84.99$	В
$78.00 \sim 79.99$	B-
$75.00 \sim 77.99$	C+
$70.00 \sim 74.99$	С
$68.00 \sim 69.99$	C-
<= 67.99	F

Δ



Who is the instructor?

Class Policies . Schedule

Tell us about yoursel

Simple Linea Regression

SLH Paramete Stimation

## **Tentative Topics and Dates**

## Part I: Regression Analysis (August 20 – September 24)

- Review of Simple Linear Regression
- Multiple Linear Regression: Statistical Inference; Model Selection and Diagnostics
- Regression Models with Quantitative and Qualitative Predictors
- Nonlinear and Non-parametric Regression

## Part II: Categorical Data Analysis (September 29 – October 20)

- Review of Inference for Proportions and Contingency Tables
- Relative Risk and Odds Ratio
- Logistic Regression and Poisson Regression



Who is the instructor?

Class Policies / Schedule

Tell us about yourself

nple Linear gression

Estimation

## **Tentative Topics and Dates cont'd**

#### CLEMS N U N I V E R S I T Y

Who is the instructor?

Class Policies / Schedule

Tell us about yourself

nple Linear gression

SLR Parameter Stimation

Residual Analys

## Part III: Experimental Design (October 22 – November 12)

- Introduction to Experimental Design: Principles and Techniques
- Completely randomized Designs, Block Designs, Latin Square Designs, Nested and Split-Plot Designs
- Computer experiments

# Part IV: Multivariate, Spatial and Time Series Analysis (November 17 – December 3)

- Discriminate Analysis, Principle Components Analysis, and Cluster Analysis
- Basic of time series and spatial data analysis

## Computing

We will use software to perform statistical analyses. The recommended software for this course are  ${\tt JASP}$  and

- R/Rstudio
   JASP
  - a free/open-source graphical program for statistical analysis
  - available at https://jasp-stats.org/
  - R Studio
    - a free/open-source programming language for statistical analysis
    - available at https://www.r-project.org/(R); https://rstudio.com/(Rstudio)

You are welcome to use a different package (e.g. SAS, JMP, SPSS, Minitab) if you prefer



Who is the instructor?

Class Policies / Schedule

Tell us about yourself

gression

R Paramete



Who is the instructor

Class Policies / Schedule

#### Tell us about yourself

Simple Linear Regression

SLR Parameter Estimation

Residual Analysis

# Tell us about yourself

## Tell us about yourself

- CLEMS N
  - Who is the instructor?
    - Class Policies Schedule

#### Tell us about yourself

- Simple Linear Regression
- SLR Paramete Estimation
- Residual Analysis

- Your name
- Degree program
- Your background in Statistics/Computing



Who is the instructor

Class Policie Schedule

Tell us about yoursel

#### Simple Linear Regression

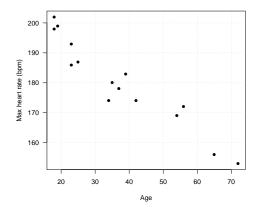
LR Parameter stimation

Residual Analysis

# Review of Simple Linear Regression

## What is Regression Analysis?

**Regression analysis**: A set of statistical procedures for estimating the relationship between response variable and predictor variable(s)



We will focus on simple linear regression in the next few lectures



Who is the instructor?

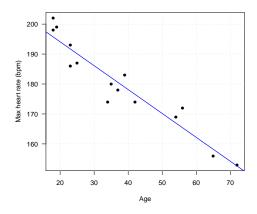
lass Policies . chedule

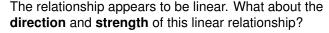
Tell us about yourself

Simple Linear Regression

LH Parameter stimation

## Scatterplot: Is Linear Trend Reasonable?





```
> cov(age, maxHeartRate)
[1] -243.9524
```



Who is the instructor?

Class Policies Schedule

Tell us about yourself

Simple Linear Regression

SLR Parameter Estimation

## **Scatterplot: Is Linear Trend Reasonable?**





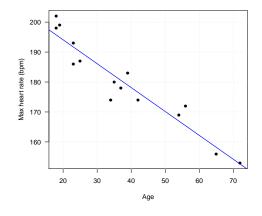
Class Policies Schedule

Tell us about yourself

#### Simple Linear Regression

SLR Parameter Estimation

Residual Analy



The relationship appears to be linear. What about the **direction** and **strength** of this linear relationship?

## Simple Linear Regression (SLR)

*Y*: dependent (response) variable; *X*: independent (predictor) variable

 In SLR we assume there is a linear relationship between X and Y:

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

- We need to estimate  $\beta_0$  (intercept) and  $\beta_1$  (slope)
- We can use the estimated regression equation to
  - make predictions
  - study the relationship between response and predictor
  - control the response
- Yet we need to quantify our estimation uncertainty regarding the linear relationship (will talk about this next time)



Who is the instructor?

Class Policies Schedule

Tell us about yourself

Simple Linear Regression

> LR Paramete stimation

## **Regression equation:** $Y = \beta_0 + \beta_1 X$



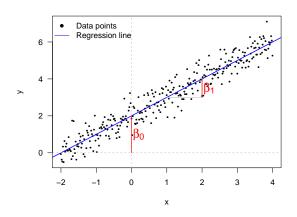


Class Policies Schedule

Tell us about voursel

#### Simple Linear Regression

SLR Paramete



- $\beta_0$ : E[Y] when X = 0
- $\beta_1$ : E[ $\Delta Y$ ] when X increases by 1

## Assumptions about the Random Error $\varepsilon$

CLEMS N UNIVERSITY

In order to estimate  $\beta_0$  and  $\beta_1$ , we make the following assumptions about  $\varepsilon$ 

- $\bullet \ \mathrm{E}[\varepsilon_i] = 0$
- $Var[\varepsilon_i] = \sigma^2$
- $Cov[\varepsilon_i, \varepsilon_j] = 0, \quad i \neq j$

Therefore, we have

The regression line  $\beta_0 + \beta_1 X$  represents the **conditional mean curve** whereas  $\sigma^2$  measures the magnitude of the **variation** around the regression curve

Who is the instructor?

lass Policies / chedule

Tell us about yourself

Simple Linear Regression

R Parameter stimation

For the given observations  $(x_i, y_i)_{i=1}^n$ , choose  $\beta_0$  and  $\beta_1$  to minimize the *sum of squared errors*:

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

Solving the above minimization problem requires some knowledge from Calculus....



who is the instructor?

Class Policies . Schedule

Tell us about yourself

imple Linear egression

SLR Parameter Estimation

For the given observations  $(x_i, y_i)_{i=1}^n$ , choose  $\beta_0$  and  $\beta_1$  to minimize the *sum of squared errors*:

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

Solving the above minimization problem requires some knowledge from Calculus....

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{(X_i - \bar{X})^2}$$



who is the instructor?

Class Policies Schedule

Tell us about yourself

legression

Estimation

For the given observations  $(x_i, y_i)_{i=1}^n$ , choose  $\beta_0$  and  $\beta_1$  to minimize the *sum of squared errors*:

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

Solving the above minimization problem requires some knowledge from Calculus....

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{(X_i - \bar{X})^2}$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$



Who is the instructor?

Class Policies Schedule

Tell us about yourself

mple Linear egression

SLR Parameter Estimation

For the given observations  $(x_i, y_i)_{i=1}^n$ , choose  $\beta_0$  and  $\beta_1$  to minimize the *sum of squared errors*:

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

Solving the above minimization problem requires some knowledge from Calculus....

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{(X_i - \bar{X})^2}$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$



Who is the instructor?

Class Policies Schedule

Tell us about yourself

mple Linear egression

SLR Parameter Estimation

For the given observations  $(x_i, y_i)_{i=1}^n$ , choose  $\beta_0$  and  $\beta_1$  to minimize the *sum of squared errors*:

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

Solving the above minimization problem requires some knowledge from Calculus....

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{(X_i - \bar{X})^2}$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta_1} \bar{X}$$

We also need to **estimate**  $\sigma^2$ 



who is the instructor?

Class Policies Schedule

Tell us about yourself

egression

SLR Parameter Estimation

For the given observations  $(x_i, y_i)_{i=1}^n$ , choose  $\beta_0$  and  $\beta_1$  to minimize the *sum of squared errors*:

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

Solving the above minimization problem requires some knowledge from Calculus....

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{(X_i - \bar{X})^2}$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta_1} \bar{X}$$

We also need to **estimate**  $\sigma^2$ 

$$\hat{\sigma}^2=rac{\sum_{i=1}^n(Y_i-\hat{Y}_i)^2}{n-2},$$
 where  $\hat{Y}_i=\hat{eta}_0+\hat{eta}_1X_i$ 



who is the instructor?

Class Policies Schedule

Tell us about yourself

imple Linear egression

LR Parameter stimation

## **Properties of Least Squares Estimates**



- Gauss-Markov theorem states that in a linear regression these least squares estimators
  - Are unbiased, i.e.,
    - $E[\hat{\beta}_1] = \beta_1; E[\hat{\beta}_0] = \beta_0$
    - $\bullet \ \mathrm{E}[\hat{\sigma}^2] = \sigma^2$
  - Have minimum variance among all unbiased linear estimators

Note that we do not make any distributional assumption on  $\varepsilon_i$ 

Who is the instructor?

lass Policies / chedule

Tell us about yourself

mple Linear gression

stimation

## **Example: Maximum Heart Rate vs. Age**

CLEMS N

The maximum heart rate MaxHeartRate of a person is often said to be related to age Age by the equation:

MaxHeartRate = 220 - Age.

Suppose we have 15 people of varying ages are tested for their maximum heart rate (bpm) (link to the "dataset":

whitneyhuang83.github.io/STAT8010/Data/
maxHeartRate.csv)

- Compute the estimates for the regression coefficients
- Compute the fitted values
- **1** Compute the estimate for  $\sigma$

Who is the instructor?

lass Policies / chedule

Tell us about yourself

mple Linear egression

Estimation

## Estimate the Parameters $\beta_1$ , $\beta_0$ , and $\sigma^2$

 $Y_i$  and  $X_i$  are the Maximum Heart Rate and Age of the i<sup>th</sup> individual

- To obtain  $\hat{\beta}_1$ 
  - Ompute  $\bar{Y} = \frac{\sum_{i=1}^{n} Y_i}{n}$ ,  $\bar{X} = \frac{\sum_{i=1}^{n} X_i}{n}$
  - ② Compute  $Y_i \bar{Y}$ ,  $X_i \bar{X}$ , and  $(X_i \bar{X})^2$  for each observation
  - Ompute  $\sum_{i=1}^{n} (X_i \bar{X})(Y_i \bar{Y})$  divived by  $\sum_{i=1}^{n} (X_i \bar{X})^2$
- $\hat{\beta}_0$ : Compute  $\bar{Y} \hat{\beta}_1 \bar{X}$
- $\hat{\sigma}^2$ 
  - Ompute the fitted values:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$ ,  $i = 1, \dots, n$
  - Ompute the **residuals**  $e_i = Y_i \hat{Y}_i, i = 1, \dots, n$
  - Occupance of Squares (RSS)  $= \sum_{i=1}^{n} (Y_i \hat{Y}_i)^2 \text{ and divided by } n 2 \text{ (why?)}$



Who is the instructor?

lass Policies / chedule

Tell us about yourself

imple Linear egression

stimation

### Let's Do the Calculations



Who is the instructor?

lass Policies chedule

Tell us about yourself

nple Linear gression

SLR Parameter Estimation

$\bar{X} = \sum_{i=1}^{15}$	$\frac{18 + 23 + \dots + 39 + 37}{15} = 37.33$
$\bar{Y} = \sum_{i=1}^{15}$	$\frac{202 + 186 + \dots + 183 + 178}{15} = 180.27$

X $Y$	18 202	23 186	25 187	35 180	65 156	54 169	34 174	56 172	72 153	19 199	23 193	42 174	18 198	39 183	37 178
	-19.33	-14.33	-12.33	-2.33	27.67	16.67	-3.33	18.67	34.67	-18.33	-14.33	4.67	-19.33	1.67	-0.33
	21.73	5.73	6.73	-0.27	-24.27	-11.27	-6.27	-8.27	-27.27	18.73	12.73	-6.27	17.73	2.73	-2.27
	-420.18	-82.18	-83.04	0.62	-671.38	-187.78	20.89	-154.31	-945.24	-343.44	-182.51	-29.24	-342.84	4.56	0.76
	373.78	205.44	152.11	5.44	765.44	277.78	11.11	348.44	1201.78	336.11	205.44	21.78	373.78	2.78	0.11
_	195.69	191.70	190.11	182.13	158.20	166.97	182.93	165.38	152.61	194.89	191.70	176.54	195.69	178.94	180.53

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = -0.7977$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} = 210.0485$$

• 
$$\hat{\sigma}^2 = \frac{\sum_{i=1}^{15} (Y_i - \hat{Y}_i)^2}{13} = 20.9563 \Rightarrow \hat{\sigma} = 4.5778$$

#### Let's Double Check

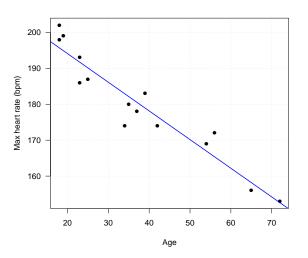
## 



```
> fit <- lm(MaxHeartRate ~ Age)</pre>
> summary(fit)
Call:
lm(formula = MaxHeartRate \sim Age)
Residuals:
    Min
            10 Median ____
                            30
                                   Max
<u>-8.9258 -2.5383</u> 0.3879 3.1867 6.6242
Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept) 210.04846 2.86694 73.27 < 2e-16 ***
             -0.79773 0.06996 -11.40 3.85e-08 ***
Age
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Residual standard error: 4.578 on 13 degrees of freedom
Multiple R-squared: 0.9091, Adjusted R-squared: 0.9021
F-statistic: 130 on 1 and 13 DF, p-value: 3.848e-08
```



## **Linear Regression Fit**



**Question:** Is linear relationship between max heart rate and age reasonable?  $\Rightarrow$  Residual Analysis



Who is the instructor?

Class Policies Schedule

Tell us about yourself

Simple Linear

SLR Parameter Estimation

### Residuals



 The residuals are the differences between the observed and fitted values:

$$e_i = Y_i - \hat{Y}_i,$$

where 
$$\hat{Y}_i = \hat{eta}_0 + \hat{eta}_1 X_i$$

- $e_i$  is NOT the error term  $\varepsilon_i = Y_i \mathrm{E}[Y_i]$
- Residuals are very useful in assessing the appropriateness of the assumptions on  $\varepsilon_i$ . Recall
  - $E[\varepsilon_i] = 0$
  - $\operatorname{Var}[\varepsilon_i] = \sigma^2$
  - $Cov[\varepsilon_i, \varepsilon_j] = 0, \quad i \neq j$

who is the instructor?

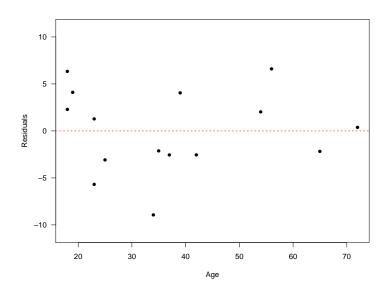
Class Policies / Schedule

Tell us about yourself

egression

Estimation

## Maximum Heart Rate vs. Age Residual Plot: $\varepsilon$ vs. X





Who is the instructor?

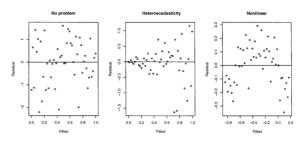
Class Policies Schedule

Tell us about yoursel

Regression

SLR Parameter Estimation

## **Interpreting Residual Plots**





Who is the instructor?

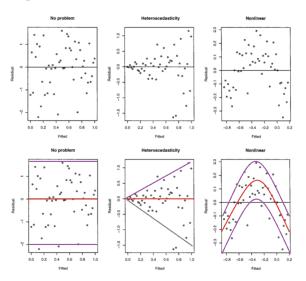
Class Policies Schedule

Tell us about yourself

Regression

Estimation Estimater

## **Interpreting Residual Plots**



**Figure:** Figure courtesy of Faraway's Linear Models with R (2005, p. 59).



Who is the instructor?

Class Policies Schedule

Tell us about yourself

egression

stimation

## **Summary**

## CLEMS N

In this lecture, we reviewed

- Simple Linear Regression:  $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$
- Method of Least Square for parameter estimation
- Residual analysis to check model assumptions
   Next time we will talk about
  - More on residual analysis
  - Normal Error Regression Model and statistical inference for  $\beta_0$ ,  $\beta_1$ , and  $\sigma^2$
  - Prediction

Who is the instructor?

Slass Policies Schedule

Tell us about yourself

mple Linear egression

Estimation