

# Lecture 2

## Simple Linear Regression I

Reading: Chapter 11

*STAT 8020 Statistical Methods II*  
August 23, 2019

Announcements

What is regression  
analysis

Simple Linear  
Regression

Whitney Huang  
Clemson University

Announcements

What is regression  
analysis

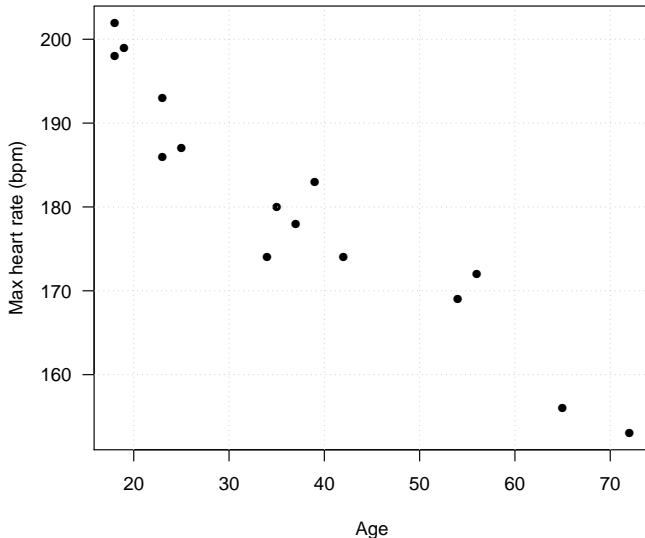
Simple Linear  
Regression

- 1 **Announcements**
- 2 **What is regression analysis**
- 3 **Simple Linear Regression**

- Syllabus and lecture notes are in CANVAS and my personal website (link: [https://whitneyhuang83.github.io/stat8020\\_2019Fall.html](https://whitneyhuang83.github.io/stat8020_2019Fall.html))
- **Academic Continuity Statement** is added in the updated syllabus (link: [https://whitneyhuang83.github.io/STAT8010\\_Syllabus\\_2019\\_Fall.pdf](https://whitneyhuang83.github.io/STAT8010_Syllabus_2019_Fall.pdf))
- Please talk to me if you would like to share your data set to be used for this class

## What is Regression Analysis?

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



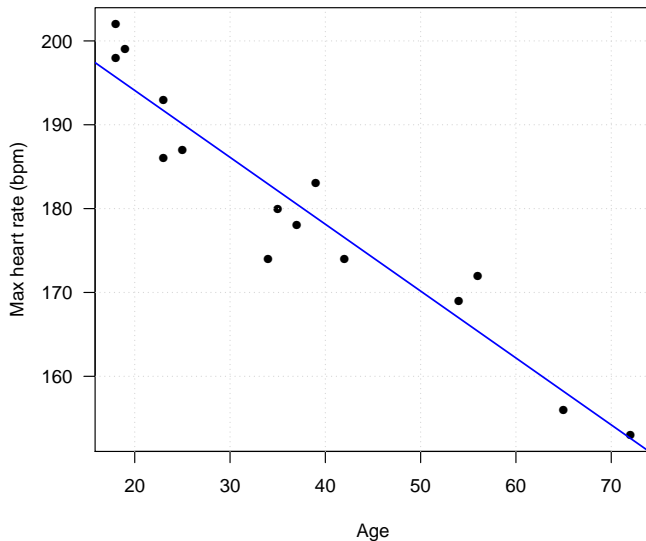
# Simple linear regression

# Scatterplot: Is Linear Trend Reasonable?

Announcements

What is regression  
analysis

Simple Linear  
Regression



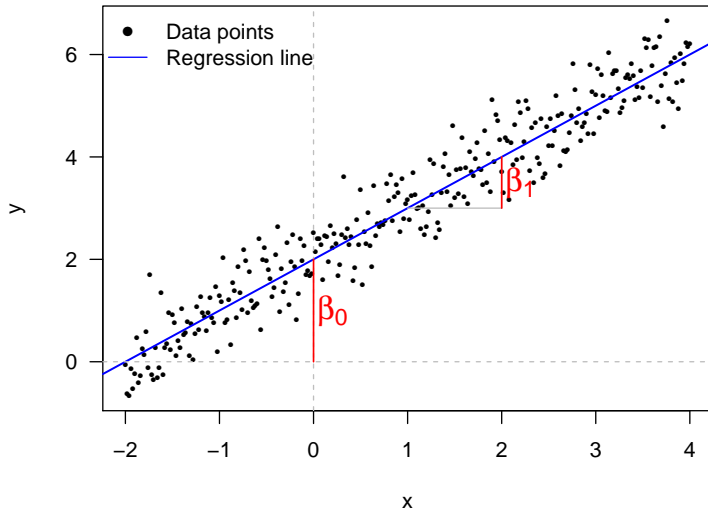
$Y$ : dependent (response) variable;  $X$ : independent (predictor) variable

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

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

- We will need to estimate  $\beta_0$  (intercept) and  $\beta_1$  (slope)
- Then 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 uncertainty regarding the linear relationship

## Regression equation: $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$





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

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

Therefore, we have

$$E[Y_i] = \beta_0 + \beta_1 X_i, \text{ and}$$

$$\text{Var}[Y_i] = \sigma^2$$

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

## Estimation: Method of Least Square

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....

- $\hat{\beta}_{1,LS} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$
- $\hat{\beta}_{0,LS} = \bar{Y} - \hat{\beta}_{1,LS} \bar{X}$

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

- $\hat{\sigma}_{LS}^2 = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n-2}$ , where  $\hat{Y}_i = \hat{\beta}_{0,LS} + \hat{\beta}_{1,LS} X_i$

- **Gauss-Markov** theorem states that in a linear regression these least squares estimators

1 Are unbiased, i.e.,

- $E[\hat{\beta}_{1,LS}] = \beta_1; E[\hat{\beta}_{0,LS}] = \beta_0$
- $E[\hat{\sigma}_{LS}^2] = \sigma^2$

2 Have **minimum variance** among all unbiased linear estimators

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

## Example: Maximum Heart Rate vs. Age

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

$$\text{MaxHeartRate} = 220 - \text{Age}.$$

Suppose we have 15 people of varying ages are tested for their maximum heart rate (bpm) (link to the “dataset”: <http://whitneyhuang83.github.io/maxHeartRate.csv>)

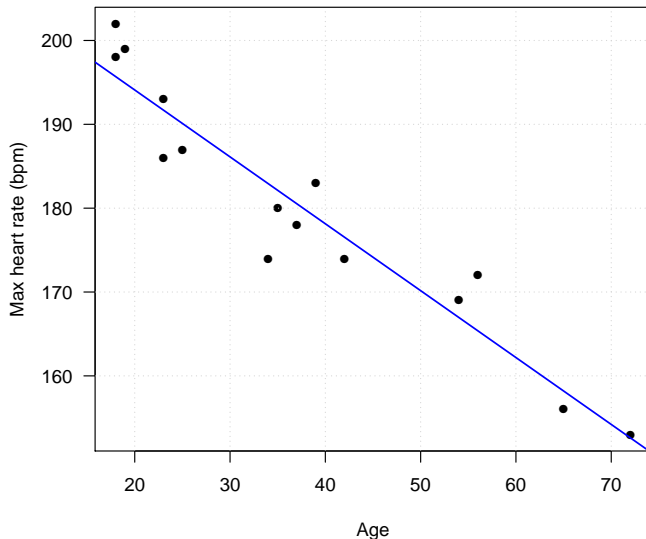
- 1 Compute the estimates for the regression coefficients
- 2 Compute the fitted values
- 3 Compute the estimate for  $\sigma$

Announcements

What is regression  
analysis

Simple Linear  
Regression

# Linear Regression Fit



**Question:** Is linear relationship between max heart rate and age reasonable?  $\Rightarrow$  [Residual Analysis](#)

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

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

$$\text{where } \hat{Y}_i = \hat{\beta}_{0,LS} + \hat{\beta}_{1,LS}X_i$$

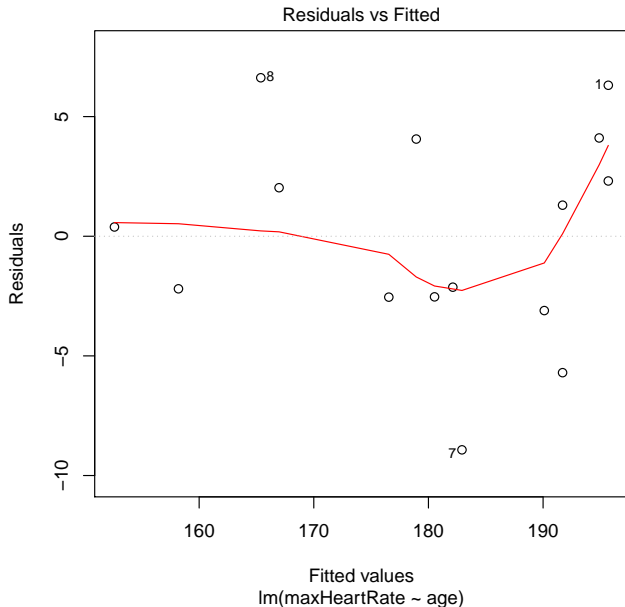
- $e_i$  is NOT the error term  $\varepsilon_i = Y_i - E[Y_i]$
- Residuals are very useful in assessing the appropriateness of the assumptions on  $\varepsilon_i$ . Recall
  - $E[\varepsilon_i] = 0$
  - $\text{Var}[\varepsilon_i] = \sigma^2$
  - $\text{Cov}[\varepsilon_i, \varepsilon_j] = 0, \quad i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

# Residual Analysis

[Announcements](#)

[What is regression  
analysis](#)

[Simple Linear  
Regression](#)

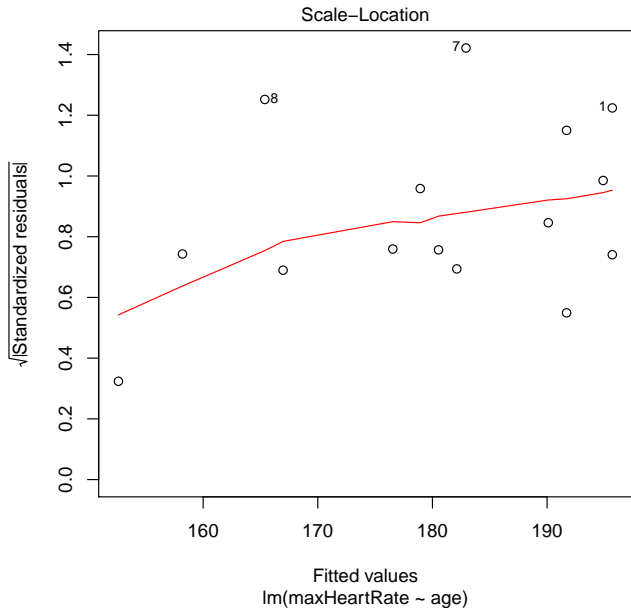


# Residual Analysis

[Announcements](#)

[What is regression  
analysis](#)

[Simple Linear  
Regression](#)





In this lecture, we learned

- **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

- 1 More on residual analysis
- 2 Normal Error Regression Model and statistical inference for  $\beta_0$ ,  $\beta_1$ , and  $\sigma^2$
- 3 Prediction