# Lecture 1
## Review of Simple Linear Regression

*DSA 8020 Statistical Methods II*
January 9-13, 2023

Whitney Huang
Clemson University

1.1

---

## Agenda

1 **Simple Linear Regression**

2 **Parameter Estimation**

3 **Residual Analysis**

4 **Confidence/Prediction Intervals**
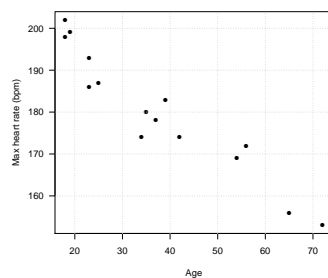
5 **Hypothesis Testing**

1.2

---

## 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: The relationship between the response variable and the predictor variable is approximately linear

1.3

## Simple Linear Regression (SLR)
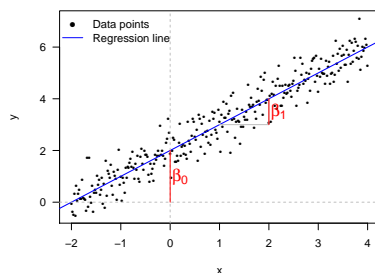
$Y$: response variable; $x$: predictor variable

- 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) based on observed data $\{x_i, y_i\}_{i=1}^n$

- 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

Simple Linear Regression
Parameter Estimation
Residual Analysis
Confidence/Prediction Intervals
Hypothesis Testing

1.4

Notes

---

## Regression equation: $Y = \beta_0 + \beta_1 x$



- $\beta_0$: $\mathrm{E}[Y]$ when $x = 0$

- $\beta_1$: $\mathrm{E}[\Delta Y]$ when $x$ increases by 1

Simple Linear Regression
Parameter Estimation
Residual Analysis
Confidence/Prediction Intervals
Hypothesis Testing

1.5

Notes

---

## Assumptions about the Random Error $\varepsilon$

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

- $\mathrm{E}[\varepsilon_i] = 0$

- $\mathrm{Var}[\varepsilon_i] = \sigma^2$

- $\mathrm{Cov}[\varepsilon_i, \varepsilon_j] = 0, \quad i \neq j$

Therefore, we have

$$\mathrm{E}[Y_i] = \beta_0 + \beta_1 x_i, \text{ and}$$
$$\mathrm{Var}[Y_i] = \sigma^2$$

> 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

Simple Linear Regression
Parameter Estimation
Residual Analysis
Confidence/Prediction Intervals
Hypothesis Testing

1.6

Notes

## Estimation: Method of Least Squares

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

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

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

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

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

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

$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n-2}$, where $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$

CLEMSON
UNIVERSITY

Simple Linear
Regression

Parameter
Estimation

Residual Analysis

Confidence/Prediction
Intervals

Hypothesis Testing

1.7

Notes

---

## Example: Maximum Heart Rate vs. Age

The maximum heart rate `MaxHeartRate` of a person is often said to be related to age `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)

1. Compute the estimates for the regression coefficients

2. Compute the fitted values

3. Compute the estimate for $\sigma$

CLEMSON
UNIVERSITY

Simple Linear
Regression

Parameter
Estimation

Residual Analysis

Confidence/Prediction
Intervals

Hypothesis Testing

1.8

Notes

---

## Maximum Heart Rate vs. Age

Output from R ( RStudio)

```
> fit <- lm(MaxHeartRate ~ Age)
> summary(fit)

Call:
lm(formula = MaxHeartRate ~ Age)

Residuals:
    Min      1Q  Median      3Q     Max
-8.9258 -2.5383  0.3879  3.1867  6.6242

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 210.04846    2.86694   73.27  < 2e-16 ***
Age          -0.79773    0.06996  -11.40 3.85e-08 ***
---
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
```

CLEMSON
UNIVERSITY

Simple Linear
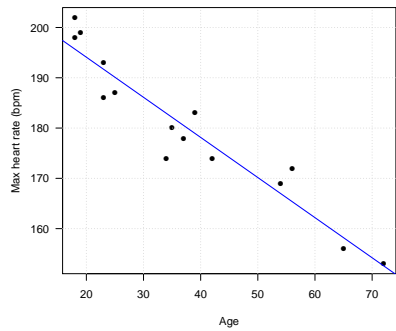Regression

Parameter
Estimation

Residual Analysis

Confidence/Prediction
Intervals

Hypothesis Testing

1.9

Notes

## Assessing Linear Regression Fit

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

---

## Residuals

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

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

where $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$

- Residuals are very useful in assessing the appropriateness of the assumptions on $\varepsilon_i$. Recall

  - $\mathrm{E}[\varepsilon_i] = 0$

  - $\mathrm{Var}[\varepsilon_i] = \sigma^2$

  - $\mathrm{Cov}[\varepsilon_i, \varepsilon_j] = 0, \quad i \neq j$

---

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

## Interpreting Residual Plots



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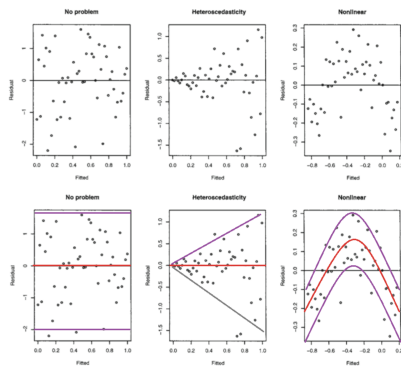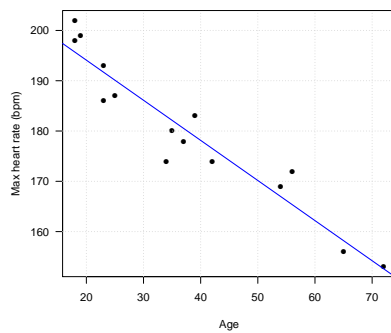
---

## How (Un)certain We Are?



**Can we formally quantify our estimation uncertainty?**
$\Rightarrow$ We need additional (distributional) assumption on $\varepsilon$

---

## Normal Error Regression Model

Recall
$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

- Further assume
  $\varepsilon_i \sim N(0, \sigma^2) \Rightarrow Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$

- With normality assumption, we can derive the
  **sampling distribution** of $\hat{\beta}_1$ and $\hat{\beta}_0 \Rightarrow$

  $\frac{\hat{\beta}_1 - \beta_1}{\hat{SE}(\hat{\beta}_1)} \sim t_{n-2}, \quad \hat{SE}(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}}$

  $\frac{\hat{\beta}_0 - \beta_0}{\hat{SE}(\hat{\beta}_0)} \sim t_{n-2}, \quad \hat{SE}(\hat{\beta}_0) = \hat{\sigma}\sqrt{\left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2}\right)}$

  where $t_{n-2}$ denotes the Student's t distribution with
  $n-2$ degrees of freedom

## Assessing Normality Assumption on $\varepsilon$

**Histogram of fit$residuals**

**Normal Q–Q Plot**

CLEMS⬤N
U N I V E R S I T Y

Simple Linear Regression

Parameter Estimation

Residual Analysis

Confidence/Prediction Intervals

Hypothesis Testing

1.16

**Notes**

---

## Confidence Intervals for $\beta_0$ and $\beta_1$

- Recall $\frac{\hat{\beta}_1 - \beta_1}{\hat{SE}(\hat{\beta}_1)} \sim t_{n-2}$, we use this fact to construct a **confidence interval (CI)** for $\beta_1$:

$$\left[\hat{\beta}_1 - t_{\alpha/2, n-2}\hat{SE}(\hat{\beta}_1), \hat{\beta}_1 + t_{\alpha/2, n-2}\hat{SE}(\hat{\beta}_1)\right],$$

  where $\alpha$ is the **confidence level** and $t_{\alpha/2, n-2}$ denotes the $1 - \alpha/2$ percentile of a student's t distribution with $n - 2$ degrees of freedom

- Similarly, we can construct a CI for $\beta_0$:

$$\left[\hat{\beta}_0 - t_{\alpha/2, n-2}\hat{SE}(\hat{\beta}_0), \hat{\beta}_0 + t_{\alpha/2, n-2}\hat{SE}(\hat{\beta}_0)\right]$$

CLEMS⬤N
U N I V E R S I T Y

Simple Linear Regression

Parameter Estimation

Residual Analysis

Confidence/Prediction Intervals

Hypothesis Testing

1.17

**Notes**

---

## Confidence Interval of $E(Y_{new})$

- We often interested in estimating the **mean** response for an unobserved predictor value, say, $x_{new}$. Therefore we would like to construct CI for $E[Y_{new}]$, the corresponding **mean response**

- We need sampling distribution of $\widehat{E(Y_{new})}$ to form CI:

  - $\frac{\widehat{E(Y_{new})} - E(Y_{new})}{\hat{SE}(E(\widehat{Y_{new}}))} \sim t_{n-2}, \quad \hat{SE}(E(\widehat{Y_{new}})) = \hat{\sigma}\sqrt{\left(\frac{1}{n} + \frac{(x_{new} - \bar{x})^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2}\right)}$

  - CI:

$$\left[\hat{Y}_{new} - t_{\alpha/2, n-2}\hat{SE}(E(\widehat{Y_{new}})), \hat{Y}_{new} + t_{\alpha/2, n-2}\hat{SE}(E(\widehat{Y_{new}}))\right]$$

- **Quiz:** Use this formula to construct CI for $\beta_0$

CLEMS⬤N
U N I V E R S I T Y

Simple Linear Regression

Parameter Estimation

Residual Analysis

Confidence/Prediction Intervals

Hypothesis Testing

1.18

**Notes**

**Prediction Interval of $Y_{new}$**

- Suppose we want to predict the response of a future observation $Y_{new}$ given $x = x_{new}$

- We need to account for added variability as a new observation does not fall directly on the regression line (i.e., $Y_{new} = \mathrm{E}[Y_{new}] + \varepsilon_{new}$)

- Replace $\hat{SE}(\widehat{\mathrm{E}(Y_{new})})$ by
$\hat{SE}(\hat{Y}_{new}) = \hat{\sigma}\sqrt{\left(1 + \frac{1}{n} + \frac{(x_{new}-\bar{x})^2}{\sum_{i=1}^{n}(x_i-\bar{x})^2}\right)}$ to construct CIs for $Y_{new}$

Notes

1.19

---

**Maximum Heart Rate vs. Age Revisited**

The maximum heart rate `MaxHeartRate` (HR$_{max}$) of a person is often said to be related to age `Age` by the equation:
$$\mathrm{HR}_{max} = 220 - \mathrm{Age}.$$

Suppose we have 15 people of varying ages are tested for their maximum heart rate (bpm)

| Age | 18 | 23 | 25 | 35 | 65 | 54 | 34 | 56 | 72 | 19 | 23 | 42 | 18 | 39 | 37 |
|-----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| HR$_{max}$ | 202 | 186 | 187 | 180 | 156 | 169 | 174 | 172 | 153 | 199 | 193 | 174 | 198 | 183 | 178 |

- Construct the 95% CI for $\beta_1$

- Compute the estimate for mean `MaxHeartRate` given `Age` $= 40$ and construct the associated 90% CI

- Construct the prediction interval for a new observation given `Age` $= 40$

Notes

1.20

---

**Maximum Heart Rate vs. Age: Hypothesis Test for Slope**

1. $H_0 : \beta_1 = 0$ vs. $H_a : \beta_1 \neq 0$

2. Compute the **test statistic**:
$t^* = \frac{\hat{\beta}_1 - 0}{\hat{SE}(\hat{\beta}_1)} = \frac{-0.7977}{0.06996} = -11.40$

3. Compute **P-value**: $\mathrm{P}(|t^*| \geq |t_{obs}|) = 3.85 \times 10^{-8}$

4. Compare to $\alpha$ and draw conclusion:

> Reject $H_0$ at $\alpha = .05$ level, evidence suggests a negative linear relationship between `MaxHeartRate` and `Age`

Notes

1.21

**Maximum Heart Rate vs. Age: Hypothesis Test for Intercept**

CLEMS☘N
U N I V E R S I T Y

Simple Linear
Regression

Parameter
Estimation

Residual Analysis

Confidence/Prediction
Intervals

Hypothesis Testing

1. $H_0 : \beta_0 = 0$ vs. $H_a : \beta_0 \neq 0$

2. Compute the **test statistic**:
   $t^* = \frac{\hat{\beta}_0 - 0}{\hat{SE}(\hat{\beta}_0)} = \frac{210.0485}{2.86694} = 73.27$

3. Compute **P-value**: $P(|t^*| \geq |t_{obs}|) \simeq 0$

4. Compare to $\alpha$ and draw conclusion:

   > Reject $H_0$ at $\alpha = .05$ level, evidence suggests
   > evidence suggests the intercept (the expected
   > `MaxHeartRate` at age 0) is different from $0$

Notes

---

**Summary**

CLEMS☘N
U N I V E R S I T Y

Simple Linear
Regression

Parameter
Estimation

Residual Analysis

Confidence/Prediction
Intervals

Hypothesis Testing

In this lecture, we reviewed

- Simple Linear Regression: $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$

- Method of Least Squares for parameter estimation

- Residual analysis to check model assumptions

- Confidence/Prediction Intervals and Hypothesis Testing

Notes

Notes