

STAT 500

Simple Linear Regression: Model and Estimates

Research Questions

- Study the relationship of two or more quantitative variables.
 - quantitative: number, usually continuous
 - qualitative: classes, identify groups
- Is there a significant linear relationship between the response variable and the explanatory variable?
- What mean value of response would we predict for a given value of the explanatory variable?
- What value of response would we predict for a given value of the explanatory variable?

Simple Linear Regression

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad i = 1, \dots, n$$

- $i = 1, \dots, n$ is the number of observations
- Y_i is the *response* or dependent variable
- X_i is the predictor, *explanatory variable*, or independent variable, treated as known and fixed
- ϵ_i is the *random error* term representing individual variation and measurement error

Linear Model Notation

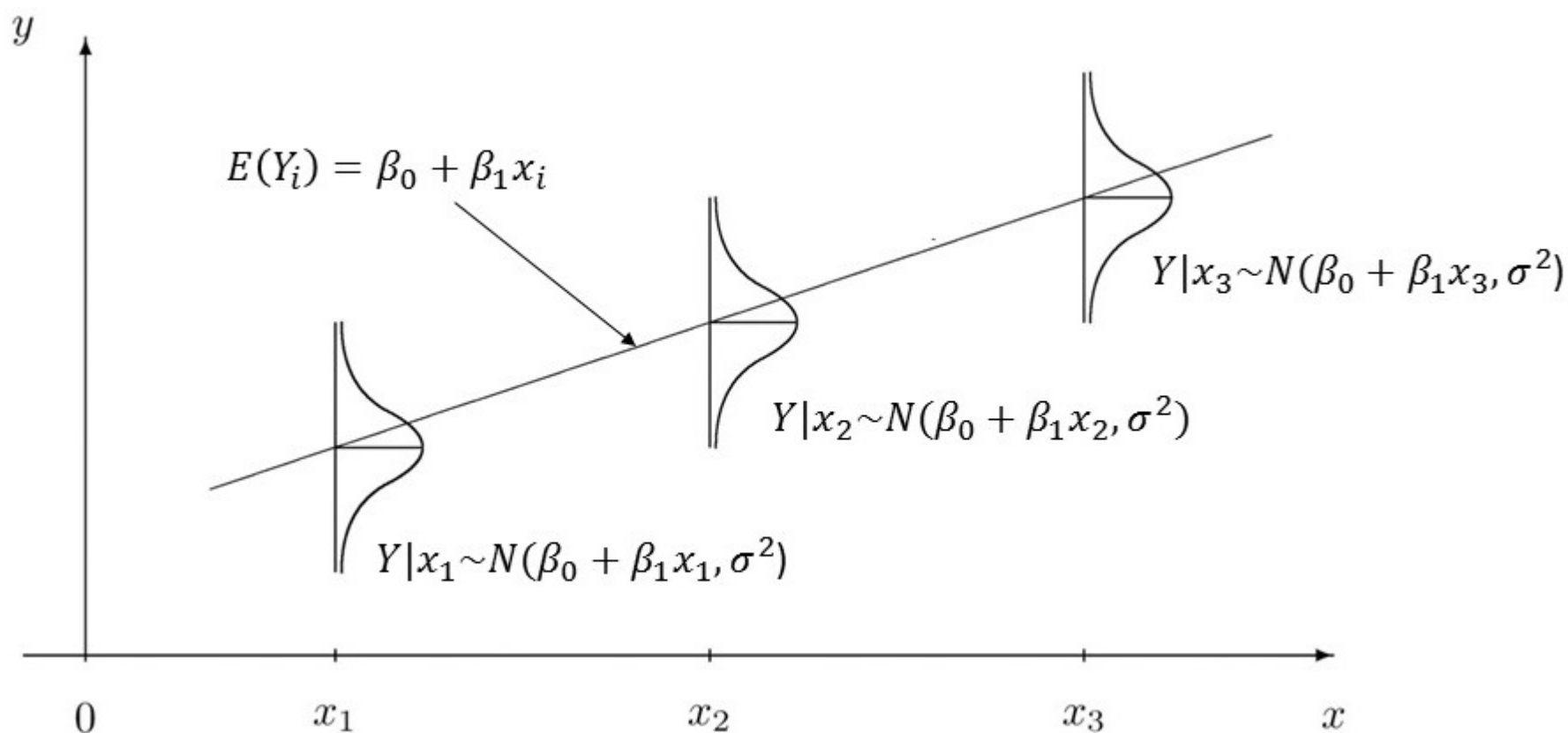
Write SLR model as a linear model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ where

$$\begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \\ \vdots \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \vdots & \vdots \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Model Assumptions

- x 's are fixed (or conditioned upon)
- The expected response is a linear function of the explanatory variable : $E(Y_i|X_i = x_i) = \beta_0 + \beta_1 x_i$
- additive random errors $Y_i = E(Y_i|X_i = x_i) + \epsilon_i$
- independent (uncorrelated) random errors
- homogeneous error variance: $Var(\epsilon_i) = \sigma^2$
- normally distributed random errors: $\epsilon_i \sim N(0, \sigma^2)$

Model and Assumptions



Model and Assumptions

The conditional distribution of Y given that $X = x$ is

$$N(\beta_o + \beta_1 x, \sigma^2)$$

- β_1 = slope, is the change in the conditional mean of Y for a one unit increase in x
- β_o is the conditional mean of Y when $X = 0$
- If we replace x by $x - x_o$ to obtain $Y = \beta_o + \beta_1(x - x_o) + \epsilon$, then β_o is the conditional mean of Y when $X = x_o$
- σ^2 is the variation of responses about the conditional mean for any specific value of the explanatory variable

Relationship to ANOVA

- ANOVA: each group (each level of explanatory variable) has its own mean
- Each x_i in regression defines its own group, but...
 - too many groups with too few observations per group
 - Linear regression analysis makes stronger assumption about the means (linear structure)

A bit of history

Sir Francis Galton coined the term “regression”.

- biometrician, geneticist, 1870-1920s
- compared the heights of children to their parents
- parents and children had similar means
- short parents had short children, tall parents had tall children
- children were closer to average than their parents
- “regression” to the mean

Least Squares Estimation

Use data $(Y_i, x_i), i = 1, 2, \dots, n$ to estimate the regression coefficients in the model

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i \text{ where } \epsilon_i \sim N(0, \sigma^2)$$

- Choose estimates b_o and b_1 to minimize

$$g(b_o, b_1) = \sum_{i=1}^n [Y_i - (b_o + b_1 x_i)]^2$$

- Why squared errors?
 - Tradition (Gauss invented least squares estimation)
 - Equivalent to maximum likelihood estimation when errors are independent and normally distributed with constant variance

Least Squares Estimation

Results:

$$b_o = \bar{Y} - b_1 \bar{x}$$

$$b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})Y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

- These are best linear unbiased estimators (blue)
- Predicted (fitted) values: $\hat{Y}_i = b_o + b_1 x_i$
- Residuals: $e_i = Y_i - \hat{Y}_i$

Least Squares Estimation

- Choose b_o, b_1 to minimize

$$g(b_o, b_1) = \sum_{i=1}^n (Y_i - (b_o + b_1 x_i))^2$$

- Taking derivatives and setting them equal to zero yields the normal equations

$$b_o n + b_1 \sum x_i = \sum Y_i$$

$$b_o \sum x_i + b_1 \sum x_i^2 = \sum x_i Y_i$$

- The normal equations can also be written as

$$\sum e_i = \sum (Y_i - (b_o + b_1 x_i)) = 0$$

$$\sum x_i e_i = \sum x_i (Y_i - (b_o + b_1 x_i)) = 0$$

Least Squares Estimation

Normal equations can be written in matrix form

$$\begin{bmatrix} b_0 n + b_1 \sum x_i \\ b_0 \sum x_i + b_1 \sum x_i^2 \end{bmatrix} = \begin{bmatrix} \sum Y_i \\ \sum x_i Y_i \end{bmatrix}$$

that is equivalent to

$$\begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \end{bmatrix} = \begin{bmatrix} \sum Y_i \\ \sum x_i Y_i \end{bmatrix}$$

and can be written as $X^T X \mathbf{b} = X^T \mathbf{Y}$

$$\text{where } \mathbf{b} = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

Least Squares Estimation

Solution to the normal equations

$$\begin{bmatrix} b_o \\ b_1 \end{bmatrix} = (X^T X)^{-1} X^T Y = \begin{bmatrix} \bar{Y} - b_1 \bar{x} \\ \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{bmatrix}$$

Least Squares Estimation

Variance-covariance matrix of the least squares estimator:

$$\begin{aligned} \text{Var} \begin{bmatrix} b_o \\ b_1 \end{bmatrix} &= \begin{bmatrix} \text{Var}(b_o) & \text{Cov}(b_o, b_1) \\ \text{Cov}(b_o, b_1) & \text{Var}(b_1) \end{bmatrix} \\ &= \text{Var} \left((X^T X)^{-1} X^T Y \right) \\ &= (X^T X)^{-1} X^T \text{Var}(Y) \left[(X^T X)^{-1} X^T \right]^T \\ &= (X^T X)^{-1} X^T \begin{bmatrix} \sigma^2 I \end{bmatrix} X (X^T X)^{-1} \\ &= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} \\ &= \sigma^2 (X^T X)^{-1} \end{aligned}$$

Least Squares Estimation

The the variance-covariance matrix of the least squares estimator for the regression coefficients is

$$\begin{aligned} \text{Var} \begin{bmatrix} b_o \\ b_1 \end{bmatrix} &= \begin{bmatrix} \text{Var}(b_o) & \text{Cov}(b_o, b_1) \\ \text{Cov}(b_o, b_1) & \text{Var}(b_1) \end{bmatrix} \\ &= \sigma^2 (X^T X)^{-1} \\ &= \sigma^2 \begin{bmatrix} \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} & \frac{-\bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ \frac{-\bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} & \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{bmatrix} \end{aligned}$$

Least Squares Estimation

- Matrix of second partial derivatives of $g(b_o, b_1)$

$$\begin{bmatrix} \frac{\partial^2 g(b_o, b_1)}{\partial b_o^2} & \frac{\partial^2 g(b_o, b_1)}{\partial b_o \partial b_1} \\ \frac{\partial^2 g(b_o, b_1)}{\partial b_o \partial b_1} & \frac{\partial^2 g(b_o, b_1)}{\partial b_1^2} \end{bmatrix} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} = X^T X$$

Since this matrix is positive definite (if we have at least two different x_i values), it guarantees we have a minimum.

- Note: least squares estimate of the slope is different if there is no intercept in the model

Multivariate Normal Distribution

Definition: Suppose $\mathbf{Z} = \begin{bmatrix} Z_1 \\ \vdots \\ Z_m \end{bmatrix}$ is a

random vector whose elements are independently distributed standard normal random variables.

For any $n \times m$ matrix \mathbf{A} , We say that

$$\mathbf{Y} = \boldsymbol{\mu} + \mathbf{AZ}$$

has a *multivariate normal distribution* with mean vector

$$\mathbf{E}(\mathbf{Y}) = \mathbf{E}(\boldsymbol{\mu} + \mathbf{AZ}) = \boldsymbol{\mu} + \mathbf{AE}(\mathbf{Z}) = \boldsymbol{\mu} + \mathbf{A}\mathbf{0} = \boldsymbol{\mu}$$

and variance-covariance matrix

$$\mathbf{Var}(\mathbf{Y}) = \mathbf{A}[\mathbf{Var}(\mathbf{Z})]\mathbf{A}^T = \mathbf{AA}^T \equiv \boldsymbol{\Sigma}$$

Multivariate Normal Distribution

We will use the notation

$$\mathbf{Y} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

When $\boldsymbol{\Sigma}$ is positive definite, the joint density function is

$$f(\mathbf{y}) = \frac{1}{(2\pi)^{n/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{y}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{y}-\boldsymbol{\mu})}$$

where

$$\boldsymbol{\Sigma} = \begin{bmatrix} \text{Var}(\mathbf{Y}_1) & \text{Cov}(\mathbf{Y}_1, \mathbf{Y}_2) & \text{Cov}(\mathbf{Y}_1, \mathbf{Y}_3) & \cdots & \text{Cov}(\mathbf{Y}_1, \mathbf{Y}_n) \\ \text{Cov}(\mathbf{Y}_2, \mathbf{Y}_1) & \text{Var}(\mathbf{Y}_2) & \text{Cov}(\mathbf{Y}_2, \mathbf{Y}_3) & \cdots & \text{Cov}(\mathbf{Y}_2, \mathbf{Y}_n) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(\mathbf{Y}_n, \mathbf{Y}_1) & \text{Cov}(\mathbf{Y}_n, \mathbf{Y}_2) & \text{Cov}(\mathbf{Y}_n, \mathbf{Y}_3) & \cdots & \text{Var}(\mathbf{Y}_n) \end{bmatrix}$$

Multivariate Normal Distribution

The multivariate normal distribution has some useful properties. One is that normality is preserved under linear transformations:

If $Y \sim N(\mu, \Sigma)$, then

$$W = c + BY \sim N(c + B\mu, B\Sigma B^T)$$

for any non-random c and B .

Forbes Data

Weisberg, Sanford, *Applied Linear Regression*, Wiley, 1980.

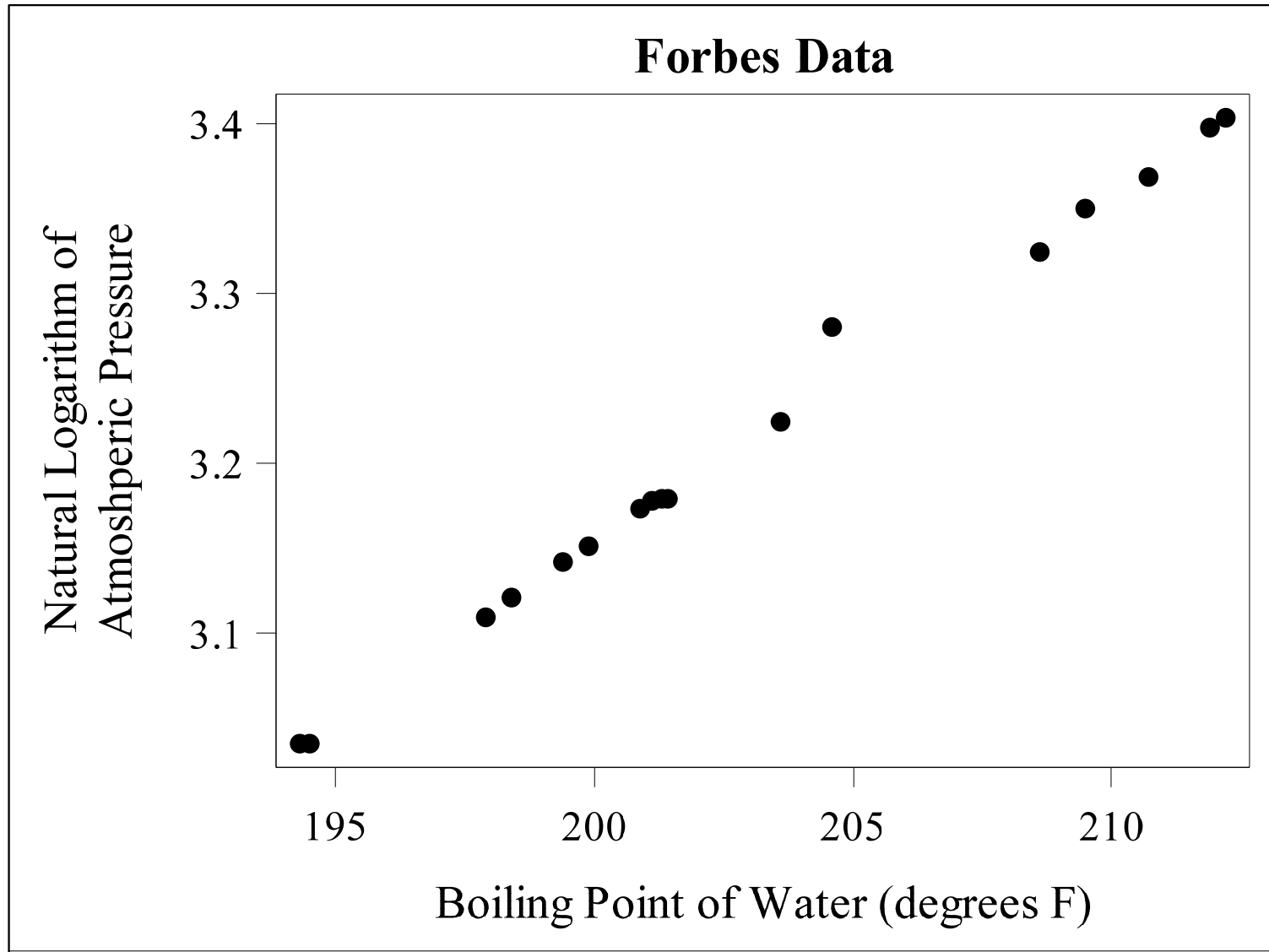
- James D. Forbes collected data in the mountains of Scotland
- $n=17$ locations (at different altitudes)
- Objective: Predict barometric pressure (in inches of mercury) from boiling point of water (X) in $^{\circ}\text{F}$.
- Use $Y=\log(\text{barometric pressure})$
- Motivation: Fragile barometers were difficult to transport

Forbes Data

	BOILING POINT OF WATER	BARAMETRIC PRESSURE	NATURAL LOG OF BARAMETRIC PRESSURE
Obs	(degrees F)	(inches Hg)	
1	194.3	20.79	3.03447
2	194.5	20.79	3.03447
3	197.9	22.40	3.10906
4	198.4	22.67	3.12104
5	199.4	23.15	3.14199
6	199.9	23.35	3.15060
7	200.9	23.89	3.17346
8	201.1	23.99	3.17764

Forbes Data

	BOILING POINT OF WATER	BARAMETRIC PRESSURE	NATURAL LOG OF BARAMETRIC PRESSURE
Obs	(degrees F)	(inches Hg)	
9	201.3	24.01	3.17847
10	201.4	24.02	3.17889
11	203.6	25.14	3.22446
12	204.6	26.57	3.27978
13	208.6	27.76	3.32360
14	209.5	28.49	3.34955
15	210.7	29.04	3.36867
16	211.9	29.88	3.39719
17	212.2	30.06	3.40320



Analysis of the Forbes Data

- Proposed regression model

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

where $\epsilon_i \sim NID(0, \sigma^2)$, $i = 1, 2, \dots, 17$

- $Y_i = \log(\text{pressure})$
- $X_i = \text{boiling point } (^{\circ}\text{F})$
- β_1 is the increase in mean $\log(\text{pressure})$ when boiling point of water increases by 1 $^{\circ}\text{F}$
- β_0 is the mean $\log(\text{pressure})$ when boiling point of water is 0 $^{\circ}\text{F}$ (Is this extrapolation realistic?)

Predicted Values and Residuals

- Predicted (fitted) values

$$\hat{Y}_i = b_0 + b_1 x_i$$

$$\hat{Y} = X\hat{\beta}$$

- Residuals

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

Analysis of the Forbes Data

- Estimated regression model

$$\hat{Y} = b_0 + b_1x = -0.97097 + 0.020623x$$

- Could have subtracted 212 °F from each boiling point. Then the estimated model is

$$\hat{Y} = b_0 + 212b_1 + b_1(x - 212)$$

$$= 3.401106 + 0.020623(x - 212)$$

- Then 3.401106 is an estimate of the mean log(pressure) at 212 °F.

Predicted Values: Example

$$\hat{Y}_i = -0.97097 + 0.020623x$$

- Values on the estimated regression line.
- Predict values of Y_i for a given value of x_i
 - $x_i = 201.1$ °F:
 $-0.97097 + 0.020623(201.1) = 3.176315$
 - $x_i = 210.7$ °F:
 $-0.97097 + 0.020623(210.7) = 3.374296$

Residuals: Example

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

- Vertical distance between observed value of Y and predicted value of Y .
- Residuals:
 - $x_i = 201.1$ °F and $Y_i = 3.17764$: $3.17764 - 3.176315 = 0.001325$
 - $x_i = 210.7$ °F and $Y_i = 3.36867$: $3.36867 - 3.374296 = -0.005626$