

>

# Linear Regression

# CI Birds

Wing span      cm/in

Beak size      mm/cm

Color      Categorical

Body mass      g

Wing span  $\rightarrow$  body mass

$X \sim \text{Wingspan} \sim N(\mu, \sigma^2)$

$Y \sim \text{Body Mass} \sim N(\mu, \sigma^2)$

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

$$\epsilon \sim N(0, \sigma^2)$$

# Learning Outcomes

- Scatter Plot
- Linear Regression
- Ordinary Least Squares
- Unbiasedness

# Scatter Plot

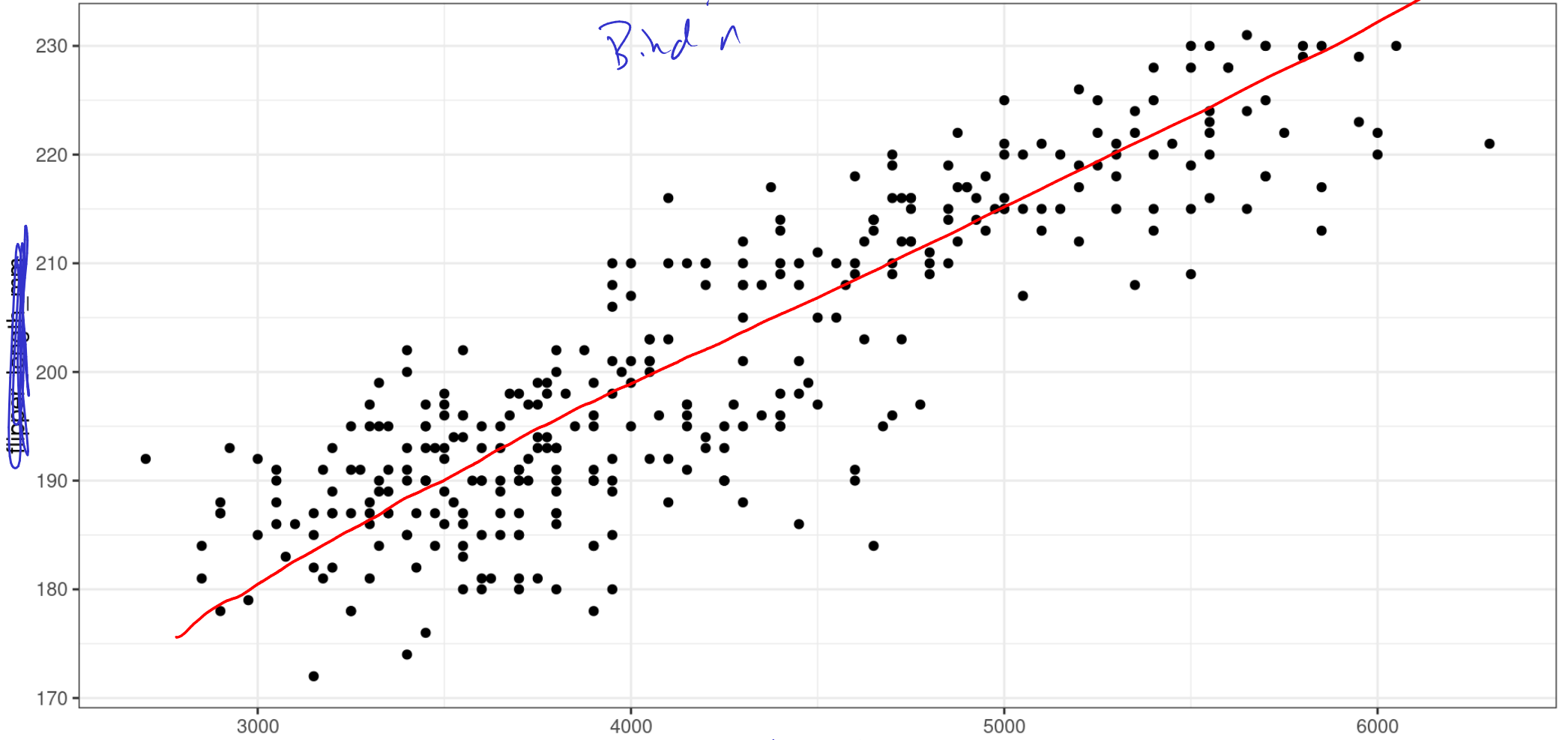
# Scatter Plot

Bird 1 ( $x_1, y_1$ )

$$y = mx + b$$

Bird n

X

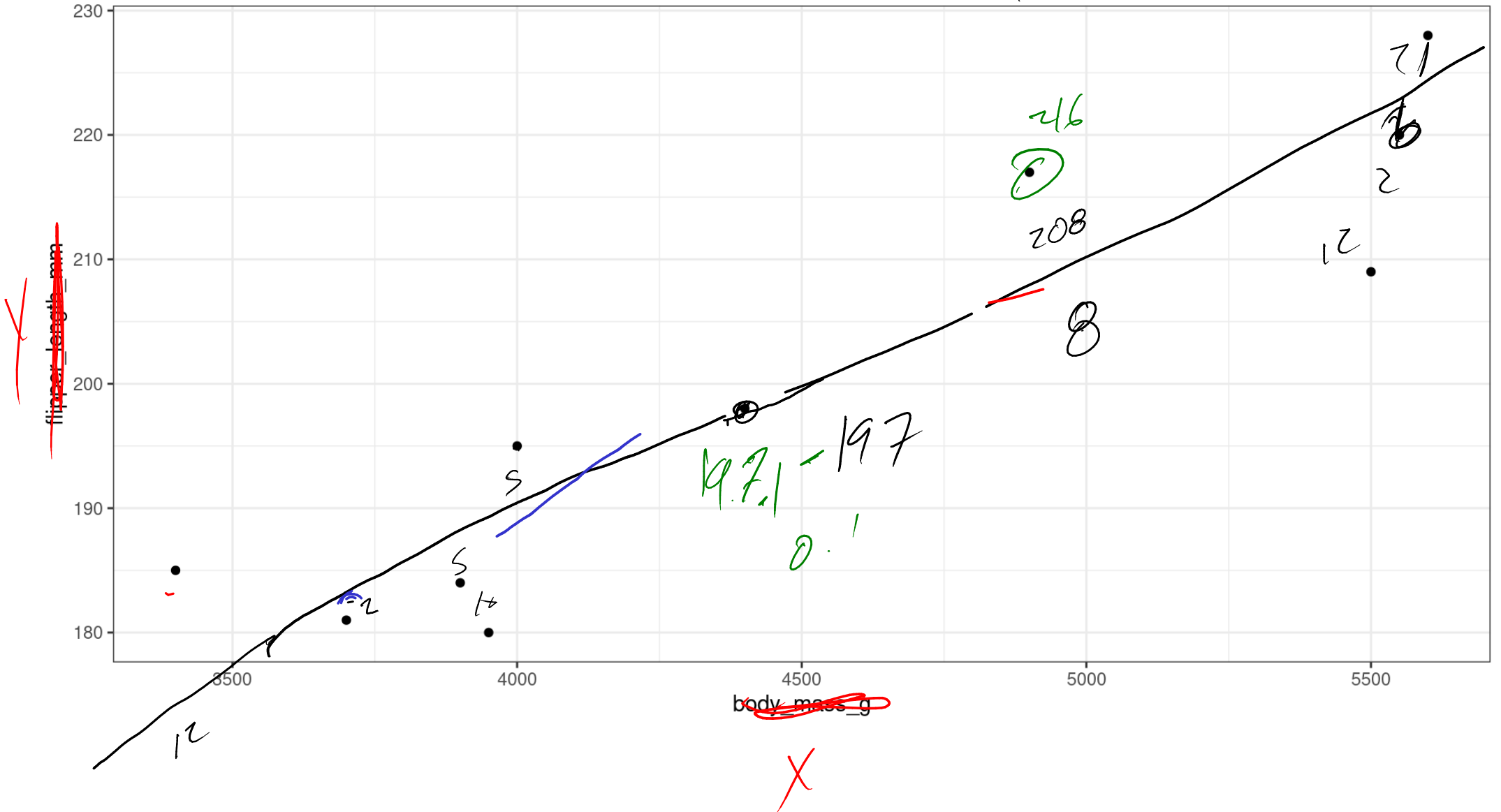


~~body mass~~  
X

X↑    Y↑  
X↓    Y↓

# Scatter Plot

$$y = mx + b$$



# Linear Regression



# Linear Regression

Linear regression is used to model the association between a set of predictor variables ( $x$ 's) and an outcome variable ( $y$ ). Linear regression will fit a line that best describes the data points.

# Simple Linear Regression

Simple linear regression will model the association between one predictor variable and an outcome:

$$\hat{Y} = \beta_0 + \beta_1 X + \epsilon$$

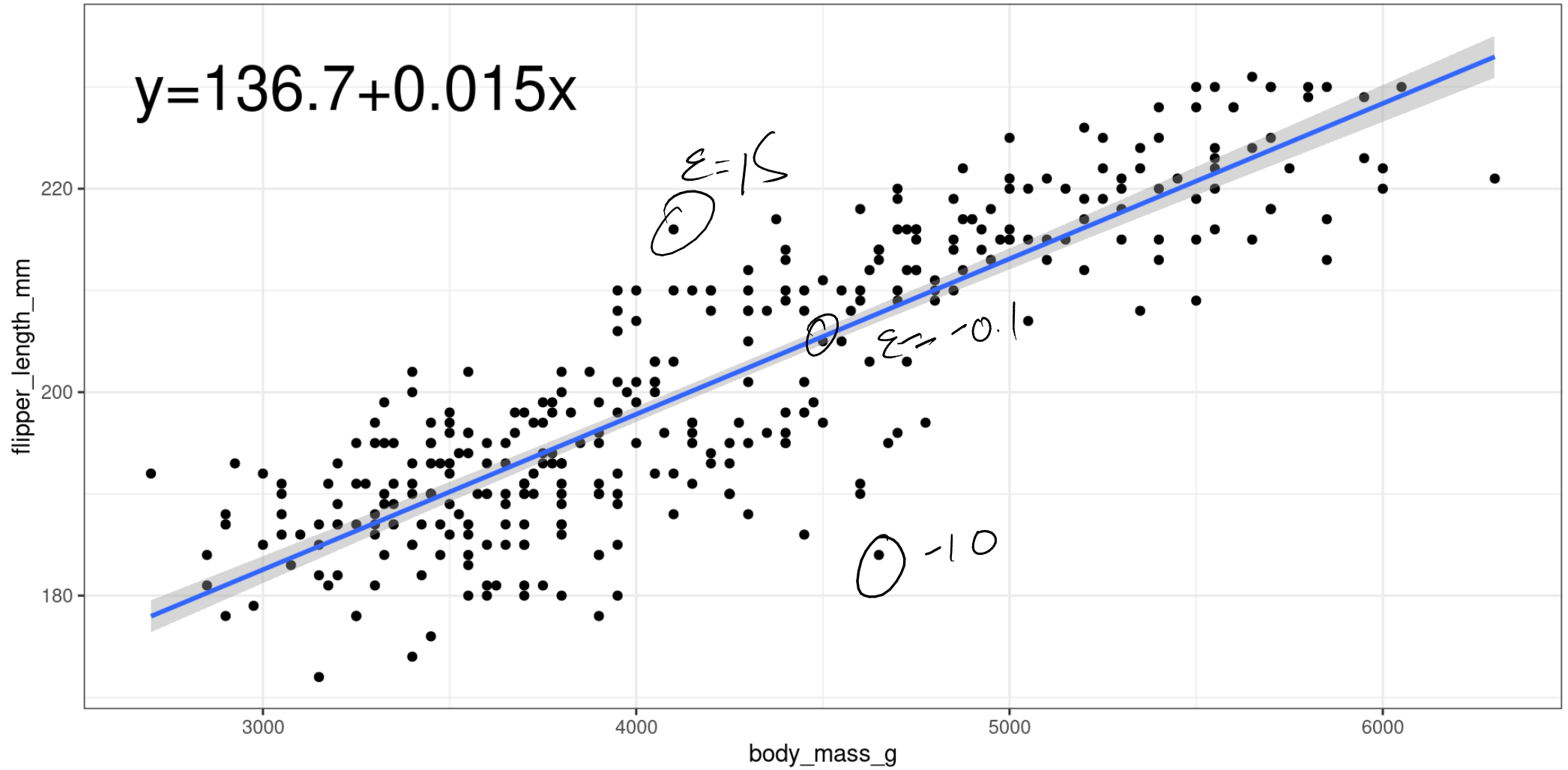
*intercept*  $\beta_0$  *slope*  $\beta_1 X$  *error term*  $\epsilon$

- $\beta_0$ : Intercept term
- $\beta_1$ : Slope term
- $\epsilon \sim N(0, \sigma^2)$



# Fitting a Line

$$y = \beta_0 + \beta_1 x$$



# Interpretation

$$\hat{y} = 136.73 + 0.015x$$

Body mass is 0, the flipper length is 136.73  
X Y

As X increase 1 unit, Y increase by an  
average of 0.015 units

# Ordinary Least Squares

# Ordinary Least Squares

For a data pair  $(X_i, Y_i)_{i=1}^n$ , the ordinary least squares estimator will find the estimates of  $\hat{\beta}_0$  and  $\hat{\beta}_1$  that minimize the following function:

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

# Estimating $\beta$ 's

$$\frac{d \ell}{d \beta_0}$$

$$\frac{d \ell}{d \beta_1}$$

# Estimating $\beta_1$

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

$$\frac{dL}{d\beta_1} = \sum_{i=1}^n 2(y_i - (\beta_0 + \beta_1 x_i)) (-x_i)$$

$$= -2 \sum (y_i - \beta_0 - \beta_1 x_i) x_i$$

$$\bar{x} = \frac{\sum x_i}{n}$$

$$= -2 \sum (x_i y_i - \beta_0 x_i - \beta_1 x_i^2)$$

$$= -2 \left[ \sum x_i y_i - \beta_0 n \bar{x} - \beta_1 \sum x_i^2 \right]$$

$$0 = -2 \left[ \sum x_i y_i - \beta_0 n \bar{x} - \beta_1 \sum x_i^2 \right]$$



$$0 = \sum X_i Y_i - \beta_0 n \bar{x} - \beta_1 \sum X_i^2$$

$$\frac{\beta_1 \sum X_i^2}{\sum X_i^2} = \frac{\sum X_i Y_i - \beta_0 n \bar{x}}{\sum X_i^2}$$

$$\hat{\beta}_1 = \frac{\sum X_i Y_i - \beta_0 n \bar{x}}{\sum X_i^2}$$

$$\hat{\beta}_1 = \frac{\sum (X_i - \bar{x})(Y_i - \bar{y})}{\sum (X_i - \bar{x})^2}$$



# Estimating $\beta_0$

# Estimates

intercept

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

slope

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

Model  
error

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

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

$$\varepsilon \sim \mathcal{N}(0, \sigma^2)$$


**Unbiasedness of  $\beta$ 's**

# Unbiasedness of $\beta$ 's

Both  $\beta_0$  and  $\beta_1$  are unbiased estimators.

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

$$\hat{\beta}_1 = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

$$\varepsilon \sim N(0, \sigma^2)$$

$$E(\varepsilon) = 0$$

$$y = \beta_0 + \beta_1 x + \varepsilon$$

$$E(y) = \beta_0 + \beta_1 x$$

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

$$\text{Var}(\beta_0 + \beta_1 x + \epsilon)$$

$$\text{Var}(\beta_0) + \text{Var}(\beta_1 x) + \text{Var}(\epsilon)$$

$$0 + 0 + \sigma^2$$

$$E(\hat{\beta}_0)$$

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

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

$$\frac{\sum x_i}{n} = \bar{x}$$

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

$$= E(\bar{y}) - \bar{x} E(\hat{\beta}_1)$$

$$E(y_i) = \beta_0 + \beta_1 x_i$$

$$E(\bar{y}) - \bar{x} \beta_1$$

$$E\left(\frac{1}{n} \sum y_i\right) - \bar{x} \beta_1$$

$$\frac{1}{n} \sum E(y_i) - \bar{x} \beta_1$$

$$\frac{1}{n} \sum (\beta_0 + \beta_1 x_i) - \bar{x} \beta_1$$



$$\frac{1}{n} [n\beta_0 + \beta_1 \sum x_i] - \bar{x} \beta_1$$

$$\beta_0 + \cancel{\frac{\beta_1}{n} \sum x_i} - \cancel{\bar{x} \beta_1}$$

$$\beta_0$$

$$E(\hat{\beta}_1) \quad \hat{\beta}_1 = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

$$E(\hat{\beta}_1) = E\left(\frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}\right)$$

$$\frac{1}{\sum (x_i - \bar{x})^2} E\left(\sum (y_i - \bar{y})(x_i - \bar{x})\right)$$

$$\frac{1}{\sum (x_i - \bar{x})^2} \sum_{i=1}^n E\left[(y_i - \bar{y})(x_i - \bar{x})\right]$$

$$\frac{1}{\sum (x_i - \bar{x})^2} \sum (x_i - \bar{x}) E(y_i - \bar{y})$$

$$\frac{1}{\sum (x_i - \bar{x})^2} \sum (x_i - \bar{x}) (E(y_i) - E(\bar{y}))$$

$$\frac{1}{\sum (x_i - \bar{x})^2} \sum (x_i - \bar{x}) (\cancel{\beta_0} + \beta_1 x_i - \cancel{\beta_0} - \beta_1 \bar{x})$$

$$\frac{1}{\sum (x_i - \bar{x})^2} \sum (x_i - \bar{x}) (\beta_1 x_i - \beta_1 \bar{x})$$

$$\frac{1}{\sum (x_i - \bar{x})^2} \sum (x_i - \bar{x}) \beta_1 (x_i - \bar{x})$$

$$\cancel{\frac{1}{\sum (x_i - \bar{x})^2}} \beta_1 \sum (x_i - \bar{x})^2$$

$$= \beta_1$$