# STAT 423/523 Statistical Methods for Engineers and Scientists

## Lecture 10: Simple Linear Regression

Chencheng Cai

Washington State University

# Simple Linear Regression

**Regression** is a statistical method for estimating the relationships among variables. THe simpest form of regression is **simple linear regression**:

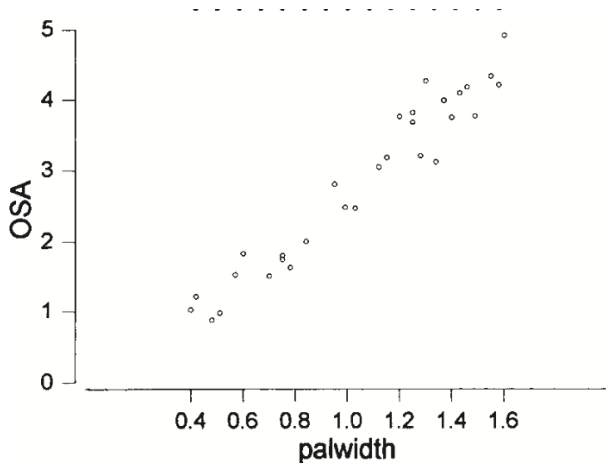$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i.$$

# Simple Linear Regression

**Regression** is a statistical method for estimating the relationships among variables. THe simpest form of regression is **simple linear regression**:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i.$$

- $y_i$ is the response variable (dependent variable).
- $x_i$ is the predictor variable (independent variable).
- $\beta_0$ is the intercept.
- $\beta_1$ is the slope.
- $\epsilon_i$ is the error term.

## Example

▶ $y$: ocular surface area

▶ $x$: width of the palprebal fissure

# Assumptions

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i.$$

- ► Linearity: The relationship between $x$ and $y$ is linear.
- ► Independence: The errors are independent.
- ► Normality: The errors are normally distributed.
- ► Equal variance: The errors have constant variance.

# Assumptions

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i.$$

- ▶ Linearity: The relationship between $x$ and $y$ is linear.
- ▶ Independence: The errors are independent.
- ▶ Normality: The errors are normally distributed.
- ▶ Equal variance: The errors have constant variance.

For short, the LINE assumptions give:

$$y_i = \beta_0 + \beta_1 x_i + N(0, \sigma^2) \quad \forall i$$

# Violations of Assumptions

- Linearity: Nonliear regression model.
- Independence: Structural equation model (SEM) in econometrics.
- Normality: $\epsilon_i$ could have a heavy-tailed distribution.
- Equal variance: Heteroscedasticity.

## Some Statistics

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i.$$

We assume all $x_i$'s are fixed and known. (not random variables!)

## Some Statistics

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i.$$

We assume all $x_i$'s are fixed and known. (not random variables!)

- $E(y_i) = \beta_0 + \beta_1 x_i$ is the mean response for a given $x_i$.

# Some Statistics

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i.$$

We assume all $x_i$'s are fixed and known. (not random variables!)

- $E(y_i) = \beta_0 + \beta_1 x_i$ is the mean response for a given $x_i$.
- $Var(y_i) = Var(\epsilon_i) = \sigma^2$ is the variance of the response.

## Some Statistics

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i.$$

We assume all $x_i$'s are fixed and known. (not random variables!)

- $E(y_i) = \beta_0 + \beta_1 x_i$ is the mean response for a given $x_i$.
- $Var(y_i) = Var(\epsilon_i) = \sigma^2$ is the variance of the response.
- $Cov(y_i, y_j) = Cov(\epsilon_i, \epsilon_j)$ for $i \neq j$. (Independence Assumption).

## Some Statistics

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i.$$

We assume all $x_i$'s are fixed and known. (not random variables!)

- $E(y_i) = \beta_0 + \beta_1 x_i$ is the mean response for a given $x_i$.
- $Var(y_i) = Var(\epsilon_i) = \sigma^2$ is the variance of the response.
- $Cov(y_i, y_j) = Cov(\epsilon_i, \epsilon_j)$ for $i \neq j$. (Independence Assumption).

If we get the estimated coefficients $\hat{\beta}_0$ and $\hat{\beta}_1$,

## Some Statistics

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i.$$

We assume all $x_i$'s are fixed and known. (not random variables!)

- $E(y_i) = \beta_0 + \beta_1 x_i$ is the mean response for a given $x_i$.
- $Var(y_i) = Var(\epsilon_i) = \sigma^2$ is the variance of the response.
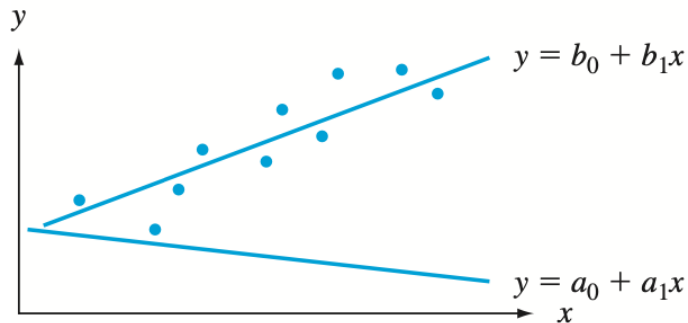- $Cov(y_i, y_j) = Cov(\epsilon_i, \epsilon_j)$ for $i \neq j$. (Independence Assumption).

If we get the estimated coefficients $\hat{\beta}_0$ and $\hat{\beta}_1$,

- The **fitted value** for $y_i$ is $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$.
- The **residual** for $y_i$ is $\hat{\epsilon}_i = y_i - \hat{y}_i$.

# Estimation

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i.$$

Given the data points



we want to find the line that **best fits** the data points.

# Ordinary Least Squares

The first approach is **Ordinary Least Squares** (OLS).

# Ordinary Least Squares

The first approach is **Ordinary Least Squares** (OLS).

▶ For each possible parameter values $\beta_0$ and $\beta_1$, we can calculate the **residual sum of squares** (RSS):

$$\text{RSS}(\beta_0, \beta_1) = \sum_{i=1}^{N}(y_i - \beta_0 - \beta_1 x_i)^2$$

# Ordinary Least Squares

The first approach is **Ordinary Least Squares** (OLS).

▶ For each possible parameter values $\beta_0$ and $\beta_1$, we can calculate the **residual sum of squares** (RSS):

$$\text{RSS}(\beta_0, \beta_1) = \sum_{i=1}^{N} (y_i - \beta_0 - \beta_1 x_i)^2$$

▶ The OLS estimates are the values of $\beta_0$ and $\beta_1$ that minimize the RSS:

$$\hat{\beta}_0, \hat{\beta}_1 = \underset{\beta_0, \beta_1}{\arg\min}\ \text{RSS}(\beta_0, \beta_1)$$

# Residual Sum of Squares

The residual sum of squares is the sum of the squared distance between the data points and the fitted line.

It is the **vertical** distance, not the orthogonal distance.
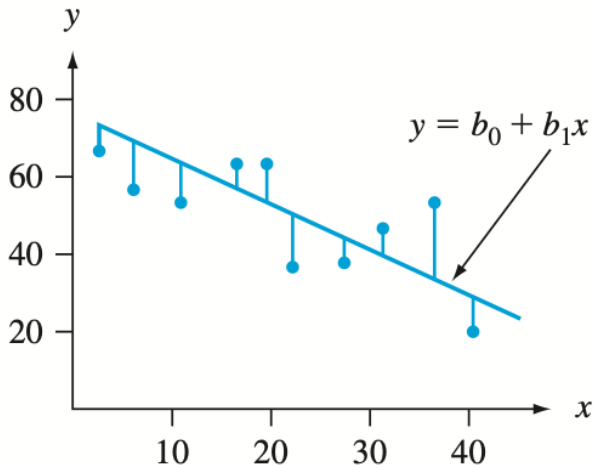
# Residual Sum of Squares

The residual sum of squares is the sum of the squared distance between the data points and the fitted line.

It is the **vertical** distance, not the orthogonal distance.

# OLS

In order to minimize the RSS, we first compute its partial derivatives.

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

$$\frac{\partial \text{RSS}}{\partial \beta_0} = -2\sum_{i=1}^{n}(y_i - \beta_0 - \beta_1 x_i) = -2\sum_{i=1}^{n} y_i + 2N\beta_0 + 2\beta_1 \sum_{i=1}^{n} x_i$$

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

To find the minimum, we set the partial derivatives to zero.

## OLS

The **estimating equations** for OLS are:

$$0 = -2\sum_{i=1}^{n} y_i + 2n\beta_0 + 2\beta_1 \sum_{i=1}^{n} x_i \tag{1}$$

$$0 = -2\sum_{i=1}^{n} y_i x_i + 2\beta_0 \sum_{i=1}^{n} x_i + 2\beta_1 \sum_{i=1}^{n} x_i^2 \tag{2}$$

## OLS

The **estimating equations** for OLS are:

$$0 = -2 \sum_{i=1}^{n} y_i + 2n\beta_0 + 2\beta_1 \sum_{i=1}^{n} x_i \tag{1}$$

$$0 = -2 \sum_{i=1}^{n} y_i x_i + 2\beta_0 \sum_{i=1}^{n} x_i + 2\beta_1 \sum_{i=1}^{n} x_i^2 \tag{2}$$

Compute $(1) \times \sum_i x_i - (2) \times n$:

$$0 = 2n \sum_i x_i y_i - 2 \sum_i x_i \sum_i y_i + \left( \left( \sum_i x_i \right)^2 - n \sum_i x_i^2 \right) \beta_1.$$

# OLS

The **estimating equations** for OLS are:

$$0 = -2\sum_{i=1}^{n} y_i + 2n\beta_0 + 2\beta_1 \sum_{i=1}^{n} x_i \qquad (1)$$

$$0 = -2\sum_{i=1}^{n} y_i x_i + 2\beta_0 \sum_{i=1}^{n} x_i + 2\beta_1 \sum_{i=1}^{n} x_i^2 \qquad (2)$$

Compute $(1) \times \sum_i x_i - (2) \times n$:

$$0 = 2n \sum_i x_i y_i - 2 \sum_i x_i \sum_i y_i + \left( \left( \sum_i x_i \right)^2 - n \sum_i x_i^2 \right) \beta_1.$$

$$\implies \hat{\beta}_1 = \frac{\sum_i x_i y_i - n^{-1} \sum_i x_i \sum_i y_i}{\sum_i x_i^2 - n^{-1} \left( \sum_i x_i \right)^2}.$$

# OLS

$$\hat{\beta}_1 = \frac{\sum_i x_i y_i - n^{-1} \sum_i x_i \sum_i y_i}{\sum_i x_i^2 - n^{-1} \left(\sum_i x_i\right)^2}.$$

▶ The numerator is

$$\sum_i x_i y_i - n^{-1} \sum_i x_i \sum_i y_i = S_{xy} = \sum_i (y_i - \bar{y})(x_i - \bar{x})$$

▶ The denominator is

$$\sum_i x_i^2 - n^{-1} \left(\sum_i x_i\right)^2 = S_{xx} = \sum_i (x_i - \bar{x})^2$$

## OLS

Therefore,

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\text{Cov}(X,Y)}{\text{Var}(X)}$$

with

$$S_{xy} = \sum_i (y_i - \bar{y})(x_i - \bar{x}) = \sum_i y_i x_i - n^{-1} \sum_i x_i \sum_i y_i$$

$$S_{xx} = \sum_i (x_i - \bar{x})^2 = \sum_i x_i^2 - n^{-1} \left( \sum_i x_i \right)^2$$

# OLS

Therefore,

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\text{Cov}(X,Y)}{\text{Var}(X)}$$

with

$$S_{xy} = \sum_i (y_i - \bar{y})(x_i - \bar{x}) = \sum_i y_i x_i - n^{-1} \sum_i x_i \sum_i y_i$$

$$S_{xx} = \sum_i (x_i - \bar{x})^2 = \sum_i x_i^2 - n^{-1} \left( \sum_i x_i \right)^2$$

From Eq. (1), we can get $\hat{\beta}_0$:

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

# OLS

We still have $\sigma^2$ to estimate.

## OLS

We still have $\sigma^2$ to estimate. The easiest way is to estimate it from the residual sum of squares:

$$\hat{\sigma}^2 = \frac{\text{RSS}(\hat{\beta}_0, \hat{\beta}_1)}{n - 2}$$

▶ $n - 2$ is the degrees of freedom.

## OLS

We still have $\sigma^2$ to estimate. The easiest way is to estimate it from the residual sum of squares:

$$\hat{\sigma}^2 = \frac{\text{RSS}(\hat{\beta}_0, \hat{\beta}_1)}{n-2}$$

- $n - 2$ is the degrees of freedom.

A quick formula in computing $\text{RSS}(\hat{\beta}_0, \hat{\beta}_1)$ is

$$\text{RSS}(\hat{\beta}_0, \hat{\beta}_1) = S_{yy} - \hat{\beta}_1 S_{xy} = S_{yy} - \hat{\beta}_1^2 S_{xx},$$

where

$$S_{yy} = \sum_i (y_i - \bar{y})^2 = \sum_i y_i^2 - n^{-1} \left( \sum_i y_i \right)^2.$$

# OLS

Summary for OLS estimators:

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

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

$$\hat{\sigma}^2 = \frac{\text{RSS}(\hat{\beta}_0, \hat{\beta}_1)}{n - 2} = \frac{S_{yy} - \hat{\beta}_1 S_{xy}}{n - 2}$$

## Example (Textbook Example 12.8)

| $x$ | 12 | 30 | 36 | 40 | 45 | 57 | 62 | 67 | 71 | 78 | 93 | 94 | 100 | 105 |
|-----|----|----|----|----|----|----|----|----|----|----|----|----|-----|-----|
| $y$ | 3.3 | 3.2 | 3.4 | 3.0 | 2.8 | 2.9 | 2.7 | 2.6 | 2.5 | 2.6 | 2.2 | 2.0 | 2.3 | 2.1 |

Some statistics:

$$n = 14 \qquad \sum x_i = 890 \qquad \sum x_i^2 = 67182$$
$$\sum y_i = 37.6 \qquad \sum y_i^2 = 103.54 \qquad \sum x_i y_i = 2234.30$$

## Example (Textbook Example 12.8)

We can compute the following statistics:

$$S_{xx} = 10603.43, \quad S_{xy} = -155.99, \quad S_{yy} = 2.557$$

The estimators are

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{-155.99}{10603.43} = -0.0147$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = \frac{37.6}{14} - (-0.0147) \times \frac{890}{14} = 3.62$$

$$\hat{\sigma}^2 = \frac{S_{yy} - \hat{\beta}_1 S_{xy}}{n - 2} = \frac{2.557 - (-0.0147) \times (-155.99)}{14 - 2} = 0.022$$

# Properties of OLS Estimators

- Because $x_i$'s are fixed, $S_{xx}$ is not a random variable.

# Properties of OLS Estimators

- Because $x_i$'s are fixed, $S_{xx}$ is not a random variable.
- $S_{xy}$ can be written as

$$S_{xy} = \sum x_i y_i - n^{-1} \sum x_i \sum y_i = \sum_i \left[ (x_i - \bar{x}) \, y_i \right]$$

The highlighted $y_i$'s are the only random variables and we have

$$y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2),$$

where $\beta_0$ and $\beta_1$ are the true parameters.

# Properties of OLS Estimators

- Because $x_i$'s are fixed, $S_{xx}$ is not a random variable.
- $S_{xy}$ can be written as

$$S_{xy} = \sum x_i y_i - n^{-1} \sum x_i \sum y_i = \sum_i \left[ (x_i - \bar{x}) \, y_i \right]$$

  The highlighted $y_i$'s are the only random variables and we have

$$y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2),$$

  where $\beta_0$ and $\beta_1$ are the true parameters.
- Therefore, $S_{xy}$ is a linear combination of normal random variables and is also normally distributed,

$$S_{xy} \sim N(\beta_1 S_{xx}, \sigma^2 S_{xx})$$

# Properties of OLS Estimators

- Because $x_i$'s are fixed, $S_{xx}$ is not a random variable.
- $S_{xy}$ can be written as

$$S_{xy} = \sum x_i y_i - n^{-1} \sum x_i \sum y_i = \sum_i \left[ (x_i - \bar{x})\, y_i \right]$$

  The highlighted $y_i$'s are the only random variables and we have

$$y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2),$$

  where $\beta_0$ and $\beta_1$ are the true parameters.
- Therefore, $S_{xy}$ is a linear combination of normal random variables and is also normally distributed,

$$S_{xy} \sim N(\beta_1 S_{xx}, \sigma^2 S_{xx})$$

- Now we have

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \sim N(\beta_1, \sigma^2 S_{xx}^{-1})$$

## Properties of OLS Estimators

▶ For the intercept estimator, we have

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \sim N(\beta_0, (n^{-1} + \bar{x}^2 S_{xx}^{-1})\sigma^2)$$

## Properties of OLS Estimators

- For the intercept estimator, we have

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \sim N(\beta_0, (n^{-1} + \bar{x}^2 S_{xx}^{-1})\sigma^2)$$

- For the variance estimator, we have

$$E(\hat{\sigma}^2) = \sigma^2.$$

## Properties of OLS Estimators

Summary:

▶ All OLS estimators are **unbiased**:

$$E(\hat{\beta}_0) = \beta_0$$
$$E(\hat{\beta}_1) = \beta_1$$
$$E(\hat{\sigma}^2) = \sigma^2$$

▶ The estimated **standard errors (se)** of the estimators are:

$$\widehat{se}(\hat{\beta}_0) = \sqrt{(n^{-1} + \bar{x}^2 S_{xx}^{-1})\hat{\sigma}^2}$$
$$\widehat{se}(\hat{\beta}_1) = \sqrt{S_{xx}^{-1}\hat{\sigma}^2}$$
$$\widehat{se}(\hat{\sigma}^2) = \sqrt{\frac{2\hat{\sigma}^4}{n-2}}$$