

# ISE 315: Engineering Statistics

*Lecture 15: Simple Linear Regression (Part 1)*

Instructor: Mansur M. Arief, PhD  
Industrial and Systems Engineering, KFUPM

Office: 22-219 — Email: [mansur.arief@kfupm.edu.sa](mailto:mansur.arief@kfupm.edu.sa)

*Based on Montgomery & Runger, Applied Statistics and Probability for Engineers, 6th Ed.*

# Lecture 15

Simple Linear Regression and Correlation (Chapter 11)

# Congratulations!

*Major Exam 1 is Done*

---

Alhamdulillah — Major Exam 1 is behind you!

# Congratulations!

*Major Exam 1 is Done*

---

Alhamdulillah — Major Exam 1 is behind you!

You have worked through four chapters of statistical inference:

- Chapter 7: Sampling distributions & the Central Limit Theorem
- Chapter 8: Confidence intervals for  $\mu$ ,  $\sigma^2$ , and  $p$
- Chapter 9: Hypothesis testing (one sample)
- Chapter 10: Two-sample inference

# Congratulations!

*Major Exam 1 is Done*

---

Alhamdulillah — Major Exam 1 is behind you!

You have worked through four chapters of statistical inference:

- Chapter 7: Sampling distributions & the Central Limit Theorem
- Chapter 8: Confidence intervals for  $\mu$ ,  $\sigma^2$ , and  $p$
- Chapter 9: Hypothesis testing (one sample)
- Chapter 10: Two-sample inference

It's a **major accomplishment**. Take a moment to appreciate how far you've come.

## Grade Status Update

---

**Exam grades** are still being calculated and finalized.

## Grade Status Update

---

**Exam grades** are still being calculated and finalized.

We are also updating **all components** accumulated so far:

- Quizzes 1–4
- Homeworks 1–5
- Attendance (Lectures 1–14)
- Major Exam 1

## Grade Status Update

---

**Exam grades** are still being calculated and finalized.

We are also updating **all components** accumulated so far:

- Quizzes 1–4
- Homeworks 1–5
- Attendance (Lectures 1–14)
- Major Exam 1

**Target:** Your grades will be updated on Blackboard by **Tuesday**, in shaa Allah.

## Grade Status Update

---

**Exam grades** are still being calculated and finalized.

We are also updating **all components** accumulated so far:

- Quizzes 1–4
- Homeworks 1–5
- Attendance (Lectures 1–14)
- Major Exam 1

**Target:** Your grades will be updated on Blackboard by **Tuesday**, in shaa Allah.

If you notice any discrepancies, please email me *after* the grades are posted.

## Where Are We in the Course?

---

<b>Part</b>	<b>Topic</b>	<b>Status</b>
Part 1	Sampling distributions (Ch. 7)	✓ Done
Part 1	Confidence intervals (Ch. 8)	✓ Done
Part 2	Hypothesis testing — one sample (Ch. 9)	✓ Done
Part 2	Hypothesis testing — two samples (Ch. 10)	✓ Done
<b>Part 3</b>	<b>Simple linear regression (Ch. 11)</b>	<b>Starting today</b>
Part 3	Multiple linear regression (Ch. 12)	Upcoming
Part 4	Design of experiments (Ch. 13–14)	Upcoming

## Where Are We in the Course?

---

Part	Topic	Status
Part 1	Sampling distributions (Ch. 7)	✓ Done
Part 1	Confidence intervals (Ch. 8)	✓ Done
Part 2	Hypothesis testing — one sample (Ch. 9)	✓ Done
Part 2	Hypothesis testing — two samples (Ch. 10)	✓ Done
<b>Part 3</b>	<b>Simple linear regression (Ch. 11)</b>	<b>Starting today</b>
Part 3	Multiple linear regression (Ch. 12)	Upcoming
Part 4	Design of experiments (Ch. 13–14)	Upcoming

We are entering the **modeling** phase of the course. Instead of testing whether a parameter equals a value, we now ask: *how are two variables related?*

## Lecture 15 Outline

---

- Why regression? From hypothesis testing to modeling

## Lecture 15 Outline

---

- Why regression? From hypothesis testing to modeling
- Empirical models and scatter diagrams (Sec. 11-1)

## Lecture 15 Outline

---

- Why regression? From hypothesis testing to modeling
- Empirical models and scatter diagrams (Sec. 11-1)
- The simple linear regression model (Sec. 11-2)

## Lecture 15 Outline

---

- Why regression? From hypothesis testing to modeling
- Empirical models and scatter diagrams (Sec. 11-1)
- The simple linear regression model (Sec. 11-2)
- Method of least squares

## Lecture 15 Outline

---

- Why regression? From hypothesis testing to modeling
- Empirical models and scatter diagrams (Sec. 11-1)
- The simple linear regression model (Sec. 11-2)
- Method of least squares
- Properties of the least squares estimators (Sec. 11-3)

## Lecture 15 Outline

---

- Why regression? From hypothesis testing to modeling
- Empirical models and scatter diagrams (Sec. 11-1)
- The simple linear regression model (Sec. 11-2)
- Method of least squares
- Properties of the least squares estimators (Sec. 11-3)
- Estimating  $\sigma^2$  and a worked example

## Why Regression?

---

So far, our statistical questions have been:

- “Is the mean equal to some value?” (hypothesis testing)
- “What is a plausible range for the parameter?” (confidence intervals)

## Why Regression?

---

So far, our statistical questions have been:

- “Is the mean equal to some value?” (hypothesis testing)
- “What is a plausible range for the parameter?” (confidence intervals)

Now we ask a **different kind of question**:

*“How does one variable **relate to** another?”*

## Why Regression?

---

So far, our statistical questions have been:

- “Is the mean equal to some value?” (hypothesis testing)
- “What is a plausible range for the parameter?” (confidence intervals)

Now we ask a **different kind of question**:

*“How does one variable **relate to** another?”*

**Examples:**

- Does increasing training data improve an AI model's accuracy?
- How does inspection time affect defect detection rate?
- Can we predict equipment failure time from operating temperature?

## Why Regression?

---

So far, our statistical questions have been:

- “Is the mean equal to some value?” (hypothesis testing)
- “What is a plausible range for the parameter?” (confidence intervals)

Now we ask a **different kind of question**:

*“How does one variable **relate to** another?”*

**Examples:**

- Does increasing training data improve an AI model's accuracy?
- How does inspection time affect defect detection rate?
- Can we predict equipment failure time from operating temperature?

**Regression analysis** is the tool for building these *empirical models*.

# Regression in ISE and AI

*Real Applications You Will Encounter*

---

**During your ISE internship**, you might need to:

- Model how production speed affects product quality (manufacturing)
- Predict demand from historical sales and price data (supply chain)
- Estimate pipeline corrosion from measurements (oil & gas)

# Regression in ISE and AI

*Real Applications You Will Encounter*

---

**During your ISE internship**, you might need to:

- Model how production speed affects product quality (manufacturing)
- Predict demand from historical sales and price data (supply chain)
- Estimate pipeline corrosion from measurements (oil & gas)

**In modern AI and safety engineering:**

- Scaling laws: predicting model performance from training compute
- Verification: does increasing test coverage reduce failure rate?
- Calibration: is a sensor's reported value linearly related to the true value?

## Empirical Models

---

An **empirical model** is built from *observed data*, not from first principles.

## Empirical Models

---

An **empirical model** is built from *observed data*, not from first principles.

Suppose we observe pairs  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ :

- $x =$  **regressor** (predictor, independent variable)
- $y =$  **response** (dependent variable)

## Empirical Models

---

An **empirical model** is built from *observed data*, not from first principles.

Suppose we observe pairs  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ :

- $x =$  **regressor** (predictor, independent variable)
- $y =$  **response** (dependent variable)

A **scatter diagram** plots each observation as a point  $(x_i, y_i)$  and helps us visualize the relationship.

## Empirical Models

---

An **empirical model** is built from *observed data*, not from first principles.

Suppose we observe pairs  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ :

- $x =$  **regressor** (predictor, independent variable)
- $y =$  **response** (dependent variable)

A **scatter diagram** plots each observation as a point  $(x_i, y_i)$  and helps us visualize the relationship.

**Key question:** Does the scatter suggest a straight-line (linear) relationship? If so, we can fit a **simple linear regression** model.

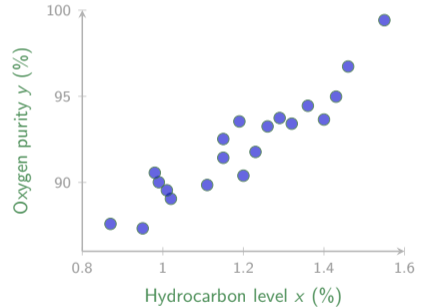
## Scatter Diagram: What Pattern Do You See?

---

**Example:** An ISE intern at a petrochemical plant measures hydrocarbon level ( $x$ , %) and oxygen purity ( $y$ , %) for 20 batches.

## Scatter Diagram: What Pattern Do You See?

**Example:** An ISE intern at a petrochemical plant measures hydrocarbon level ( $x$ , %) and oxygen purity ( $y$ , %) for 20 batches.

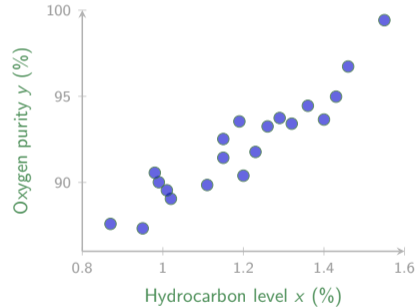


## Scatter Diagram: What Pattern Do You See?

**Example:** An ISE intern at a petrochemical plant measures hydrocarbon level ( $x$ , %) and oxygen purity ( $y$ , %) for 20 batches.

The scatter diagram suggests:

- A **positive linear** trend
- Higher hydrocarbon  $\Rightarrow$  higher purity
- Some scatter around the trend



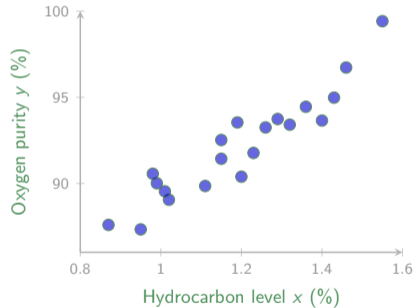
## Scatter Diagram: What Pattern Do You See?

**Example:** An ISE intern at a petrochemical plant measures hydrocarbon level ( $x$ , %) and oxygen purity ( $y$ , %) for 20 batches.

The scatter diagram suggests:

- A **positive linear** trend
- Higher hydrocarbon  $\Rightarrow$  higher purity
- Some scatter around the trend

Can we **quantify** this relationship?



## The Simple Linear Regression Model

---

We model the **expected value** of  $Y$  at each level of  $x$  as a straight line:

$$E(Y | x) = \beta_0 + \beta_1 x$$

## The Simple Linear Regression Model

---

We model the **expected value** of  $Y$  at each level of  $x$  as a straight line:

$$E(Y | x) = \beta_0 + \beta_1 x$$

Each individual observation is:

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, 2, \dots, n$$

## The Simple Linear Regression Model

---

We model the **expected value** of  $Y$  at each level of  $x$  as a straight line:

$$E(Y | x) = \beta_0 + \beta_1 x$$

Each individual observation is:

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, 2, \dots, n$$

where:

- $\beta_0 =$  **intercept** (value of  $E(Y)$  when  $x = 0$ )
- $\beta_1 =$  **slope** (change in  $E(Y)$  per one-unit increase in  $x$ )
- $\varepsilon_i =$  **random error** for observation  $i$

## The Simple Linear Regression Model

---

We model the **expected value** of  $Y$  at each level of  $x$  as a straight line:

$$E(Y | x) = \beta_0 + \beta_1 x$$

Each individual observation is:

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, 2, \dots, n$$

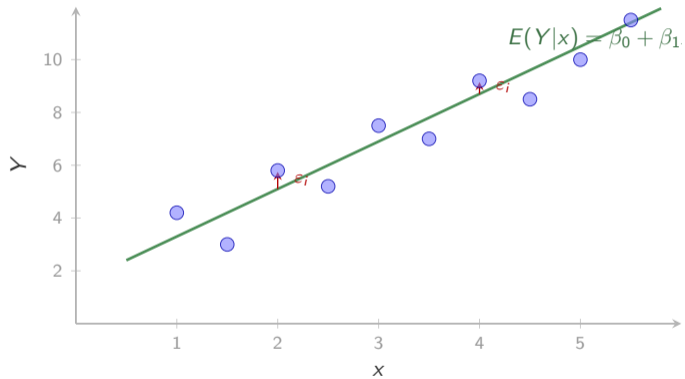
where:

- $\beta_0 =$  **intercept** (value of  $E(Y)$  when  $x = 0$ )
- $\beta_1 =$  **slope** (change in  $E(Y)$  per one-unit increase in  $x$ )
- $\varepsilon_i =$  **random error** for observation  $i$

**Assumptions on errors:**  $E(\varepsilon_i) = 0$ ,  $\text{Var}(\varepsilon_i) = \sigma^2$ , and  $\varepsilon_i$  are uncorrelated.

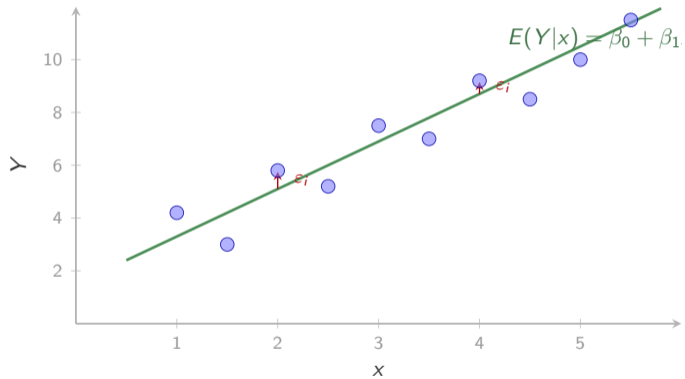
## Visualizing the Model

---



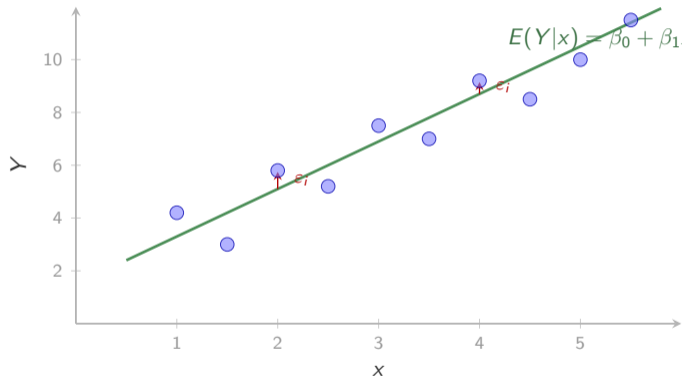
## Visualizing the Model

---

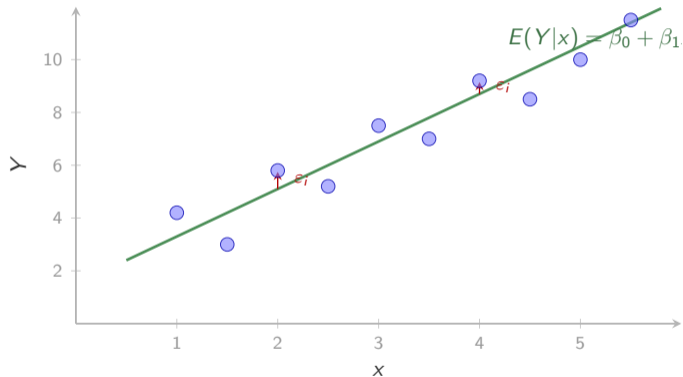


## Visualizing the Model

---



## Visualizing the Model



The **green line** is the true (unknown) regression function. The **blue points** are the observed data. The **red arrows** show the errors  $\epsilon_i$ .

## Our Goal: Estimate the Unknown Parameters

---

We do not know  $\beta_0$  and  $\beta_1$ . We need to **estimate** them from the data.

## Our Goal: Estimate the Unknown Parameters

---

We do not know  $\beta_0$  and  $\beta_1$ . We need to **estimate** them from the data.

Once we have estimates  $\hat{\beta}_0$  and  $\hat{\beta}_1$ , the **fitted regression line** is:

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

## Our Goal: Estimate the Unknown Parameters

---

We do not know  $\beta_0$  and  $\beta_1$ . We need to **estimate** them from the data.

Once we have estimates  $\hat{\beta}_0$  and  $\hat{\beta}_1$ , the **fitted regression line** is:

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

The **residual** for observation  $i$  is the difference between the actual value and the fitted value:

$$e_i = y_i - \hat{y}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$$

## Our Goal: Estimate the Unknown Parameters

---

We do not know  $\beta_0$  and  $\beta_1$ . We need to **estimate** them from the data.

Once we have estimates  $\hat{\beta}_0$  and  $\hat{\beta}_1$ , the **fitted regression line** is:

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

The **residual** for observation  $i$  is the difference between the actual value and the fitted value:

$$e_i = y_i - \hat{y}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$$

**Question:** How should we choose  $\hat{\beta}_0$  and  $\hat{\beta}_1$  so that the line fits the data “best”?

# Method of Least Squares

*Minimizing the Sum of Squared Residuals*

---

Choose  $\hat{\beta}_0$  and  $\hat{\beta}_1$  to **minimize** the sum of squared residuals:

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

## Method of Least Squares

*Minimizing the Sum of Squared Residuals*

---

Choose  $\hat{\beta}_0$  and  $\hat{\beta}_1$  to **minimize** the sum of squared residuals:

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

Take partial derivatives, set them equal to zero, and solve. This gives the **normal equations**:

$$\frac{\partial S}{\partial \beta_0} = 0 \quad \text{and} \quad \frac{\partial S}{\partial \beta_1} = 0$$

## Method of Least Squares

*Minimizing the Sum of Squared Residuals*

---

Choose  $\hat{\beta}_0$  and  $\hat{\beta}_1$  to **minimize** the sum of squared residuals:

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

Take partial derivatives, set them equal to zero, and solve. This gives the **normal equations**:

$$\frac{\partial S}{\partial \beta_0} = 0 \quad \text{and} \quad \frac{\partial S}{\partial \beta_1} = 0$$

**Why “least squares”?** We minimize the total squared distance from points to the line — the same core idea used in many machine learning methods!

# Least Squares Estimators

*The Formulas*

---

The **least squares estimators** are:

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$$

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

# Least Squares Estimators

## *The Formulas*

---

The **least squares estimators** are:

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

where:

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n}$$
$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{(\sum x_i)^2}{n}$$

# Least Squares Estimators

## The Formulas

---

The **least squares estimators** are:

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$$

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

where:

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n}$$
$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{(\sum x_i)^2}{n}$$

$\hat{\beta}_1$  measures the *co-movement* of  $x$  and  $y$  relative to the *spread* of  $x$ .

## Example: AI Safety Verification

*Does More Testing Reduce Failures?*

---

An ISE engineer is evaluating an autonomous inspection system at a refinery. She runs the system through different amounts of verification testing (measured in hours) and records the failure rate (failures per 1000 inspections) in deployment:

## Example: AI Safety Verification

*Does More Testing Reduce Failures?*

---

An ISE engineer is evaluating an autonomous inspection system at a refinery. She runs the system through different amounts of verification testing (measured in hours) and records the failure rate (failures per 1000 inspections) in deployment:

$x$ (test hours)	10	15	20	25	30	35	40	45
$y$ (failures/1000)	48	42	38	32	29	24	22	18

## Example: AI Safety Verification

*Does More Testing Reduce Failures?*

---

An ISE engineer is evaluating an autonomous inspection system at a refinery. She runs the system through different amounts of verification testing (measured in hours) and records the failure rate (failures per 1000 inspections) in deployment:

$x$ (test hours)	10	15	20	25	30	35	40	45
$y$ (failures/1000)	48	42	38	32	29	24	22	18

Useful summaries ( $n = 8$ ):

$$\bar{x} = 27.5, \quad \bar{y} = 31.625, \quad S_{xx} = 1050, \quad S_{xy} = -881.25$$

## Example: AI Safety Verification (Solution)

---

**Step 1:** Compute the slope.

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{-881.25}{1050} = -0.8393$$

## Example: AI Safety Verification (Solution)

---

**Step 1:** Compute the slope.

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{-881.25}{1050} = -0.8393$$

**Step 2:** Compute the intercept.

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 31.625 - (-0.8393)(27.5) = 54.706$$

## Example: AI Safety Verification (Solution)

---

**Step 1:** Compute the slope.

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{-881.25}{1050} = -0.8393$$

**Step 2:** Compute the intercept.

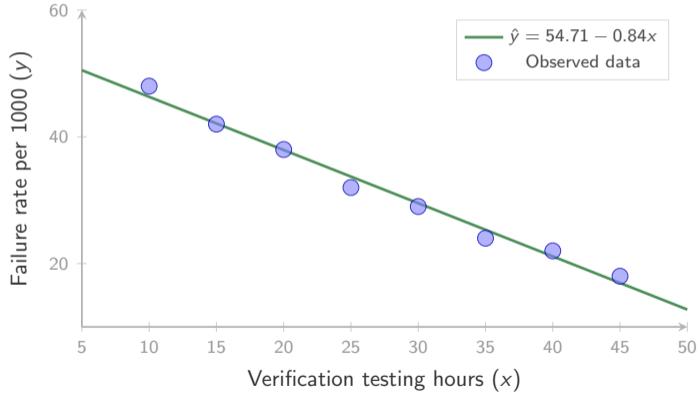
$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 31.625 - (-0.8393)(27.5) = 54.706$$

**Step 3:** Write the fitted model.

$$\hat{y} = 54.706 - 0.8393x$$

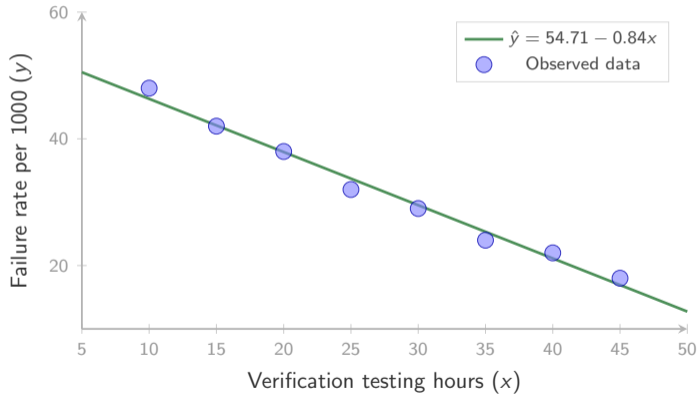
## Example: Fitted Regression Line

---



## Example: Fitted Regression Line

---



## Think About It

---

Based on the fitted model  $\hat{y} = 54.706 - 0.8393 x$ :

## Think About It

---

Based on the fitted model  $\hat{y} = 54.706 - 0.8393 x$ :

**Question 1:** What failure rate would you predict for  $x = 50$  hours of testing?

## Think About It

---

Based on the fitted model  $\hat{y} = 54.706 - 0.8393x$ :

**Question 1:** What failure rate would you predict for  $x = 50$  hours of testing?

$$\hat{y} = 54.706 - 0.8393(50) = 54.706 - 41.964 = 12.74 \text{ failures per 1000.}$$

## Think About It

---

Based on the fitted model  $\hat{y} = 54.706 - 0.8393x$ :

**Question 1:** What failure rate would you predict for  $x = 50$  hours of testing?

$$\hat{y} = 54.706 - 0.8393(50) = 54.706 - 41.964 = 12.74 \text{ failures per 1000.}$$

**Question 2:** Would you trust this prediction for  $x = 200$  hours?

## Think About It

---

Based on the fitted model  $\hat{y} = 54.706 - 0.8393x$ :

**Question 1:** What failure rate would you predict for  $x = 50$  hours of testing?

$$\hat{y} = 54.706 - 0.8393(50) = 54.706 - 41.964 = 12.74 \text{ failures per 1000.}$$

**Question 2:** Would you trust this prediction for  $x = 200$  hours?

$$\hat{y} = 54.706 - 0.8393(200) = -113.15 \quad (\text{negative — nonsensical!})$$

**Warning:** Do not **extrapolate** far beyond the range of the observed data. The linear model may not hold outside  $10 \leq x \leq 45$ .

# Properties of the Least Squares Estimators

## *Section 11-3*

---

The least squares estimators  $\hat{\beta}_0$  and  $\hat{\beta}_1$  have nice statistical properties:

# Properties of the Least Squares Estimators

## Section 11-3

---

The least squares estimators  $\hat{\beta}_0$  and  $\hat{\beta}_1$  have nice statistical properties:

**Slope:**

$$E(\hat{\beta}_1) = \beta_1 \quad \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}}$$

$\hat{\beta}_1$  is an **unbiased** estimator of  $\beta_1$ .

# Properties of the Least Squares Estimators

## Section 11-3

---

The least squares estimators  $\hat{\beta}_0$  and  $\hat{\beta}_1$  have nice statistical properties:

**Slope:**

$$E(\hat{\beta}_1) = \beta_1 \quad \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}}$$

$\hat{\beta}_1$  is an **unbiased** estimator of  $\beta_1$ .

**Intercept:**

$$E(\hat{\beta}_0) = \beta_0 \quad \text{Var}(\hat{\beta}_0) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)$$

$\hat{\beta}_0$  is an **unbiased** estimator of  $\beta_0$ .

# Properties of the Least Squares Estimators

## Section 11-3

---

The least squares estimators  $\hat{\beta}_0$  and  $\hat{\beta}_1$  have nice statistical properties:

**Slope:**

$$E(\hat{\beta}_1) = \beta_1 \quad \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}}$$

$\hat{\beta}_1$  is an **unbiased** estimator of  $\beta_1$ .

**Intercept:**

$$E(\hat{\beta}_0) = \beta_0 \quad \text{Var}(\hat{\beta}_0) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)$$

$\hat{\beta}_0$  is an **unbiased** estimator of  $\beta_0$ .

These are **exactly the same ideas** from Chapters 7–8: unbiasedness and variance of estimators. Regression builds on what you already know!

## Estimating $\sigma^2$

### *The Error Variance*

---

We also need to estimate the error variance  $\sigma^2$ . The **error sum of squares** is:

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

## Estimating $\sigma^2$

### *The Error Variance*

---

We also need to estimate the error variance  $\sigma^2$ . The **error sum of squares** is:

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

A computational shortcut:

$$SS_E = SS_T - \hat{\beta}_1 S_{xy} \quad \text{where} \quad SS_T = S_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2$$

## Estimating $\sigma^2$

### *The Error Variance*

---

We also need to estimate the error variance  $\sigma^2$ . The **error sum of squares** is:

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

A computational shortcut:

$$SS_E = SS_T - \hat{\beta}_1 S_{xy} \quad \text{where} \quad SS_T = S_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2$$

The **unbiased estimator** of  $\sigma^2$  is:

$$\hat{\sigma}^2 = \frac{SS_E}{n-2} = MS_E$$

## Estimating $\sigma^2$

### *The Error Variance*

---

We also need to estimate the error variance  $\sigma^2$ . The **error sum of squares** is:

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

A computational shortcut:

$$SS_E = SS_T - \hat{\beta}_1 S_{xy} \quad \text{where} \quad SS_T = S_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2$$

The **unbiased estimator** of  $\sigma^2$  is:

$$\hat{\sigma}^2 = \frac{SS_E}{n-2} = MS_E$$

## Example: Computing $\hat{\sigma}^2$

*AI Safety Verification (Continued)*

---

From our AI verification example:

$$S_{yy} = \sum y_i^2 - n\bar{y}^2 = 8798 - 8(31.625)^2 = 800.875$$

## Example: Computing $\hat{\sigma}^2$

*AI Safety Verification (Continued)*

---

From our AI verification example:

$$S_{yy} = \sum y_i^2 - n\bar{y}^2 = 8798 - 8(31.625)^2 = 800.875$$

$$SS_E = S_{yy} - \hat{\beta}_1 \cdot S_{xy} = 800.875 - (-0.8393)(-881.25) = 61.20$$

## Example: Computing $\hat{\sigma}^2$

*AI Safety Verification (Continued)*

---

From our AI verification example:

$$S_{yy} = \sum y_i^2 - n\bar{y}^2 = 8798 - 8(31.625)^2 = 800.875$$

$$SS_E = S_{yy} - \hat{\beta}_1 \cdot S_{xy} = 800.875 - (-0.8393)(-881.25) = 61.20$$

$$\hat{\sigma}^2 = MS_E = \frac{SS_E}{n-2} = \frac{61.20}{6} = 10.20 \quad \Rightarrow \quad \hat{\sigma} = \sqrt{10.20} = 3.19 \text{ failures per 1000}$$

## Connecting Old and New Ideas

---

Concept	Ch. 7–10	Ch. 11 (Regression)
Parameter	$\mu, \sigma^2, \rho$	$\beta_0, \beta_1, \sigma^2$
Estimator	$\bar{X}, S^2, \hat{P}$	$\hat{\beta}_0, \hat{\beta}_1, MS_E$
Unbiased?	Yes	Yes
Degrees of freedom	$n - 1$	$n - 2$
Next step	CI and hypothesis test	<i>Same — next lecture!</i>

## Connecting Old and New Ideas

---

Concept	Ch. 7–10	Ch. 11 (Regression)
Parameter	$\mu, \sigma^2, \rho$	$\beta_0, \beta_1, \sigma^2$
Estimator	$\bar{X}, S^2, \hat{P}$	$\hat{\beta}_0, \hat{\beta}_1, MS_E$
Unbiased?	Yes	Yes
Degrees of freedom	$n - 1$	$n - 2$
Next step	CI and hypothesis test	<i>Same — next lecture!</i>

The structure is **identical**: estimate parameters, quantify uncertainty, then test hypotheses and build confidence intervals.

## Practice Problem

*Exam-Style: Fitting a Simple Linear Regression*

---

A quality engineer wants to study how **production speed** ( $x$ , in units/hr) affects the **defect rate** ( $y$ , in defects per 100 units). Data was collected over  $n = 5$  production runs:

## Practice Problem

*Exam-Style: Fitting a Simple Linear Regression*

---

A quality engineer wants to study how **production speed** ( $x$ , in units/hr) affects the **defect rate** ( $y$ , in defects per 100 units). Data was collected over  $n = 5$  production runs:

<b>Run</b>	1	2	3	4	5
$x$ (units/hr)	2	4	6	8	10
$y$ (defects/100)	2	5	4	7	7

## Practice Problem

*Exam-Style: Fitting a Simple Linear Regression*

---

A quality engineer wants to study how **production speed** ( $x$ , in units/hr) affects the **defect rate** ( $y$ , in defects per 100 units). Data was collected over  $n = 5$  production runs:

Run	1	2	3	4	5
$x$ (units/hr)	2	4	6	8	10
$y$ (defects/100)	2	5	4	7	7

The following summary statistics may be useful:

$\sum x$	$\sum y$	$\sum x^2$	$\sum y^2$	$\sum xy$
30	25	220	143	174

## Practice Problem

*Exam-Style: Fitting a Simple Linear Regression*

A quality engineer wants to study how **production speed** ( $x$ , in units/hr) affects the **defect rate** ( $y$ , in defects per 100 units). Data was collected over  $n = 5$  production runs:

Run	1	2	3	4	5
$x$ (units/hr)	2	4	6	8	10
$y$ (defects/100)	2	5	4	7	7

The following summary statistics may be useful:

$\sum x$	$\sum y$	$\sum x^2$	$\sum y^2$	$\sum xy$
30	25	220	143	174

## Practice Problem — Solution

*Part (a): Fit the Model*

---

**Compute**  $S_{xx}$  and  $S_{xy}$ :

$$S_{xx} = \sum x_i^2 - \frac{(\sum x_i)^2}{n} = 220 - \frac{(30)^2}{5} = 220 - 180 = 40$$

$$S_{xy} = \sum x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n} = 174 - \frac{(30)(25)}{5} = 174 - 150 = 24$$

## Practice Problem — Solution

Part (a): Fit the Model

---

**Compute**  $S_{xx}$  and  $S_{xy}$ :

$$S_{xx} = \sum x_i^2 - \frac{(\sum x_i)^2}{n} = 220 - \frac{(30)^2}{5} = 220 - 180 = 40$$

$$S_{xy} = \sum x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n} = 174 - \frac{(30)(25)}{5} = 174 - 150 = 24$$

**Slope:**

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{24}{40} = 0.60$$

## Practice Problem — Solution

Part (a): Fit the Model

---

**Compute**  $S_{xx}$  and  $S_{xy}$ :

$$S_{xx} = \sum x_i^2 - \frac{(\sum x_i)^2}{n} = 220 - \frac{(30)^2}{5} = 220 - 180 = 40$$

$$S_{xy} = \sum x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n} = 174 - \frac{(30)(25)}{5} = 174 - 150 = 24$$

**Slope:**

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{24}{40} = 0.60$$

**Intercept:**  $\bar{x} = 30/5 = 6$ ,  $\bar{y} = 25/5 = 5$ .

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 5 - 0.60(6) = 5 - 3.6 = 1.40$$

## Practice Problem — Solution

Part (a): Fit the Model

---

**Compute**  $S_{xx}$  and  $S_{xy}$ :

$$S_{xx} = \sum x_i^2 - \frac{(\sum x_i)^2}{n} = 220 - \frac{(30)^2}{5} = 220 - 180 = 40$$

$$S_{xy} = \sum x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n} = 174 - \frac{(30)(25)}{5} = 174 - 150 = 24$$

**Slope:**

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{24}{40} = 0.60$$

**Intercept:**  $\bar{x} = 30/5 = 6$ ,  $\bar{y} = 25/5 = 5$ .

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 5 - 0.60(6) = 5 - 3.6 = 1.40$$

## Practice Problem — Solution (Continued)

*Parts (b), (c), and (d)*

---

**(b)** Estimate defect rate at  $x = 7$ :

$$\hat{y} = 1.40 + 0.60(7) = 1.40 + 4.20 = \boxed{5.60 \text{ defects per } 100}$$

## Practice Problem — Solution (Continued)

Parts (b), (c), and (d)

---

**(b)** Estimate defect rate at  $x = 7$ :

$$\hat{y} = 1.40 + 0.60(7) = 1.40 + 4.20 = \boxed{5.60 \text{ defects per } 100}$$

**(c)** Residual for observation at  $x = 4$  (where  $y = 5$ ):

$$\hat{y} = 1.40 + 0.60(4) = 3.80 \quad e = y - \hat{y} = 5 - 3.80 = \boxed{+1.20}$$

Since  $e > 0$  (actual  $>$  predicted), the model **underestimates** the true value.

## Practice Problem — Solution (Continued)

Parts (b), (c), and (d)

---

**(b)** Estimate defect rate at  $x = 7$ :

$$\hat{y} = 1.40 + 0.60(7) = 1.40 + 4.20 = \boxed{5.60 \text{ defects per } 100}$$

**(c)** Residual for observation at  $x = 4$  (where  $y = 5$ ):

$$\hat{y} = 1.40 + 0.60(4) = 3.80 \quad e = y - \hat{y} = 5 - 3.80 = \boxed{+1.20}$$

Since  $e > 0$  (actual  $>$  predicted), the model **underestimates** the true value.

**(d)** Compute  $\hat{\sigma}^2$ : First,  $S_{yy} = \sum y_i^2 - n\bar{y}^2 = 143 - 5(5)^2 = 143 - 125 = 18$ .

$$SS_E = S_{yy} - \hat{\beta}_1 \cdot S_{xy} = 18 - 0.60(24) = 18 - 14.4 = 3.60$$

$$\hat{\sigma}^2 = \frac{SS_E}{n-2} = \frac{3.60}{3} = \boxed{1.20}$$

## Lecture 15 Summary

---

- **Simple linear regression** models the relationship  $Y = \beta_0 + \beta_1 x + \varepsilon$ .

## Lecture 15 Summary

---

- **Simple linear regression** models the relationship  $Y = \beta_0 + \beta_1 x + \varepsilon$ .
- The **method of least squares** minimizes  $\sum e_i^2$  to estimate  $\beta_0$  and  $\beta_1$ .

## Lecture 15 Summary

---

- **Simple linear regression** models the relationship  $Y = \beta_0 + \beta_1 x + \varepsilon$ .
- The **method of least squares** minimizes  $\sum e_i^2$  to estimate  $\beta_0$  and  $\beta_1$ .
- Key formulas:  $\hat{\beta}_1 = S_{xy}/S_{xx}$  and  $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ .

## Lecture 15 Summary

---

- **Simple linear regression** models the relationship  $Y = \beta_0 + \beta_1 x + \varepsilon$ .
- The **method of least squares** minimizes  $\sum e_i^2$  to estimate  $\beta_0$  and  $\beta_1$ .
- Key formulas:  $\hat{\beta}_1 = S_{xy}/S_{xx}$  and  $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ .
- The error variance is estimated by  $\hat{\sigma}^2 = SS_E/(n - 2)$ .

## Lecture 15 Summary

---

- **Simple linear regression** models the relationship  $Y = \beta_0 + \beta_1 x + \varepsilon$ .
- The **method of least squares** minimizes  $\sum e_i^2$  to estimate  $\beta_0$  and  $\beta_1$ .
- Key formulas:  $\hat{\beta}_1 = S_{xy}/S_{xx}$  and  $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ .
- The error variance is estimated by  $\hat{\sigma}^2 = SS_E/(n - 2)$ .
- **Do not extrapolate** the model far beyond the range of observed  $x$  values.

## Lecture 15 Summary

---

- **Simple linear regression** models the relationship  $Y = \beta_0 + \beta_1 x + \varepsilon$ .
- The **method of least squares** minimizes  $\sum e_i^2$  to estimate  $\beta_0$  and  $\beta_1$ .
- Key formulas:  $\hat{\beta}_1 = S_{xy}/S_{xx}$  and  $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ .
- The error variance is estimated by  $\hat{\sigma}^2 = SS_E/(n - 2)$ .
- **Do not extrapolate** the model far beyond the range of observed  $x$  values.
- **Next lecture:** Hypothesis tests on  $\beta_0$  and  $\beta_1$ , ANOVA for regression, confidence intervals, prediction intervals, and  $R^2$ .

## Reminders

---

- **Major Exam 1 grades** and grades will be on Blackboard by **Tuesday**

## Reminders

---

- **Major Exam 1 grades** and grades will be on Blackboard by **Tuesday**
- If you have grade questions, please come next class or the office hour

## Reminders

---

- **Major Exam 1 grades** and grades will be on Blackboard by **Tuesday**
- If you have grade questions, please come next class or the office hour
- Office hours: Tuesdays 9–10 AM, Room 22-219 or Zoom.

## Reminders

---

- **Major Exam 1 grades** and grades will be on Blackboard by **Tuesday**
- If you have grade questions, please come next class or the office hour
- Office hours: Tuesdays 9–10 AM, Room 22-219 or Zoom.
- Start reading **Chapter 11** (Sections 11-1 through 11-3 for now).