

## Linear Regression and Correlation

### Introduction

Many real-world situations involve understanding relationships between two or more numerical variables including:

- Education & Income: Does a higher education (e.g.) lead to a higher salary?

Exam Performance: Does a student's grade on a midterm predict their final exam score?

- Repair Costs: How does the total cost of a repair depend on hourly labor fees?

The data in these examples are called bivariate data because they involve two variables. More complex analysis use multivariate data (three or more variables). This focuses on the simplest type of regression: linear regression with one independent variable ( $X$ ). We also study correlation, which measures how strong the relationship is between two variables.

### Key Concepts

Measures the strength and direction of the relationship between  $X$  and  $y$ .

The correlation coefficient ( $r$ ) ranges from  $(-1)$  to  $(1)$ .

$r = +1 \rightarrow$  Perfect positive relationship (as  $X$  increases,  $y$  increases)

$r = -1 \rightarrow$  Perfect negative relationship (as  $X$  increases,  $y$  decreases)

$r = 0 \rightarrow$  No correlation

### 3. Outlier

Outliers are data points that do not follow the pattern of the rest of the data.

They can affect regression results and correlation calculation.

Example: A student who scores low on midterms but high on the final exam may be an outlier.

### Linear Equations in Regression

Linear regression is based on the equation of a straight line,

written as:  $y = a + bx$

where

$x$  = independent variable,  $y$  = dependent variable,  $a$  =

$y$ -Intercept (value of  $y$  when  $x=0$ )  $b$  = slope (rate of change)

## Examples of Linear Equations

1.  $y = 3 + 2x$

When  $x = 0, y = 3$

for every increase of 1 in  $x$ ,  $y$  increases by 2

2.  $y = -0.01 + 1.2x$

When  $x = 0, y = -0.01$

for every increase of 1 in  $x$ ,  $y$  increases by 1.2

Key fact: The graph of any linear equation is a straight line (unless it is vertical)

Real-life Example: Aaron's word processing service (AWPS)

AWPS charges \$32 per hour

One-time setup fee of \$31.50

Step 1: Define variable

Let  $x$  = number of hours needed for the job

Let  $y$  = total cost

Step 2: Write the equation

fixed cost = \$31.50

Hourly rate = \$32 per hour

$$y = 31.50 + 32x$$

This equation models total cost based on the number of hours worked.

Example calculation:

If the job takes 5 hours, the total cost is

$$y = 31.50 + (32 \times 5) = 31.50 + 160 = 191.50$$

so the total charge for 5 hours is \$191.50

Why are lines & Equation useful?

Predicts a value based on given inputs

Analyze trends in business, economics, psychology & sciences

Create models for decision making

Interpreting the slope ( $b$ )

The slope  $b$  determines the direction and steepness of the line

if  $b > 0$ : The line slopes upward (positive relationship)

if  $b = 0$  The line is horizontal (no change)

if  $b < 0$ : The line slopes downward (negative relationship)

Understanding the correlation coefficient ( $r$ ) and coefficient of determination ( $r^2$ )

The correlation coefficient ( $r$ ) measures the strength and direction of a linear relationship between two variables

Interpreting the correlation coefficient ( $r$ )

Range:  $-1 \leq r \leq 1$ ,  $r = 1 \rightarrow$  perfect positive correlation

$r = -1 \rightarrow$  perfect negative correlation

$r = 0 \rightarrow$  No linear correlation

Strength of relationship:  $|r| \approx 1 \rightarrow$  strong correlation

$|r| \approx 0.5 \rightarrow$  Moderate correlation

$|r| \approx 0 \rightarrow$  weak or no correlation

Coefficient of Determination ( $r^2$ )

What does  $r^2$  tell us?

$r^2$  is the square of  $r$  and is expressed as a percentage

Represents the percentage of variation in  $y$  explained by  $x$

Example

given  $\hat{y} = 0.663x + 1.165$

$$R^2 = (0.663)^2 = 0.4397 \approx 44\%$$

Interpretation

- 44% of the variation in final exam scores is explained by this exam scores
- The remaining 56% is due to other factors (e.g. study/habit, sleep, stress)

Summary of Partitioning: the sum of squares

### 1 Understanding the sum of squares (SSY)

The total variation in  $y$  is measured using sum of squares  $Y$  (SSY), which represent the sum of squared deviation of each  $y$  value from the mean of  $y$

$$SSY = \sum (y - \bar{y})^2$$

Eg Dataset

$$\text{Mean of } y = 2.06$$

$$SSY = 4.397 \text{ (sum of squared deviation from the mean)}$$

### 2 Converting raw scores to deviation scores

Instead of using raw scores  $y$ , it is useful to compute deviation scores, represented as  $y = y - \bar{y}$

Conventionally, the small letters ( $y$ ) denote deviations from the mean

### 3 Computing predicted scores from the Regression Equation

The predicted value  $\hat{y}_1$  is obtained using the regression equation  $\hat{y}_1 = 0.425x + 0.785$

$$\text{Eg. for } x = 1 \quad \hat{y}_1 = 0.425(1) + 0.785 = 1.210$$

4 Partitioning  $SSY$  into  $SSY'$  and  $SS$   
 $SSY$  (Total variation) can be divided into  
 $SSY'$  (sum of squares Predicted): Variation of Predicted  $Y$  values  
 from the mean:

$SSE$  (sum of squares Errors): Validation of actual  $Y$  values from  
 $Y'$  (unexplained variation)

Formula:  $SSY = SSY' + SSE$

Eg calculation  $4.597 = 1.806 + 2.79$

5 Interpretation of  $r^2$

$r^2$  represents the proportion of explained variation

$$r^2 = \frac{SSY'}{SSY}$$

Eg: If  $r=0.4$ , then  $r^2=0.16$ , meaning 16% of variation is explained by the predictor variable

The proportion of unexplained variation =  $\frac{SSE}{SSY}$

Summary of standard errors of the estimate

- 1) Understanding the standard error of the estimate (or est)  
 The standard error of the estimate (or est) measures the accuracy of regression predictions  
 A smaller standard error means the predictions are more accurate (points closer to the regression line), while a larger error indicates more spread (less accuracy)

2) Formula for standard error of the estimate

The standard error is calculated using

$$est = \sqrt{\frac{\sum (Y - Y')^2}{N}}$$

3 Example Calculation: (Population Case,  $N=5$ )

from the dataset

$x$	$y$	$y_i$	$y - y_i$	$(y - y_i)^2$
1.0	1.00	1.210	-0.210	0.044
2.0	2.00	1.635	0.365	0.133
3.0	1.30	2.060	-0.760	0.578
4.0	3.75	2.485	-1.265	1.600
5.0	2.25	2.910	-0.660	0.436
sum	15.00	10.30	0.00	2.797

4 Computing standard error using Pearson's correlation ( $\rho$ )

An alternative formula relates  $\sigma_{\text{est}}$  to Pearson's Correlation

$$\sigma_{\text{est}} = \sqrt{(1 - \rho^2) \frac{SSY}{N}}$$

where;  $\rho = 0.6768$  (Pearson Correlation)

$SSY = 4.597$  (total sum of squares)

$N = 5$

substituting values =  $\sigma_{\text{est}} = 0.747$

5 Standard error of the estimate for a sample

When estimating from a sample, the denominator changes from  $N$  to  $N-2$

Why use  $N-2$ ?

Two parameters (slope and intercept) are estimated before computing  $SSY$

1) Inferential statistics for regression slope ( $b$ ) and correlation ( $r$ )

2) Assumption  $\rightarrow$  the inferential statistics in regression

Inferential statistics for regression rely on several assumptions about the population (not just the sample):

i. Linearity: - The relationship between  $X$  (predictor) and  $Y$  (criterion) is linear

ii. Homoscedasticity: - The variance of residuals ( $\text{errors}$ ) is constant across all values of  $X$

- If the spread of residuals changes (eg wider spread at one end), the assumption is violated (heteroscedasticity)

- Example (Figure 3): Predictions are more accurate for students with higher GPAs, but less accurate for lower GPAs  $\rightarrow$  heteroscedasticity

iii. Normality of errors: - The residuals (deviations from the regression lines) follow a normal distribution

Note: This does not mean  $X$  or  $y$  is normally distributed - only the errors should be

2) Significance test for the slope ( $b$ )

It test whether the regression slope ( $b$ ) is significantly different from zero. (ie whether  $X$  significantly predicts  $Y$ )

Formula for  $t$ -test on the slope ( $b$ )  $t = \frac{b - 0}{s_b}$

Standard error of the slope ( $s_b$ )

$$s_b = \frac{\text{S.E. of } b}{\sqrt{S.S.X}}$$

## 3 Confidence Interval for the Slope (b)

A 95% confidence interval for the slope is  $b \pm (t_{95\%} s_b)$

## 4 Significance Test for Pearson's Correlation (r)

To test whether the correlation (r) is significantly different from zero, we use

$$t = \frac{r\sqrt{N-2}}{\sqrt{1-r^2}}$$