
 **Marmara University, 2021**

Probability and Statistics

Subject 3
Describing Bivariate Data


Mujdat Soyuturk, Ph.D.
Associate Professor

 **Contents**

- Bivariate Data
- Graphs for Categorical Variables
- Scatterplots for Two Quantitative Variables
- Numerical Measures for Quantitative Bivariate Data
 - Correlation Coefficient
 - Covariance
 - Regression


Most parts of the slides are derived from the textbook: "Mendenhall, Beaver, Beaver, Introduction to Probability and Statistics, 14th Ed., Brooks/Cole, Cengage Learning, 2013"

3 - 2 Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

 **Bivariate Data**

- When two variables are measured on a single experimental unit, the resulting data are called **bivariate data**.
- You can describe each variable individually, and you can also explore the **relationship** between the two variables.
- Bivariate data can be described with
 - Graphs
 - Numerical Measures

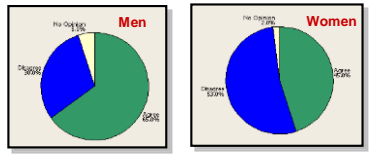
3 - 3 Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

 **Graphs for Qualitative Variables**


When at least one of the variables is qualitative, you can use **comparative pie charts** or **bar charts**.

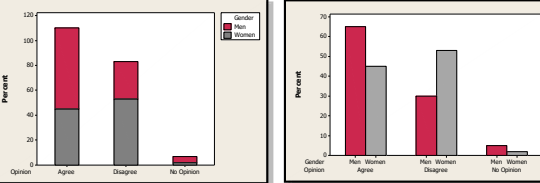
Do you think that men and women are treated equally in the workplace?

Variable #1 = Opinion
Variable #2 = Gender



3 - 4 Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

 **Comparative Bar Charts**




- Stacked Bar Chart
- Side-by-Side Bar Chart

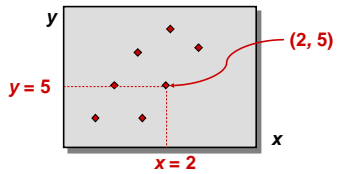
Describe the relationship between opinion and gender:

More women than men feel that they are not treated equally in the workplace.

3 - 5 Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

 **Two Quantitative Variables**

When both of the variables are quantitative, call one variable **x** and the other **y**. A single measurement is a pair of numbers (x, y) that can be plotted using a two-dimensional graph called a **scatterplot**.



3 - 6 Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

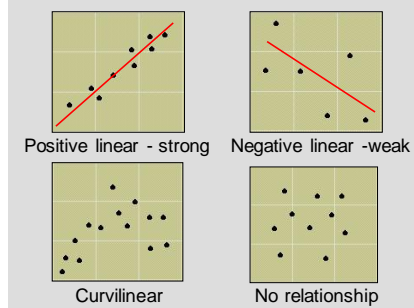
Describing the Scatterplot

- Relationship between two variables
 - What **pattern** or **form** do you see?
 - Straight line upward or downward
 - Curve or
 - No pattern at all
 - How **strong** is the pattern?
 - Strong or weak
 - Are there any **unusual observations**?
 - Clusters or outliers

3 - 7

Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

Examples



3 - 8

Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

Numerical Measures for Two Quantitative Variables

- Assume that the two variables x and y exhibit a **linear pattern** or **form**.
- There are two numerical measures to describe
 - The **strength** and **direction** of the relationship between x and y .
 - The **form** of the relationship.

3 - 9

Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

The Correlation Coefficient

- The strength and direction of the relationship between x and y are measured using the **correlation coefficient, r** .

$$r = \frac{s_{xy}}{s_x s_y} \quad \text{where} \quad s_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

$$s_x = \text{standard deviation of the } x\text{'s} \quad s_y = \text{standard deviation of the } y\text{'s}$$

$$s_{xy} = \frac{\sum x_i y_i - (\sum x_i)(\sum y_i)}{n - 1}$$

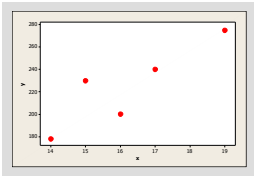
3 - 10

Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

Example

- Living area x and selling price y of 5 homes.

Residence	1	2	3	4	5
x (thousand sq ft)	14	15	17	19	16
y (\$000)	178	230	240	275	200



•The scatterplot indicates a positive linear relationship.

3 - 11

Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

Example

Calculate

$$\bar{x} = 16.2 \quad s_x = 1.924$$

$$\bar{y} = 224.6 \quad s_y = 37.360$$

$$s_{xy} = \frac{\sum x_i y_i - (\sum x_i)(\sum y_i)}{n - 1}$$

$$= \frac{18447 - (81)(1123)}{5} = 63.6$$

x	y	xy
14	178	2492
15	230	3450
17	240	4080
19	275	5225
16	200	3200
81	1123	18447

$$r = \frac{s_{xy}}{s_x s_y}$$

$$= \frac{63.6}{1.924(37.36)} = .885$$

3 - 12

Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

Interpreting r

$$-1 \leq r \leq 1$$

Sign of r indicates direction of the linear relationship.

$$r \approx 0$$

Weak relationship; random scatter of points

$$r \approx 1 \text{ or } -1$$

Strong relationship; either positive or negative

$$r = 1 \text{ or } -1$$

All points fall exactly on a straight line.

3 - 13

Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

The Regression Line

- Sometimes x and y are related in a particular way—the value of y depends on the value of x .
 - y = dependent variable
 - x = independent variable
- The form of the linear relationship between x and y can be described by fitting a line as best we can through the points. This is the regression line, (also known as least-squares line)

$$y = a + bx.$$

- a = y -intercept of the line
- b = slope of the line

3 - 14

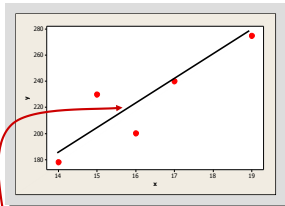
Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

The Regression Line

To find the slope and y -intercept of the best fitting line, use:

$$b = r \frac{s_y}{s_x}$$

$$a = \bar{y} - b\bar{x}$$



The least squares (regression) line is

$$y = a + bx$$

3 - 15

Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

Example

Recall

$$\bar{x} = 16.2 \quad s_x = 1.9235$$

$$\bar{y} = 224.6 \quad s_y = 37.3604$$

$$r = .885$$

x	y	xy
14	178	2492
15	230	3450
17	240	4080
19	275	5225
16	200	3200
81	1123	18447

$$b = r \frac{s_y}{s_x} = (.885) \frac{37.3604}{1.9235} = 17.189$$

$$a = \bar{y} - b\bar{x} = 224.6 - 17.189(16.2) = -53.86$$

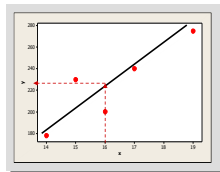
$$\text{RegressionLine: } y = -53.86 + 17.189x$$

3 - 16

Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

Example

Predict the selling price for another residence with 1600 square feet of living area.



$$\text{Predict: } y = -53.86 + 17.189x$$

$$= -53.86 + 17.189(16) = 221.16 \text{ or } \$221,160$$

3 - 17

Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

Key Concepts

I. Bivariate Data

- Both qualitative and quantitative variables
- Describing each variable separately
- Describing the relationship between the variables

II. Describing Two Qualitative Variables

- Side-by-Side pie charts
- Comparative line charts
- Comparative bar charts
 - ✓ Side-by-Side
 - ✓ Stacked
- Relative frequencies to describe the relationship between the two variables.

3 - 18

Mujdat Soyuturk, Probability and Statistics, Spring 2021, Marmara University

Key Concepts



III. Describing Two Quantitative Variables

1. Scatterplots

- ✓ Linear or nonlinear pattern
- ✓ Strength of relationship
- ✓ Unusual observations; clusters and outliers

2. Covariance and correlation coefficient

3. The best fitting line

- ✓ Calculating the slope and y-intercept
- ✓ Graphing the line
- ✓ Using the line for prediction