

College of Sciences  
Department of Statistics & Operations Research

# STAT 109

# Biostatistics

August 2024

**NOTE:** These class notes were originally prepared by Professor Abdullah al-Shiha, and have been subsequently revised and improved by: Dr. Razan Alsehibani, Sana AbuNasrah, Eidah Al-Enazi, Zahra M Kaabi, Kholoud Basalim, and Maram Al-Shammari.

# Content

- ❑ Chapter 1: Getting Acquainted with Biostatistics
- ❑ Chapter 2: Strategies for Understanding the Meaning of Data
- ❑ Chapter 3: Probability as the Basis of Statistical Inference
- ❑ Chapter 4: Probabilistic Features of Certain Data Distribution (Probability Distributions)
- ❑ Chapter 5: Probabilistic Features of the Distributions of Certain Sample Statistics
- ❑ Chapter 6: Using Sample Data to Make Estimations about Population Parameters
- ❑ Chapter 7: Using Sample Statistics to Test Hypotheses about Population Parameters

# Chapter 1:

## Getting Acquainted with Biostatistics

1.1 Introduction

1.2 Some Basic Concept

1.3 Sampling and Statistical Inference

# Introduction

1) How to organize, summarize, and describe data.  
(Descriptive Statistics)

2) How to reach decisions about a large body of data by examine only a small part of the data.  
(Inferential Statistics)

# Some Basic Concepts

## Data:

- 1) Quantitative data (numbers: weights, ages, ...).
- 2) Qualitative data (words: nationalities, occupations, ...).

## Statistics:

- 1) Collection, organization, summarization, and analysis of data (Descriptive Statistics).
- 2) Drawing of inferences and conclusions about a body of data (population) when only a part of the data (sample) is observed (Inferential Statistics).

## Biostatistics:

When the data is obtained from the biological sciences and medicine, we use the term "biostatistics".

## Sources of Data:

- 1) Routinely kept records.
- 2) Surveys.
- 3) Experiments.
- 4) External sources (published reports, data bank, ...).

## Population:

- A population is the largest collection of entities (elements or individuals) in which we are interested at a particular time and about which we want to draw some conclusions.
- When we take a measurement of some variable on each of the entities in a population, we generate a population of values of that variable.

## Population Size ( $N$ ):

The number of elements in the population is called the population size and is denoted by  $N$ .

## Sample:

- A sample is a part of a population.
- From the population, we select various elements on which we collect our data. This part of the population on which we collect data is called the sample.

## Sample Size ( $n$ ):

The number of elements in the sample is called the sample size and is denoted by  $n$ .

## Example :

Suppose that we are interested in studying the characteristics of the weights of the students enrolled in the college of engineering at KSU. If we randomly select 50 students among the students of the college of engineering at KSU and measure their weights, then the weights of these 50 students form our sample.

### Population :

All students enrolled in the college of engineering at KSU.

### Sample :

50 students.

### Variable :

Weight of students .

## Variables:

The characteristic to be measured on the elements is called variable. The value of the variable varies from element to element.

## Example of Variables:

- 1) No. of patients
- 2) Height
- 3) Sex
- 4) Educational Level

## Types of Variables:

- 1) Quantitative Variables.
- 2) Qualitative Variables.

# 1) Quantitative Variables:

A quantitative variable is a characteristic that can be measured. The values of a quantitative variable are numbers indicating how much or how many of something.

## (a) Discrete Variables:

There are jumps or gaps between the values.

- Family size ( $x = 1, 2, 3, \dots$ )
- Number of patients ( $x = 0, 1, 2, 3, \dots$ )

## (b) Continuous Variables:

There are no gaps between the values. A continuous variable can have any value within a certain interval of values.

- Height ( $140 < x < 190$ )
- Blood sugar level ( $10 < x < 15$ )

## 2) Qualitative Variables.

The values of a qualitative variable are words or attributes indicating to which category an element belong.

### (a) Nominal Qualitative Variables:

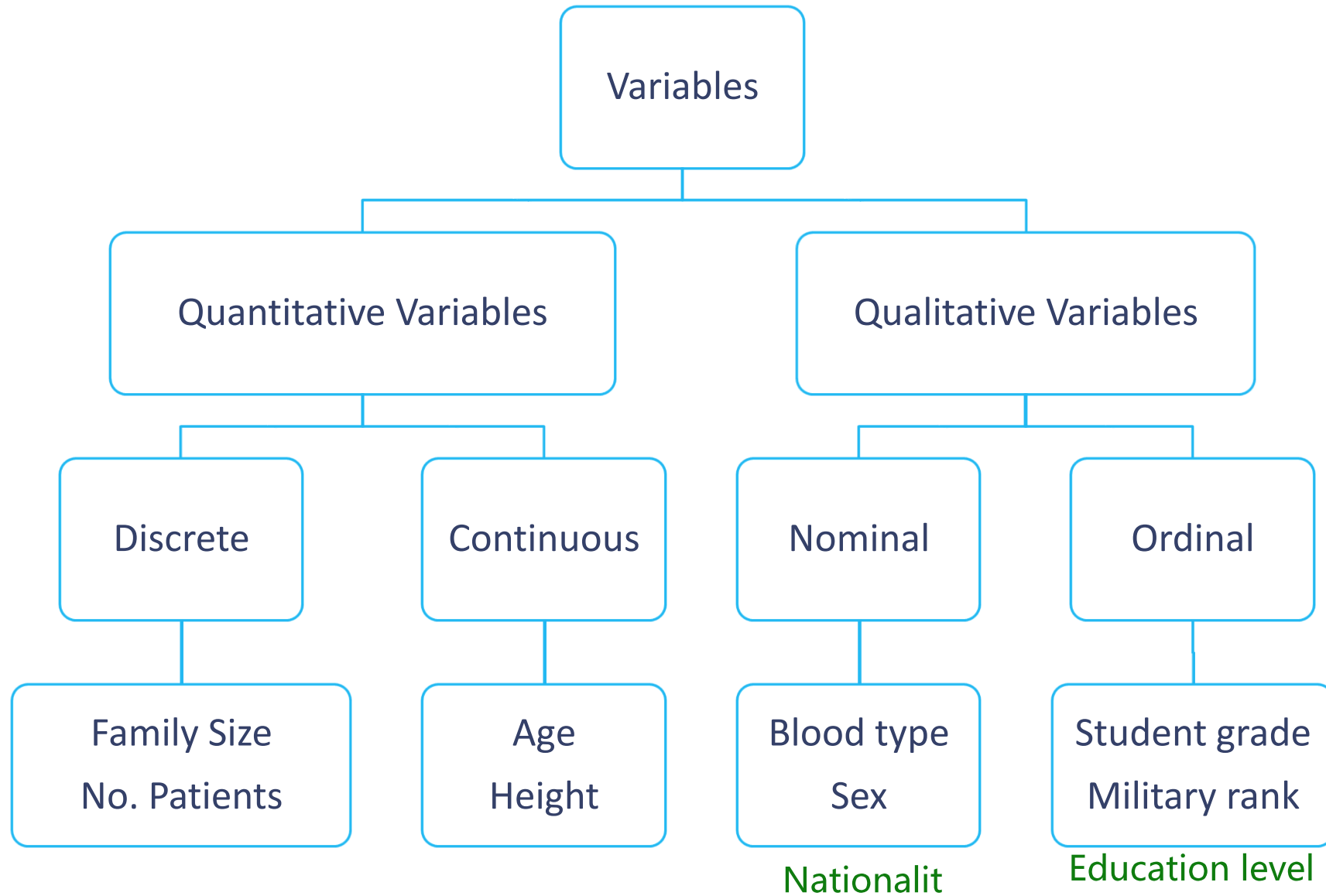
A nominal variable classifies the observations into various mutually exclusive and collectively non-ranked categories.

- Blood type (O, AB, A, B)
- Nationality (Saudi, Egyptian, British, ...)
- Sex (male, female)

## (b) Ordinal Qualitative Variables:

An ordinal variable classifies the observations into various mutually exclusive and collectively ranked categories. The values of an ordinal variable are categories that can be

- Educational level (elementary, intermediate, ...)
- Students grade (A, B, C, D, F)
- Military rank



# Sampling and Statistical Inference

## (1) Simple Random Sampling:

If a sample of size ( $n$ ) is selected from a population of size ( $N$ ) in such a way that each element in the population has the same chance to be selected, the sample is called a simple random sample.

## (2) Stratified Random Sampling:

In this type of sampling, the elements of the population are classified into several homogenous groups (strata). From each group, an independent simple random sample is drawn. The sample resulting from combining these samples is called a stratified random Sample.

# Chapter 2:

## Strategies for Understanding the Meaning of Data

2.1 Introduction

2.2 The Ordered Array

2.3 Grouped Data: The Frequency Distribution

2.4 Descriptive Statistics: Measures of Central Tendency

2.5 Descriptive Statistics: Measures of Dispersion (Measures of Variation)

# Introduction

In this chapter, we learn several techniques for organizing and summarizing data so that we may more easily determine what information they contain.

Summarization techniques involve:

- Frequency distributions
- Descriptive measures

# The Ordered Array

A first step in organizing data is the preparation of an ordered array.

An ordered array is a listing of the values in order of magnitude from the smallest to the largest value.

Example:

Ages of subjects who participate in a study on diabetic:

55 46 58 54 52 69 40 65 53 58

The ordered array is:

40 46 52 53 54 55 58 58 65 69

# Grouped Data: The Frequency Distribution

To group a set of observations, we select a suitable set of contiguous, non overlapping intervals such that each value in the set of observations can be placed in one, and only one, of the intervals. These intervals are called "class intervals".

Example:

Study of the hemoglobin level (g/dl) of a sample of 50 men.

17.0	17.7	15.9	15.2	16.2	17.1	15.7	17.3	<u>13.5</u>	16.3
14.6	15.8	15.3	16.4	13.7	16.2	16.4	16.1	17.0	15.9
14.0	16.2	16.4	14.9	17.8	16.1	15.5	<u>18.3</u>	15.8	16.7
15.9	15.3	13.9	16.8	15.9	16.3	17.4	15.0	17.5	16.1
14.2	16.1	15.7	15.1	17.4	16.5	14.4	16.3	17.3	15.8

Class intervals:

13.0 – 13.9,    14.0 – 14.9,    15.0 – 15.9,

16.0 – 16.9,    17.0 – 17.9,    18.0 – 18.9.

Variable =  $X$  = hemoglobin level (continuous, quantitative).

Sample size =  $n = 50$ .

Min = 13.5

Max = 18.3

The grouped frequency distribution for the hemoglobin level of the 50 men is:

Class Interval (Hemoglobin level)	Frequency (No. of men)
13.0 – 13.9	3
14.0 – 14.9	5
15.0 – 15.9	15
16.0 – 16.9	16
17.0 – 17.9	10
18.0 – 18.9	1
Total	$n = 50$

Notes:

1. Minimum value  $\in$  first interval.
2. Maximum value  $\in$  last interval.
3. The intervals are not overlapped.
4. Each value belongs to one, and only one, interval.
5. Total of the frequencies = the sample size =  $n$ .

Mid-Points of Class Intervals:

$$\text{Mid-point} = \frac{\text{upper limit} + \text{lower limit}}{2}$$

True Class Intervals:

- d is the gap between class intervals
- $d = \text{lower limit} - \text{upper limit of the preceding class intervals}$
- True upper limit = upper limit +  $d/2$
- True lower limit = lower limit -  $d/2$

Example:

- Mid-point of the 1<sup>st</sup> interval =  $(13.0 + 13.9)/2 = 13.45$
- Mid-point of the last interval =  $(18.0 + 18.9)/2 = 18.45$

Class Interval	True Class Interval	Mid-point	Frequency
13.0 – 13.9	12.95 – 13.95	13.45	3
14.0 – 14.9	13.95 – 14.95	14.45	5
15.0 – 15.9	14.95 – 15.95	15.45	15
16.0 – 16.9	15.95 – 16.95	16.45	16
17.0 – 17.9	16.95 – 17.95	17.45	10
18.0 – 18.9	17.95 – 18.95	18.45	1

Note:

1. Mid-point of a class interval is considered as a typical value for all values in the class interval.

For example:

Approximately we may say that: There are 5 observations with the value of 14.45.

2. There are no gaps between true class intervals.

The true upper limit of each true class interval equals to the true lower limit of the following true class interval.

3. Calculating the Mid-point using class Interval or True class Interval yields the same result.

## Cumulative frequency:

- Cumulative frequency of the 1<sup>st</sup> class interval = Frequency.
- Cumulative frequency of a class interval = Frequency + Cumulative frequency of the preceding class interval.

## Relative frequency and Percentage frequency:

- Relative frequency =  $\text{Frequency}/n$ .
- Percentage frequency =  $\text{Relative frequency} \times 100\%$ .

Frq /50

(Frq /50)\*100

Class Interval	Frequency	Cumulative Frequency	Relative Frequency	Cumulative Relative Frequency	Percentage Frequency	Cumulative Percentage Frequency
13.0 – 13.9	3	3	0.06	0.06	6%	6%
14.0 – 14.9	5	8	0.10	0.16	10%	16%
15.0 – 15.9	15	23	0.30	0.46	30%	46%
16.0 – 16.9	16	39	0.32	0.78	32%	78%
17.0 – 17.9	10	49	0.20	0.98	20%	98%
18.0 – 18.9	1	50	0.02	1.00	2%	100%

width of the interval:  $W = \text{lower limit} - \text{lower limit of the preceding interval}$ .

➤ From frequencies:

The number of people whose hemoglobin levels are between 17.0 and 17.9 = 10.

➤ From cumulative frequencies:

The number of people whose hemoglobin levels are less than or equal to 15.9 = 23.

➤ From percentage frequencies:

The percentage of people whose hemoglobin levels are between 17.0 and 17.9 = 20%.

➤ From cumulative percentage frequencies:

The percentage of people whose hemoglobin levels are less than or equal to = 14.9 = 16%.

The percentage of people whose hemoglobin levels are more than 16.9 = 22%.

## Displaying Grouped Frequency Distributions:

For representing frequencies, we may use one of the following graphs:

- The Histogram
- The Frequency Polygon

Example:

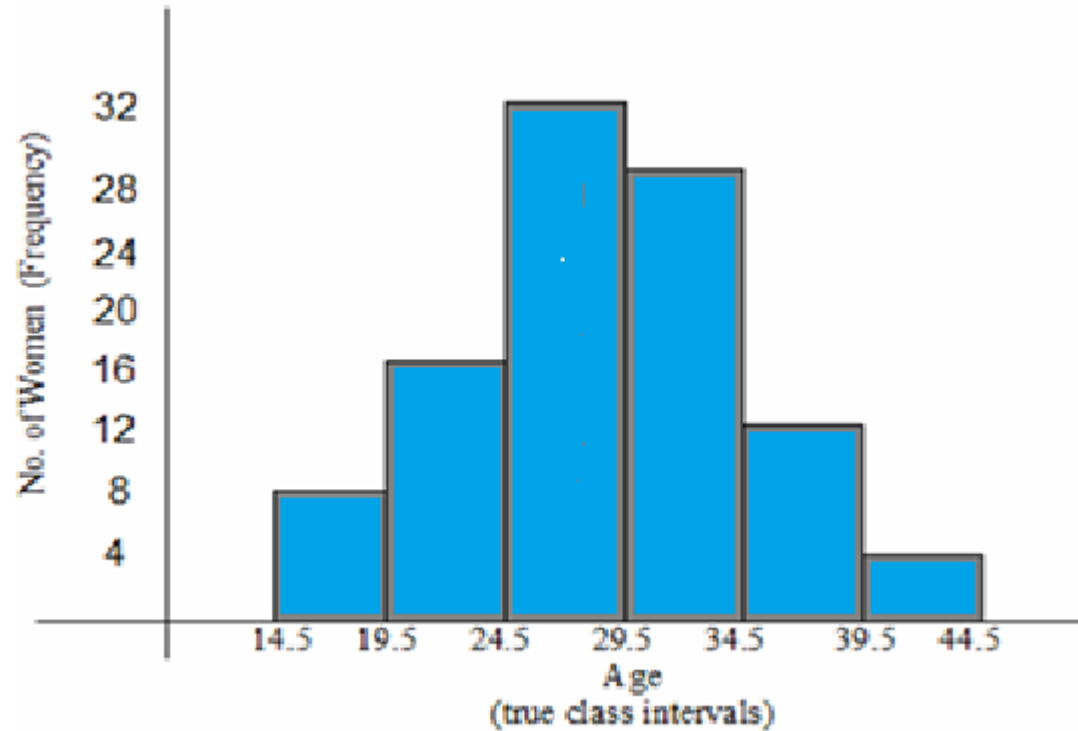
Frequency distribution of the ages of 100 women.

True Class Interval (Age)	Frequency (No. of women)	Cumulative Frequency	Mid-point
14.5 – 19.5	8	8	17
19.5 – 24.5	16	24	22
24.5 – 29.5	32	56	27
29.5 – 34.5	28	84	32
34.5 – 39.5	12	96	37
39.5 – 44.5	4	100	42
Total	$n = 100$		

Width of the interval = true upper limit – true lower limit =  $19.5 - 14.5 = 5$

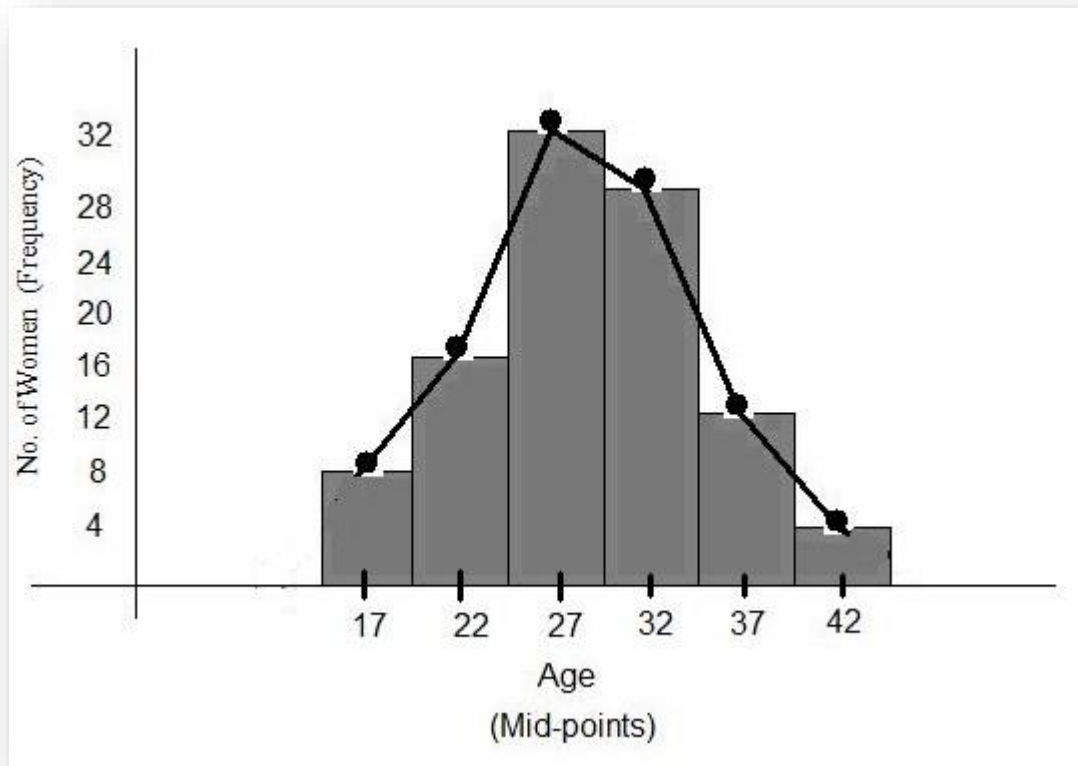
# (1) Histogram: Organizing and Displaying Data using Histogram:

frequency  
or  
Percentage  
frequency  
or  
Relative  
frequency

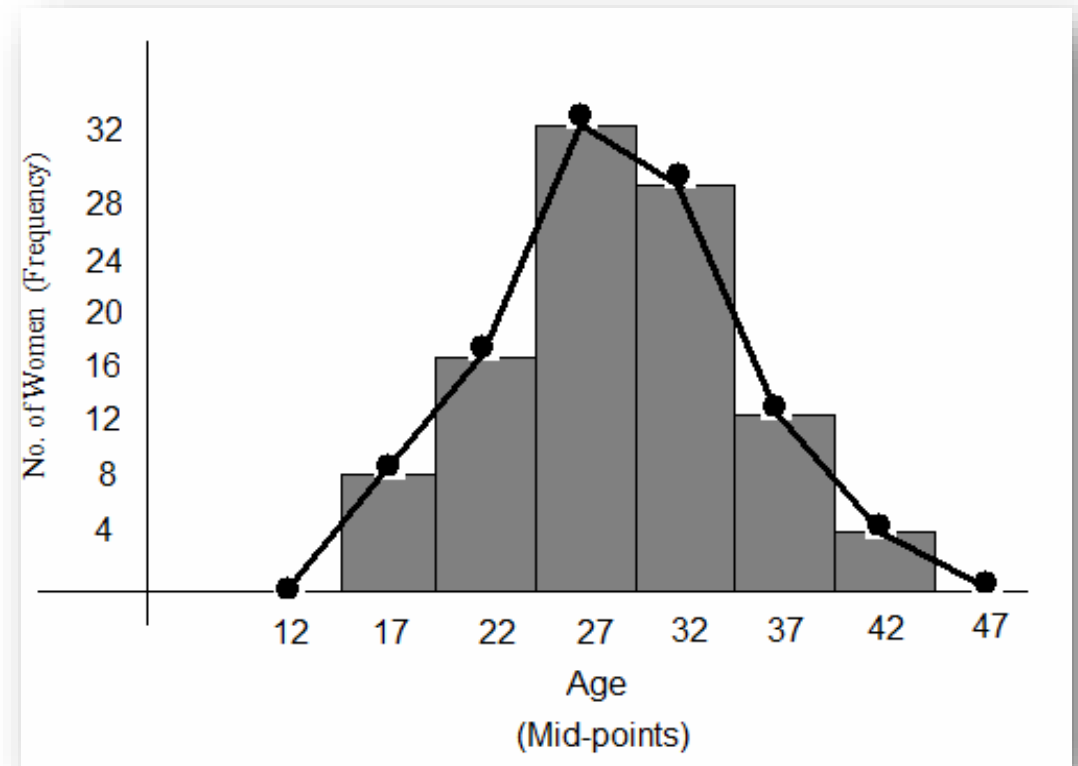


## (2) Frequency Polygon: Organizing and Displaying Data using Polygon

Frequency Polygon(open)



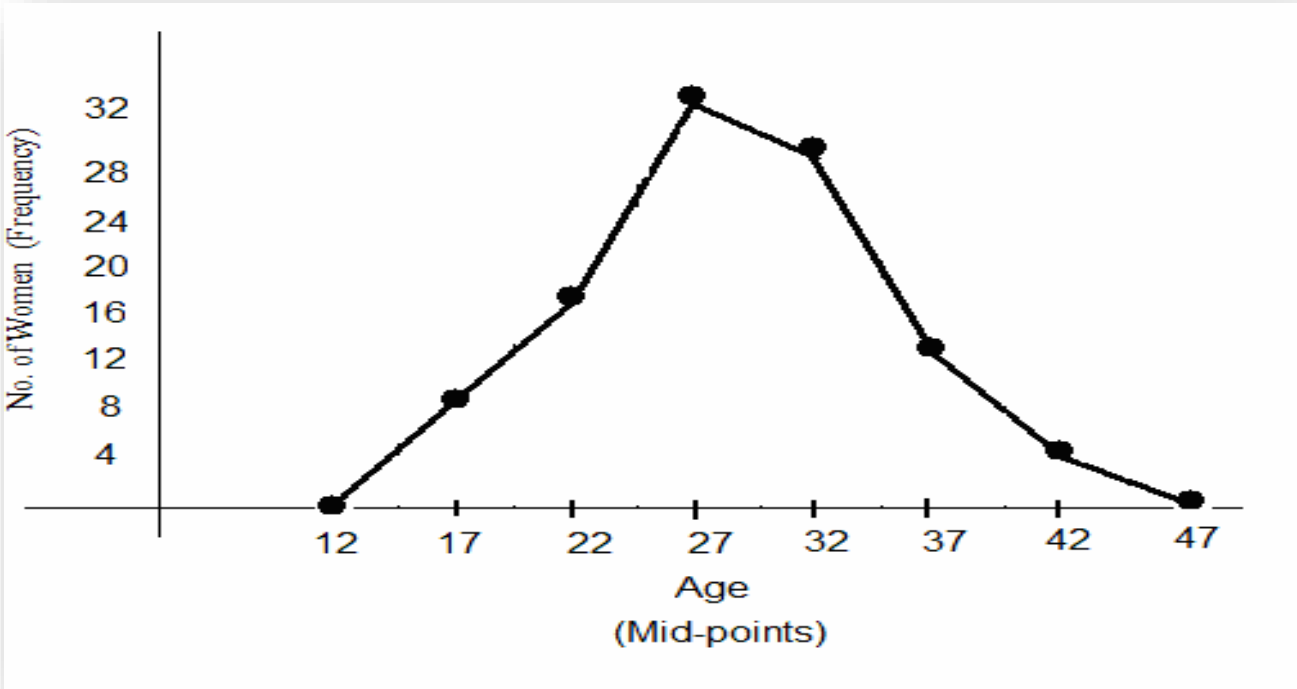
Frequency Polygon(closed)



first Mid point - W  
17- 5= 12

last Mid point + W  
42+ 5= 47

# Frequency Polygon(closed)



# Descriptive Statistics:

Measures of Central Tendency (or location)

Mean ; Mode ; Median

Measures of Dispersion (or Variation)

Range ; Variance ; Standard Deviation ; Coefficient of Variation

We introduce the concept of summarization of the data by means of a single number called "a descriptive measure".

- A descriptive measure computed from the values of a sample is called a "statistic".

- A descriptive measure computed from the values of a population is called a "parameter".

For the variable of interest there are:

- Population values " $N$ ".

Let  $X_1, X_2, \dots, X_N$  be the population values of the variable of interest  
(in general, they are unknown).

The population size =  $N$ .

- Sample of values " $n$ ".

Let  $x_1, x_2, \dots, x_n$  be the sample values of the variable of interest (these values are known).

The sample size =  $n$ .

- A parameter is a measure (or number) obtained from the population values:  $X_1, X_2, \dots, X_N$ .
  - Values of the parameters are unknown in general.
  - We are interested to know true values of the parameters.
  
- A statistic is a measure (or number) obtained from the sample values:  $x_1, x_2, \dots, x_n$ .
  - Values of statistic are known in general.
  - Since parameters are unknown, statistics are used to approximate (estimate) parameters.

# Measures of Central Tendency: (or measures of location):

The commonly used measures of central tendency are:

- The mean
  - The median
  - The mode.
- 
- The values of a variable often tend to be concentrated around the center of the data.
  - The center of the data can be determined by the measures of central tendency.
  - A measure of central tendency is a typical (or a representative) value of the set of data.

# Mean

The Population mean ( $\mu$ ):

If  $X_1, X_2, \dots, X_N$  are the population values, then the population mean is:

$$\mu = \frac{X_1 + X_2 + \dots + X_N}{N} = \frac{\sum_{i=1}^N X_i}{N} \quad (\text{unit})$$

- The population mean  $\mu$  is a parameter (it is usually unknown, and we are interested to know its value)

The Sample mean ( $\bar{x}$ ):

If  $x_1, x_2, \dots, x_n$  are the sample values, then the sample mean is:

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{\sum_{i=1}^n x_i}{n} \quad (\text{unit})$$

- The sample mean  $\bar{x}$  is a statistic (it is known – we can calculate it from the sample).
- The sample mean  $\bar{x}$  is used to approximate (estimate) the population mean  $\mu$ .

# Example

Suppose that we have a population of 5 population values:

$$X_1 = 41, \quad X_2 = 30, \quad X_3 = 35, \quad X_4 = 22, \quad X_5 = 27. \quad (N = 5)$$

Suppose that we randomly select a sample of size 3, and the sample values we obtained are:

$$x_1 = 30, \quad x_2 = 35, \quad x_3 = 27. \quad (n = 3)$$

- Calculate the population mean.
- Calculate the sample mean.

# Solution

The population mean is:

$$\mu = \frac{41 + 30 + 35 + 22 + 27}{5} = \frac{155}{5} = 31 \quad (\text{unit})$$

The sample mean is:

$$\bar{x} = \frac{30 + 35 + 27}{3} = \frac{92}{3} = 30.67 \quad (\text{unit})$$

- Notice that  $\bar{x} = 30.67$  is approximately equals to  $\mu = 31$ .
- Note: The unit of the mean is the same as the unit of the data.

# Advantages and disadvantages of the mean:

## Advantages:

- **Simplicity:** The mean is easily understood and easy to compute.
- **Uniqueness:** There is one and only one mean for a given set of data.
- The mean considers all values of the data.

## Disadvantages:

- Extreme values have an influence on the mean. Therefore, the mean may be distorted by extreme values.

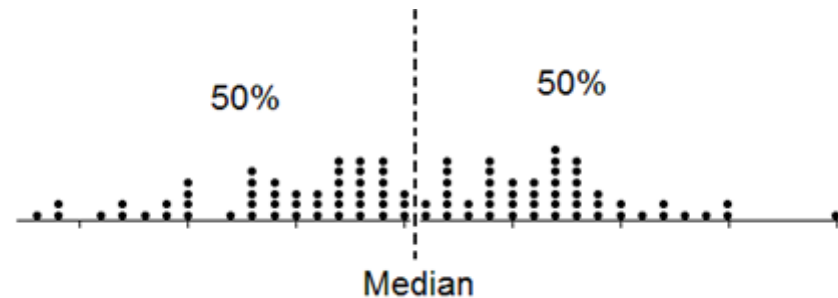
For example:

Sample	Data	Mean
A	2 4 5 7 7 10	5.83
B	2 4 5 7 7 100	20.83

- The mean can only be found for quantitative variables.

# Median

The median of a finite set of numbers is that value which divides the **ordered array** into two equal parts. The numbers in the first part are less than or equal to the median and the numbers in the second part are greater than or equal to the median.



Notice that:

50% (or less) of the data is  $\leq$  Median

50% (or less) of the data is  $\geq$  Median

# Calculating the Median

Let  $x_1, x_2, \dots, x_n$  be the sample values. The sample size ( $n$ ) can be odd or even.

- First we order the sample to obtain the ordered array.
- Suppose that the ordered array is:  $y_1, y_2, \dots, y_n$ .
- We compute the rank of the middle value ( $s$ ):

$$rank = \frac{n + 1}{2}$$

- If the sample size ( $n$ ) is an odd number, there is only one value in the middle, and the rank will be an integer:

$$rank = \frac{n + 1}{2} = m \quad (m \text{ is integer})$$

The median is the middle value of the ordered observations, which is:

$$\text{Median} = y_m$$

- If the sample size ( $n$ ) is an even number, there are two values in the middle, and the rank will be an integer plus 0.5:

$$rank = \frac{n + 1}{2} = m + 0.5$$

Therefore, the ranks of the middle values are ( $m$ ) and ( $m + 1$ ).

				Middle value			
Ordered set (smallest to largest)	→	$y_1$	$y_2$	...	$y_m$	...	$y_n$
Rank (or order)	→	1	2	...	$m$	...	n

The median is the mean (average) of the two middle values of the ordered observations:

$$\text{Median} = \frac{y_m + y_{m+1}}{2}$$

				Middle value	Middle value			
Ordered set	→	$y_1$	$y_2$	...	$y_m$	$y_{m+1}$	...	$y_n$
Rank (or order)	→	1	2	...	$m$	$m + 1$	...	$n$

# Example (odd number):

Find the median for the sample values: 10, 54, 21, 38, 53.

**Solution:**

$n = 5$  (odd number).

There is only one value in the middle.

The rank of the middle value is:

$$rank = \frac{n + 1}{2} = \frac{5 + 1}{2} = 3. \quad (m = 3)$$

Ordered set	→	10	21	38	53	54
Rank (or order)	→	1	2	3 ( <i>m</i> )	4	5

Middle value

The median = 38 (unit)

# Example (even number):

Find the median for the sample values: 10, 35, 41, 16, 20, 32.

Solution:

$n = 6$  (even number).

There are two values in the middle.

The rank is:

$$\text{rank} = \frac{n + 1}{2} = \frac{6 + 1}{2} = 3.5 = 3 + 0.5 = m + 0.5. \quad (m = 3)$$

Therefore, the ranks of the middle values are:

$$m = 3 \quad \text{and} \quad m + 1 = 4.$$

Ordered set	→	10	16	20	32	35	41
Rank (or order)	→	1	2	3 ( $m$ )	4 ( $m+1$ )	5	6

The middle values are 20 and 32.

$$\text{The median} = \frac{20 + 32}{2} = \frac{52}{2} = 26 \quad (\text{unit}).$$

Note: The unit of the median is the same as the unit of the data.

# Advantages of the median

Advantages:

- **Simplicity:** The median is easily understood and easy to compute.
- **Uniqueness:** There is only one median for a given set of data.
- The median is not as drastically affected by extreme values as is the mean. (i.e., the median is not affected too much by extreme values).

For example:

Sample	Data	Median
A	9 4 5 9 2 10	7
B	9 4 5 9 2 100	7

# Disadvantages of the median

## Disadvantages:

- The median does not consider all values of the sample. The mean can only be found for quantitative variables.
- In general, the median can only be found for quantitative variables.

However, in some cases, the median can be found for ordinal qualitative variables (with odd sample size).

Fore example: student grades (A, A, B, C, D)

# Mode

The mode of a set of values is that value which occurs most frequently (i.e., with the highest frequency).

- If all values are different or have the same frequencies, there will be no mode.
- A set of data may have more than one mode.

# Example

Data set	Type	Mode(s)
26, 25, 25, 34	Quantitative	25
3, 7, 12, 6, 19	Quantitative	No mode
3, 3, 7, 7, 12, 12, 6, 6, 19, 19	Quantitative	No mode
3, 3, 12, 6, 8, 8	Quantitative	3 and 8
B C A B B B C B B	Qualitative	B
B C A B A B C A C	Qualitative	No mode
B C A B B C B C C	Qualitative	B and C

Note: The unit of the mode is the same as the unit of the data.

# Advantages of the mode

Advantages:

- **Simplicity:** The mode is easily understood and easy to compute.
- The mode is not as drastically affected by extreme values as is the mean. (i.e., the mode is not affected too much by extreme values).

For example:

Sample	Data	Median
A	7 4 5 7 2 10	7
B	7 4 5 7 2 100	7

- The mode may be found for both quantitative and qualitative variables.

# Disadvantages of the mode

## Disadvantages:

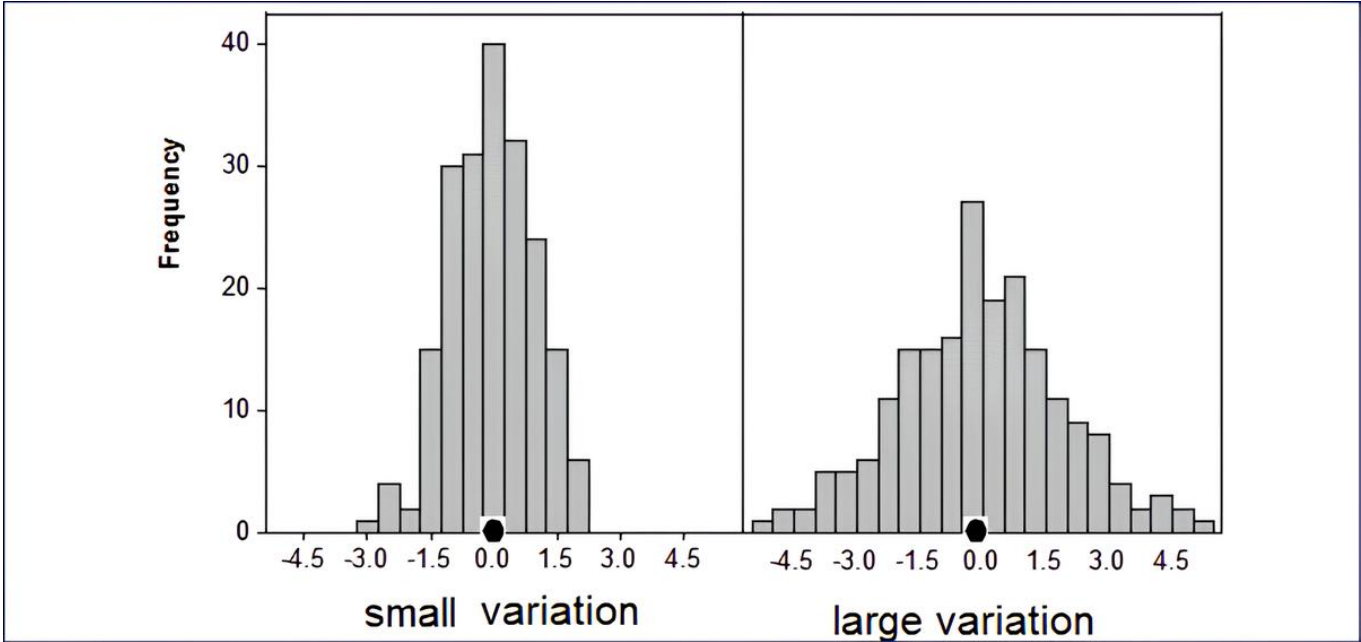
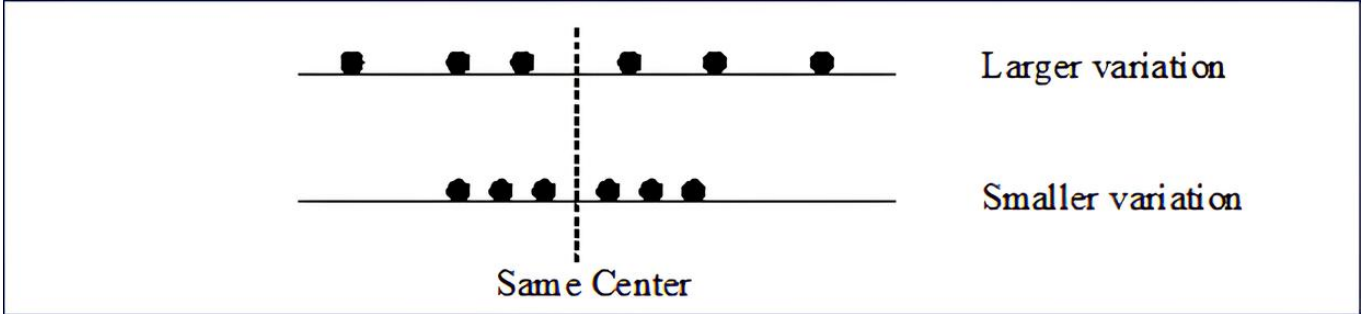
- The mode is not a “good” measure of location, because it depends on a few values of the data.
- The mode does not consider all values of the sample.
- There might be no mode for a data set.
- There might be more than one mode for a data set.

# Measures of Dispersion (Measures of Variation):

The dispersion (variation) of a set of observations refers to the variety that they exhibit. A measure of dispersion conveys information regarding the amount of variability present in a set of data. There are several measures of dispersion, some of which are: Range, Variance, Standard Deviation, and Coefficient of Variation.

The variation or dispersion in a set of values refers to how spread out the values is from each other.

- The dispersion (variation) is small when the values are close together.
- There is no dispersion (no variation) if the values are the same.



# The Range

The Range is the difference between the largest value (Max) and the smallest value (Min).

$$\text{Range (R)} = \text{Max} - \text{Min}$$

Notes:

1. The unit of the range is the same as the unit of the data.
2. The usefulness of the range is limited. The range is a poor measure of the dispersion because it only considers two of the values; however, it plays a significant role in many applications.

# Example

Find the range for the sample values: 26, 25, 35, 27, 29, 29.

**Solution:**

$$\text{Max} = 35.$$

$$\text{Min} = 25.$$

$$\text{Range (R)} = 35 - 25 = 10 \quad (\text{unit}).$$

# The Variance

The variance is one of the most important measures of dispersion.

The variance is a measure that uses the mean as a point of reference.

- The variance of the data is small when the observations are close to the mean.
- The variance of the data is large when the observations are spread out from the mean.
- The variance of the data is zero (no variation) when all observations have the same value (concentrated at the mean).

# Deviations of sample values from the sample mean:

Let  $x_1, x_2, \dots, x_n$  be the sample values, and  $\bar{x}$  be the sample mean.

The deviation of the value from the sample mean  $\bar{x}$  is:

$$x_i - \bar{x}$$

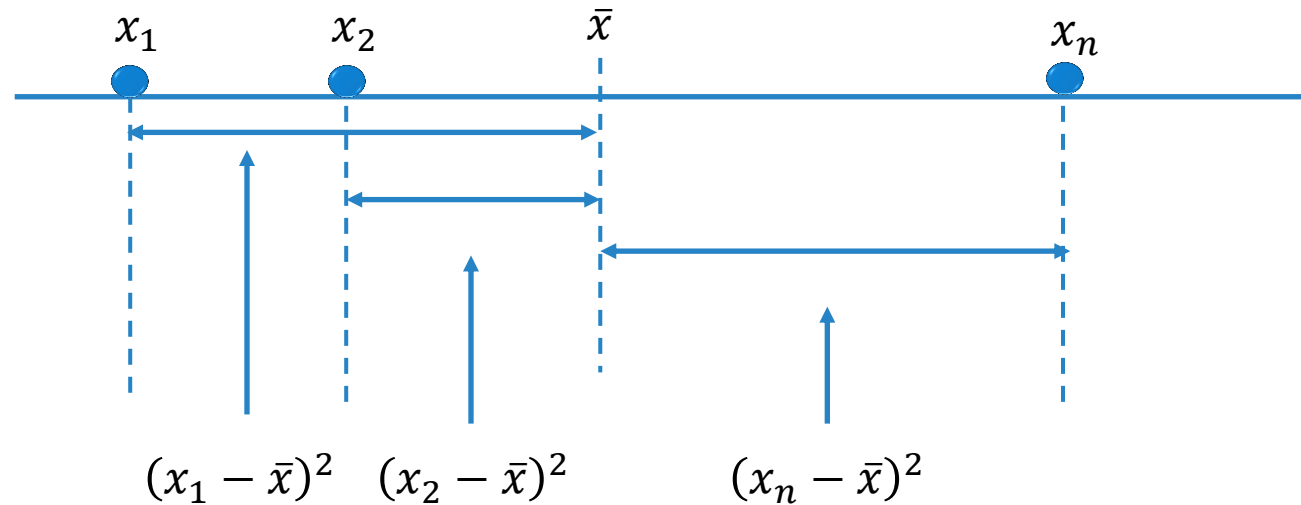
The squared deviation is:

$$(x_i - \bar{x})^2$$

The sum of squared deviations is:

$$\sum_{i=1}^n (x_i - \bar{x})^2$$

The following graph shows the squared deviations of the values from their mean:



# The Population Variance $\sigma^2$ :

(Variance computed from the population)

Let  $X_1, X_2, \dots, X_N$  be the population values. The population variance ( $\sigma^2$ ) is defined by:

$$\sigma^2 = \frac{\sum_{i=1}^N (X_i - \mu)^2}{N} = \frac{(X_1 - \mu)^2 + (X_2 - \mu)^2 + \dots + (X_N - \mu)^2}{N} \quad (\text{unit})^2$$

where,  $\mu = \frac{\sum_{i=1}^N X_i}{N}$  is the population mean, and ( $N$ ) is the population size.

Notes:

- $\sigma^2$  is a parameter because it is obtained from the population values (it is unknown in general).
- $\sigma^2 \geq 0$

# The Sample Variance $S^2$ :

(Variance computed from the sample)

Let  $x_1, x_2, \dots, x_n$  be the sample values. The sample variance ( $S^2$ ) is defined by:

$$S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1} = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n - 1} \quad (\text{unit})^2$$

where,  $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$  is the sample mean, and (n) is the sample size.

## Notes:

- $S^2$  is a statistic because it is obtained from the sample values (it is known).
- $S^2$  is used to approximate (estimate)  $\sigma^2$ .
- $S^2 \geq 0$
- $S^2 = 0$   $\left\{ \begin{array}{l} \text{all observation have the same value} \\ \text{there is no disperation (no variation)} \end{array} \right.$

# Example

We want to compute the sample variance of the following sample values: 10, 21, 33, 53, 54.

Solution:

$$n = 5$$

$$\bar{x} = \frac{\sum_{i=1}^5 x_i}{5} = \frac{10+21+33+53+54}{5} = \frac{171}{5} = 34.2$$

$$S^2 = \frac{\sum_{i=1}^5 (x_i - \bar{x})^2}{5-1} = \frac{(10-34.2)^2 + (21-34.2)^2 + (33-34.2)^2 + (53-34.2)^2 + (54-34.2)^2}{4} = \frac{1506.8}{4} = 376.7 \quad (\text{unit})^2$$

Another Method for calculating sample variance:

$x_i$	$x_i - \bar{x} = (x_i - 34.2)$	$(x_i - \bar{x})^2 = (x_i - 34.2)^2$
10	-24.2	585.64
21	-13.2	174.24
33	-1.2	1.44
53	18.8	353.44
54	19.8	392.04
$\sum_{i=1}^5 x_i = 171$	$\sum_{i=1}^5 (x_i - \bar{x}) = 0$	$\sum_{i=1}^5 (x_i - \bar{x})^2 = 1506.8$

$$\bar{x} = \frac{\sum_{i=1}^5 x_i}{5} = \frac{171}{5} = 34.2 \quad \text{and} \quad s^2 = \frac{1506.8}{4} = 376.7$$

# Standard Deviation

The variance represents squared units, therefore, is not appropriate measure of dispersion when we wish to express the concept of dispersion in terms of the original unit.

- The standard deviation is another measure of dispersion.
- The standard deviation is the square root of the variance.
- The standard deviation is expressed in the original unit of the data.

(1) Population standard deviation is:  $\sigma = \sqrt{\sigma^2}$  (unit)

(2) Sample standard deviation is:  $S = \sqrt{S^2} = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$  (unit)

# Example

For the previous example, the sample standard deviation is:

$$S = \sqrt{S^2} = \sqrt{376.7} = 19.41$$

# Coefficient of Variation (C.V.)

- The variance and the standard deviation are useful as measures of variation of the values of a single variable for a single population.
- If we want to compare the variation of two variables, we cannot use the variance or the standard deviation because:
  1. The variables might have different units.
  2. The variables might have different means.
- We need a measure of the relative variation that will not depend on either the units or on how large the values are. This measure is the coefficient of variation (C.V.).
- The coefficient of variation is defined by:

$$C.V. = \frac{S}{\bar{x}} \times 100\%$$

- The C.V. is free of unit (unit-less).

- To compare the variability of two sets of data (i.e., to determine which set is more variable), we need to calculate the following quantities:

	Mean	Standard Deviation	C.V.
1 <sup>st</sup> data set	$\bar{x}_1$	$S_1$	$C.V_1 = \frac{S_1}{\bar{x}_1} 100\%$
2 <sup>nd</sup> data set	$\bar{x}_2$	$S_2$	$C.V_2 = \frac{S_2}{\bar{x}_2} 100\%$

- The data set with the larger value of CV has larger variation.
- The relative variability of the 1st data set is larger than the relative variability of the 2nd data set if  $C.V_1 > C.V_2$  (and vice versa).

# Example

Suppose we have two data sets:

	Mean	Standard Deviation	C.V.
1 <sup>st</sup> data set	$\bar{x}_1 = 66 \text{ kg}$	$S_1 = 4.5 \text{ kg}$	$C.V_1 = \frac{4.5}{66} 100\% = 6.8\%$
2 <sup>nd</sup> data set	$\bar{x}_2 = 36 \text{ kg}$	$S_2 = 4.5 \text{ kg}$	$C.V_2 = \frac{4.5}{36} 100\% = 12.5\%$

Since  $C.V_2 > C.V_1$ , the relative variability of the 2<sup>nd</sup> data set is larger than the relative variability of the 1<sup>st</sup> data set.

If we use the standard deviation to compare the variability of the two data sets, we will wrongly conclude that the two data sets have the same variability because the standard deviation of both sets is 4.5 kg.

# Chapter 3:

## Probability the Basis of Statistical Inference

3.1 Introduction

3.2 Probability

3.3 Elementary Properties of Probability

3.4 Calculating the Probability of an Event

# General Definitions and Concepts

Probability:

Probability is a measure (or number) used to measure the chance of the occurrence of some event. This number is between 0 and 1.

An Experiment:

An experiment is some procedure (or process) that we do.

Sample Space:

The sample space of an experiment is the set of all possible outcomes of an experiment. Also, it is called the universal set, and is denoted by  $\Omega$ .

## An Event:

Any subset of the sample space  $\Omega$  is called an event.

- $\emptyset \subseteq \Omega$  is an event (impossible event).
- $\Omega \subseteq \Omega$  is an event (sure event).

# Example

Experiment: Selecting a ball from a box containing 6 balls numbered from 1 to 6 and observing the number on the selected ball. This experiment has 6 possible outcomes.

The sample space is:  $\Omega = \{1, 2, 3, 4, 5, 6\}$ .

Consider the following events:

$E_1 =$  getting an even number  $= \{2, 4, 6\} \subseteq \Omega$ .

$E_2 =$  getting a number less than 4  $= \{1, 2, 3\} \subseteq \Omega$ .

$E_3 =$  getting 1 or 3  $= \{1, 3\} \subseteq \Omega$ .

$E_4 =$  getting an odd number  $= \{1, 3, 5\} \subseteq \Omega$ .

$E_5 =$  getting a negative number  $= \{\} = \emptyset \subseteq \Omega$ .

$E_6 =$  getting a number less than 10  $= \{1, 2, 3, 4, 5, 6\} = \Omega \subseteq \Omega$ .

Notation:

$n(\Omega)$  = no. of outcomes (elements) in  $\Omega$ .

$n(E)$  = no. of outcomes (elements) in  $E$ .

## Equally Likely Outcomes:

The outcomes of an experiment are equally likely if the outcomes have the same chance of occurrence.

## Probability of An Event:

If the experiment has  $n(\Omega)$  equally likely outcomes, then the probability of the event  $E$  is denoted by  $P(E)$  and is defined by:

$$P(E) = \frac{n(E)}{n(\Omega)} = \frac{\text{no. of outcomes in } E}{\text{no. of outcomes in } \Omega}$$

# Example

In the ball experiment in the previous example, suppose the ball is selected at random. Determine the probabilities of the following events:

$E_1$  = getting an even number

$E_2$  = getting a number less than 4

$E_3$  = getting 1 or 3

# Solution

$$\Omega = \{1, 2, 3, 4, 5, 6\} ; \quad n(\Omega) = 6.$$

$$E_1 = \{2, 4, 6\} \quad ; \quad n(E_1) = 3$$

$$E_2 = \{1, 2, 3\} \quad ; \quad n(E_2) = 3$$

$$E_3 = \{1, 3\} \quad ; \quad n(E_3) = 2$$

The outcomes are equally likely.

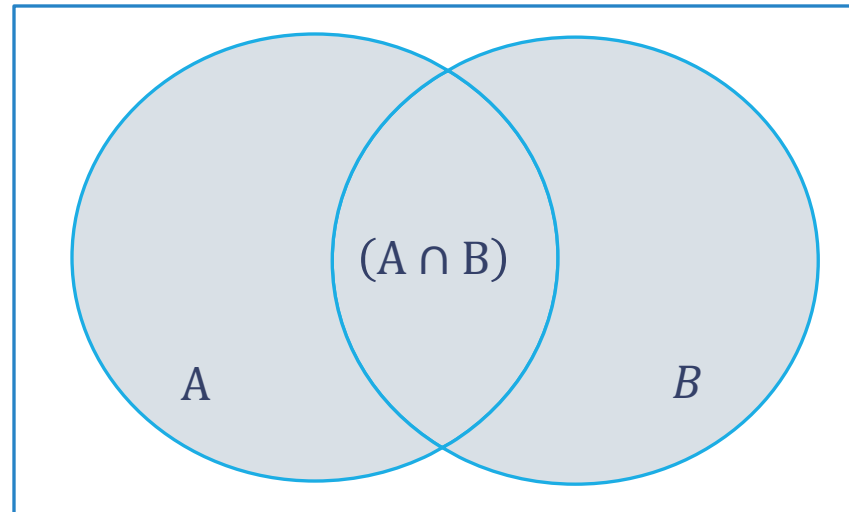
$$\therefore P(E_1) = \frac{3}{6}, \quad P(E_2) = \frac{3}{6}, \quad P(E_3) = \frac{2}{6}.$$

# Some Operations on Events

Let  $A$  and  $B$  be two events defined on the sample space  $\Omega$ .

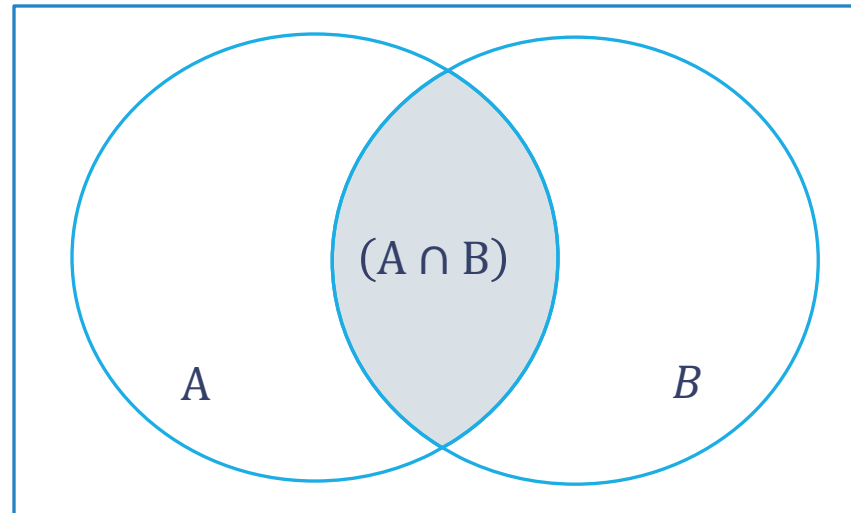
Union of two events:  $(A \cup B)$  or  $(A + B)$

The event  $A \cup B$  consists of all outcomes in  $A$  **or** in  $B$  **or** in both  $A$  and  $B$ . The event  $A \cup B$  occurs if  $A$  occurs, **or**  $B$  occurs, **or** both  $A$  and  $B$  occur.



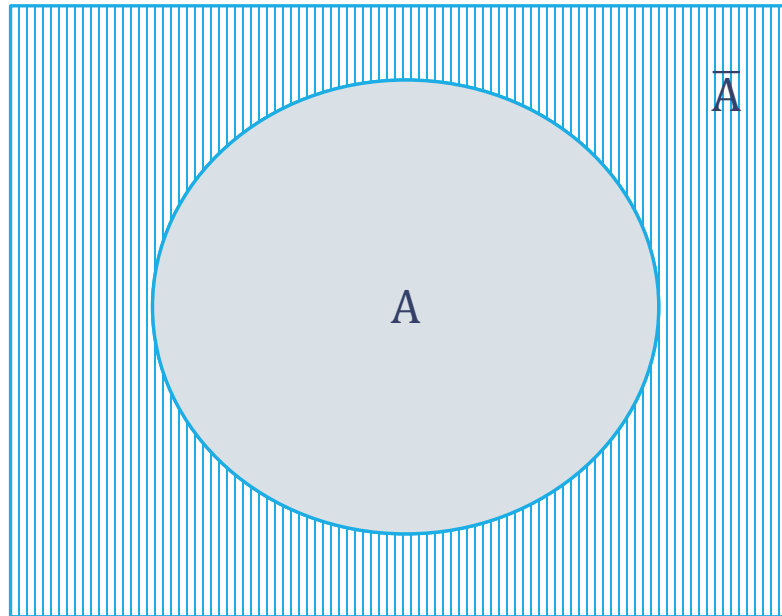
Intersection of two events:  $(A \cap B)$

The event  $A \cap B$  consists of all outcomes in both **A and B**. The event  $A \cap B$  occurs if both **A and B** occur.



Complement of an Event:  $(\bar{A})$  or  $(A^C)$  or  $(A')$

The complement of the event  $A$  is denoted by  $\bar{A}$ . The event  $\bar{A}$  consists of all outcomes of  $\Omega$  but are not in  $A$ . The event  $\bar{A}$  occurs if  $A$  does not.



# Example

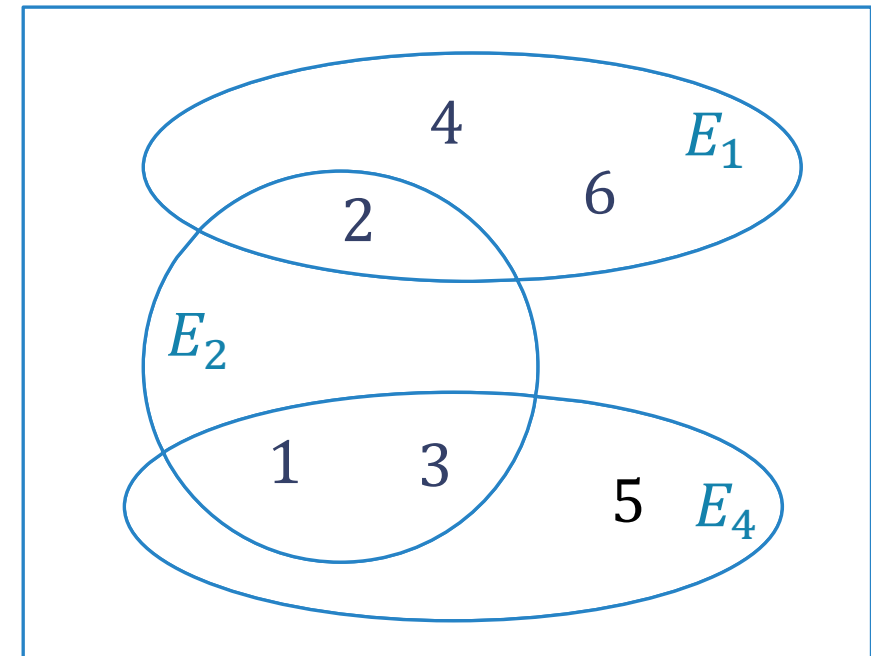
Experiment: Selecting a ball from a box containing 6 balls numbered 1, 2, 3, 4, 5, and 6 randomly.

Define the following events:

$E_1 = \{2, 4, 6\}$  = getting an even number.

$E_2 = \{1, 2, 3\}$  = getting a number less than 4.

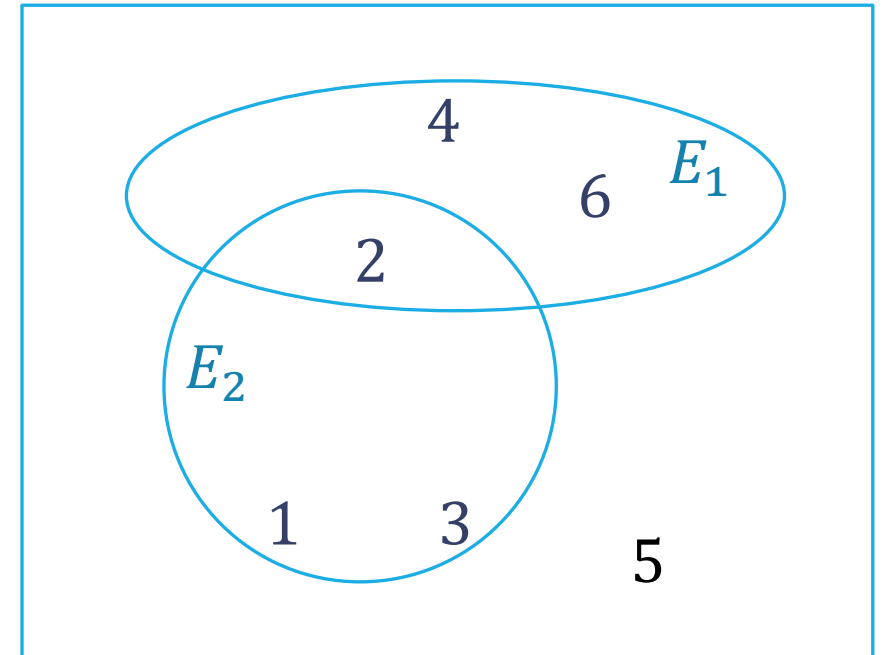
$E_4 = \{1, 3, 5\}$  = getting an odd number.



(1)  $E_1 \cup E_2 =$  getting an even number **or** a number less than 4.  
 $= \{1, 2, 3, 4, 6\}$

what is the probability of getting ....

$$P(E_1 \cup E_2) = \frac{n(E_1 \cup E_2)}{n(\Omega)} = \frac{5}{6}$$



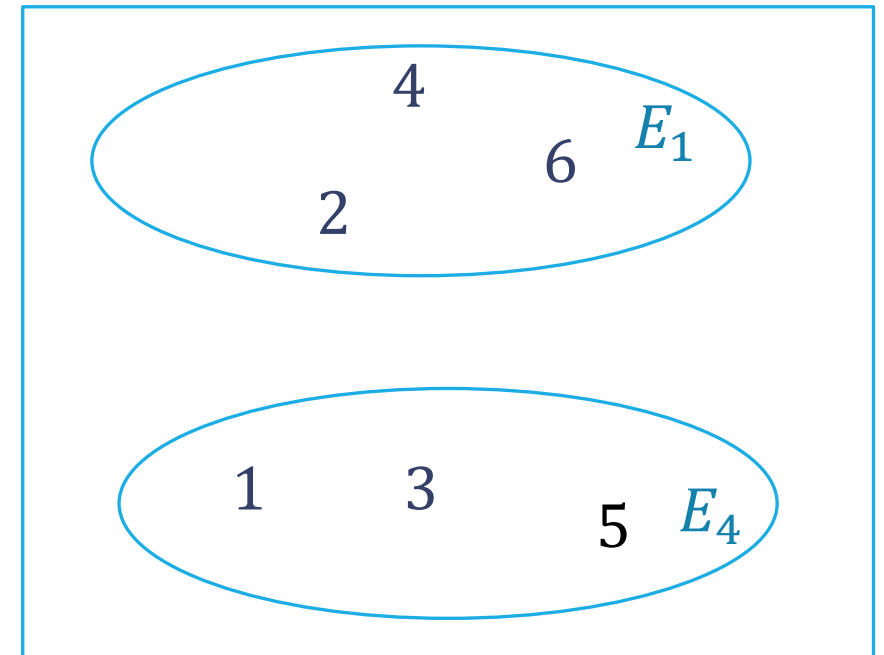
(2)  $E_1 \cup E_4 =$  getting an even number **or** an odd number.  
 $= \{1, 2, 3, 4, 5, 6\}$

$$P(E_1 \cup E_4) = \frac{n(E_1 \cup E_4)}{n(\Omega)} = \frac{6}{6}$$

Note:  $E_1 \cup E_4 = \Omega$ .

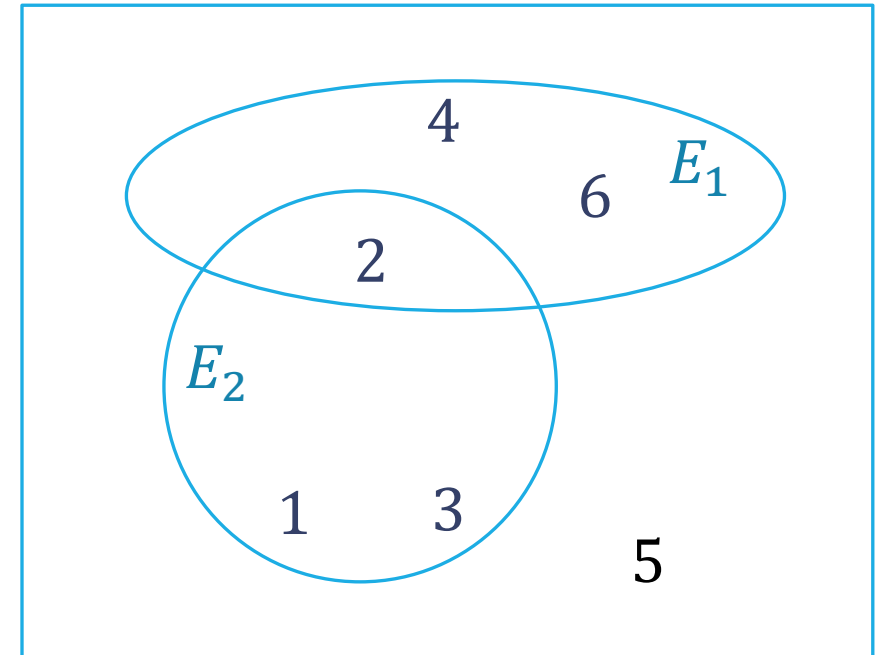
$E_1$  and  $E_4$  are called exhaustive events.

The union of these events gives the whole sample space.



(3)  $E_1 \cap E_2 =$  getting an even number **and** a number less than 4.  
 $= \{2\}$

$$P(E_1 \cap E_2) = \frac{n(E_1 \cap E_2)}{n(\Omega)} = \frac{1}{6}$$



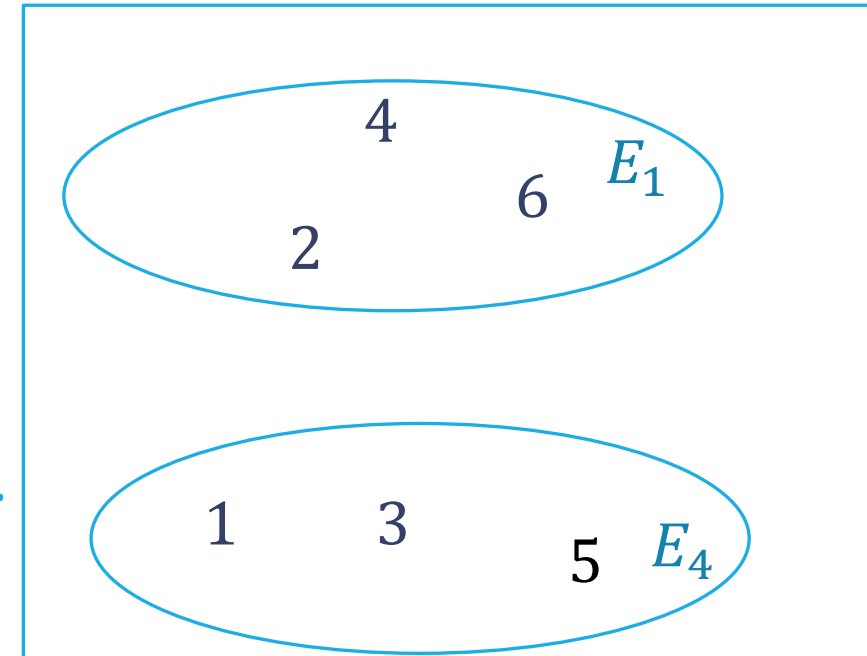
(4)  $E_1 \cap E_4 =$  getting an even number **and** an odd number.  
 $= \{\emptyset\}$

$$P(E_1 \cap E_4) = \frac{n(E_1 \cap E_4)}{n(\Omega)} = \frac{n(\emptyset)}{6} = \frac{0}{6} = 0$$

Note:  $E_1 \cap E_4 = \emptyset$ .

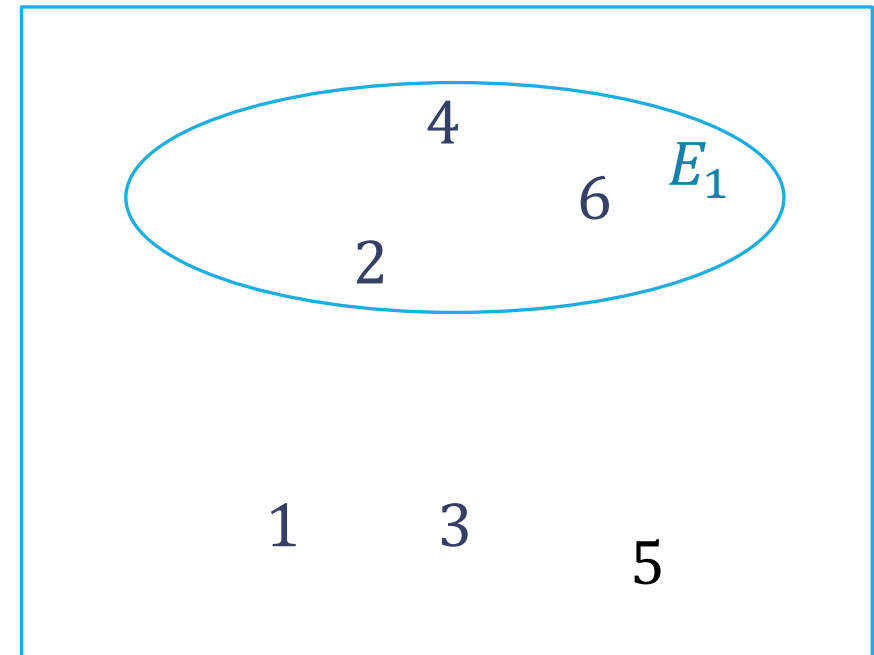
$E_1$  and  $E_4$  are called disjoint (or mutually exclusive) events.

These kinds of events can not occur simultaneously (together at the same time).



(5) The complement of  $E_1$

$$\begin{aligned} \overline{E_1} &= \text{not getting an even number} = \overline{\{2, 4, 6\}} = \{1, 3, 5\} \\ &= \text{not getting an odd number} \\ &= E_4 \end{aligned}$$



Mutually exclusive (disjoint) events:

The events  $A$  and  $B$  are disjoint (or mutually exclusive) if:

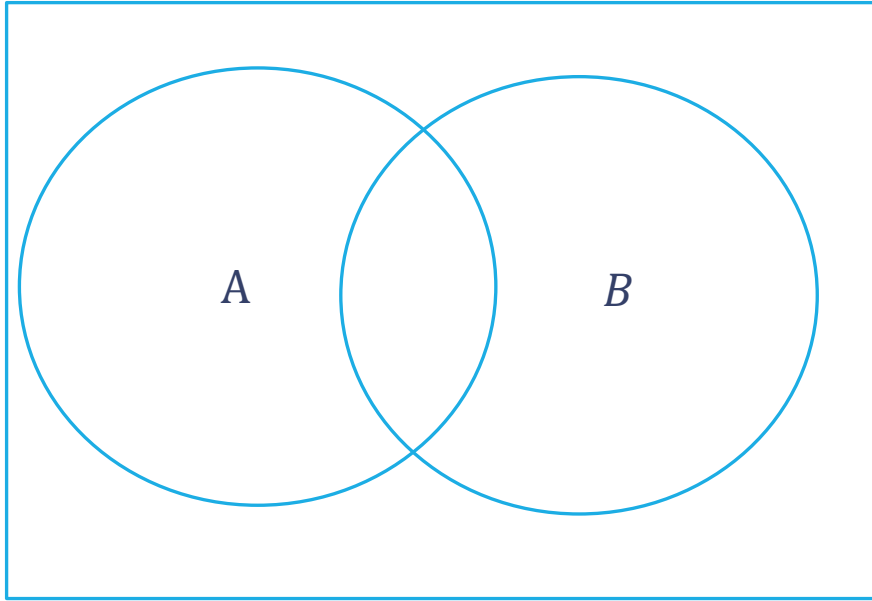
$$A \cap B = \emptyset$$

For this case, it is impossible that both events occur simultaneously (i.e., together in the same time). In this case:

$$(i) P(A \cap B) = 0$$

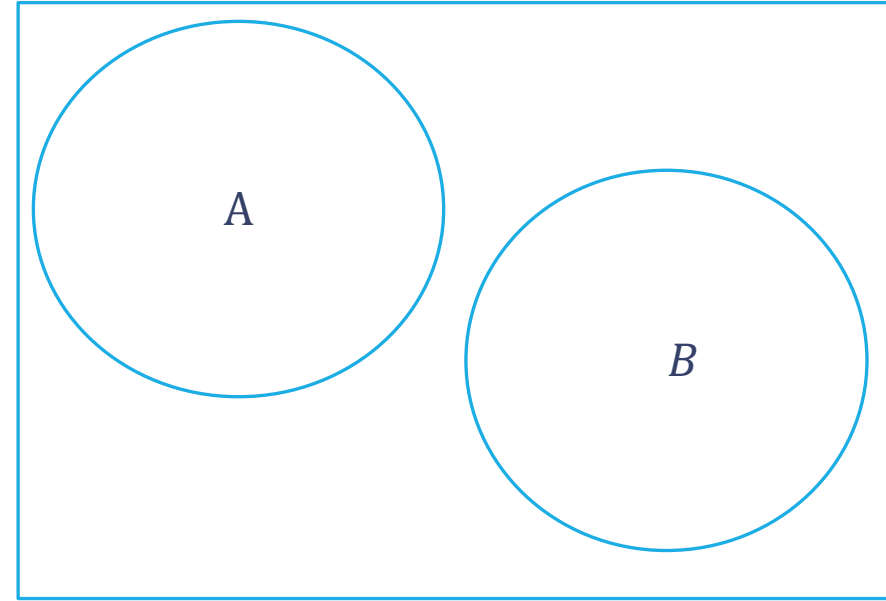
$$(ii) P(A \cup B) = P(A) + P(B)$$

If  $A \cap B \neq \emptyset$ , then  $A$  and  $B$  are not mutually exclusive (not disjoint).



$$A \cap B \neq \emptyset$$

A and B are not mutually exclusive  
(it is possible that both events occur in the same time)



$$A \cap B = \emptyset$$

A and B are mutually exclusive  
(it is impossible that both events occur in the same time)

Exhaustive Events:

The events  $A_1, A_2, \dots, A_n$  are exhaustive events if:

$$A_1 \cup A_2 \cup \dots \cup A_n = \Omega$$

For this case,  $P(A_1 \cup A_2 \cup \dots \cup A_n) = P(\Omega) = 1$

Note:

1)  $A \cup \bar{A} = \Omega$  ( $A$  and  $\bar{A}$  are exhaustive events).

2)  $A \cap \bar{A} = \emptyset$  ( $A$  and  $\bar{A}$  are mutually exclusive (disjoint) events).

3)  $n(\bar{A}) = n(\Omega) - n(A)$

4)  $P(\bar{A}) = 1 - P(A)$

## General Probability Rules:

$$1) 0 \leq P(A) \leq 1$$

$$2) P(\Omega) = 1$$

$$3) P(\emptyset) = 0$$

$$4) P(\bar{A}) = 1 - P(A)$$

## The Addition Rule:

For any two events A and B:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Special Cases:

1) For mutually exclusive (disjoint) events A and B:

$$P(A \cup B) = P(A) + P(B)$$

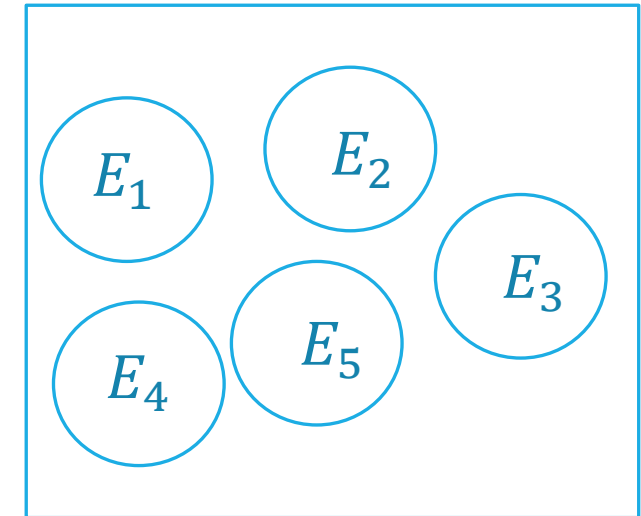
2) For mutually exclusive (disjoint) events  $E_1, E_2, \dots, E_n$ :

$$P(E_1 \cup E_2 \cup \dots \cup E_n) = P(E_1) + P(E_2) + \dots + P(E_n)$$

Note:

If the events  $A_1, A_2, \dots, A_n$  are exhaustive and mutually exclusive (disjoint) events, then:

$$\begin{aligned} P(A_1 \cup A_2 \cup \dots \cup A_n) &= P(A_1) + P(A_2) + \dots + P(A_n) \\ &= P(\Omega) = 1 \end{aligned}$$



# Marginal Probability

Given some variable that can be broken down into (m) categories designated by  $A_1, A_2, \dots, A_m$  and another jointly occurring variable that is broken down into (n) categories designated by  $B_1, B_2, \dots, B_n$ .

	$B_1$	$B_2$	...	$B_n$	Total
$A_1$	$n(A_1 \cap B_1)$	$n(A_1 \cap B_2)$	...	$n(A_1 \cap B_n)$	$n(A_1)$
$A_2$	$n(A_2 \cap B_1)$	$n(A_2 \cap B_2)$	...	$n(A_2 \cap B_n)$	$n(A_2)$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$A_m$	$n(A_m \cap B_1)$	$n(A_m \cap B_2)$	...	$n(A_m \cap B_n)$	$n(A_m)$
Total	$n(B_1)$	$n(B_2)$	...	$n(B_n)$	$n(\Omega)$

( This table contains the number of elements in each event)

	$B_1$	$B_2$	...	$B_n$	Marginal Probability
$A_1$	$P(A_1 \cap B_1)$	$P(A_1 \cap B_2)$	...	$P(A_1 \cap B_n)$	$P(A_1)$
$A_2$	$P(A_2 \cap B_1)$	$P(A_2 \cap B_2)$	...	$P(A_2 \cap B_n)$	$P(A_2)$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$A_m$	$P(A_m \cap B_1)$	$P(A_m \cap B_2)$	...	$P(A_m \cap B_n)$	$P(A_m)$
Marginal Probability	$P(B_1)$	$P(B_2)$	...	$P(B_n)$	1.00

( This table contains the probability of each event)

Joint probability of A and B

The marginal probability of  $A_i$ ,  $P(A_i)$ , is equal to the sum of the joint probabilities of  $A_i$  with all categories of  $B$ . That is:

$$\begin{aligned} P(A_i) &= P(A_i \cap B_1) + P(A_i \cap B_2) + \cdots + P(A_i \cap B_n) \\ &= \sum_{j=1}^n P(A_i \cap B_j) \end{aligned}$$

For example,

$$\begin{aligned} P(A_2) &= P(A_2 \cap B_1) + P(A_2 \cap B_2) + \cdots + P(A_2 \cap B_n) \\ &= \sum_{j=1}^n P(A_2 \cap B_j) \end{aligned}$$

We define the marginal probability of  $B_j$ ,  $P(B_j)$  in a similar way.

# Example

Table of number of elements in each event:

	$B_1$	$B_2$	$B_3$	Total
$A_1$	50	30	70	150
$A_2$	20	70	10	100
$A_3$	30	100	120	250
Total	100	200	200	500

Table of probability of each event:

	$B_1$	$B_2$	$B_3$	Marginal Probability
$A_1$	0.1	0.06	0.14	0.3
$A_2$	0.04	0.14	0.02	0.2
$A_3$	0.06	0.2	0.24	0.5
Marginal Probability	0.2	0.4	0.4	1

$$\begin{aligned}P(A_2) &= P(A_2 \cap B_1) + P(A_2 \cap B_2) + P(A_2 \cap B_3) \\ &= 0.04 + 0.14 + 0.02 \\ &= 0.2\end{aligned}$$

# Example

630 patients are classified as follows:

Blood Type	O ( $E_1$ )	A ( $E_2$ )	B ( $E_3$ )	AB ( $E_4$ )	Total
No. of patients	284	258	63	25	630

**Experiment:** Selecting a patient at random and observe his/her blood type.

This experiment has 630 equally likely outcomes  $n(\Omega) = 630$ .

Define the events:

$E_1$  = The blood type of the selected patient is "O"

$E_2$  = The blood type of the selected patient is "A"

$E_3$  = The blood type of the selected patient is "B"

$E_4$  = The blood type of the selected patient is "AB"

Number of elements in each event:

$$n(E_1) = 284, \quad n(E_2) = 258, \quad n(E_3) = 63, \quad n(E_4) = 25.$$

probability

~~Number~~ of elements in each event:

$$P(E_1) = \frac{284}{630} = 0.4508, \quad P(E_2) = \frac{258}{630} = 0.4095$$

$$P(E_3) = \frac{63}{630} = 0.1, \quad P(E_4) = \frac{25}{630} = 0.0397.$$

# Some Operations on the Events

What is the probability of getting blood type of the selected patients is "A" and "AB" ??

1)  $E_2 \cap E_4 =$  The blood type of the selected patients is "A" and "AB".

$E_2 \cap E_4 = \emptyset$  (disjoint event/mutually exclusive events)

$$P(E_2 \cap E_4) = P(\emptyset) = 0$$

2)  $E_2 \cup E_4 =$  The blood type of the selected patients is "A" or "AB".

$$P(E_2 \cup E_4) = \begin{cases} \frac{n(E_2 \cup E_4)}{n(\Omega)} = \frac{258+25}{630} = \frac{283}{630} = 0.4492 \\ \text{or} \\ P(E_2) + P(E_4) = \frac{258}{630} + \frac{25}{630} = \frac{283}{630} = 0.4492 \end{cases}$$

3)  $\bar{E}_1$  = The blood type of the selected patients **is not** "O".

$$n(\bar{E}_1) = n(\Omega) - n(\cancel{E_1})^{E_1} = 630 - 284 = 346$$

$$P(\bar{E}_1) = \frac{n(\bar{E}_1)}{n(\Omega)} = \frac{346}{630} = 0.5492.$$

Another solute

$$P(\bar{E}_1) = 1 - P(E_1) = 1 - 0.4508 = 0.5492.$$

Notes:

1)  $E_1, E_2, E_3, E_4$  are mutually disjoint,  $E_i \cap E_j = \emptyset$  ( $i \neq j$ ).

2)  $E_1, E_2, E_3, E_4$  are exhaustive events,  $E_1 \cup E_2 \cup E_3 \cup E_4 = \Omega$

# Example

339 physicians are classified based on their ages and smoking habits as follows:

Age	Smoking Habit			Total
	Daily ( $B_1$ )	Occasionally ( $B_2$ )	Not at all ( $B_3$ )	
20 – 29 ( $A_1$ )	31	9	7	47
30 – 39 ( $A_2$ )	110	30	49	189
40 – 49 ( $A_3$ )	29	21	29	79
50+ ( $A_4$ )	6	0	18	24
Total	176	60	103	339

Experiment: Selecting a physician at random.

The number of elements of the sample space is  $n(\Omega) = 339$ .

The outcomes of the experiment are equally likely.

# Some Events

What is the probability that the selected physician's age is 50 years

- $A_3$  = the selected physician is aged (40 – 49).

$$P(A_3) = \frac{n(A_3)}{n(\Omega)} = \frac{79}{339} = 0.2330$$

- $B_2$  = the selected physician smokes occasionally.

$$P(B_2) = \frac{n(B_2)}{n(\Omega)} = \frac{60}{339} = 0.1770$$

- $A_3 \cap B_2$  = the selected physician is aged (40 – 49) **and** smokes occasionally.

$$P(A_3 \cap B_2) = \frac{n(A_3 \cap B_2)}{n(\Omega)} = \frac{21}{339} = 0.06195$$

# Some Events

- $A_3 \cup B_2$  = the selected physician is aged(40 – 49) **or** smokes occasionally (**or** both)

$$\begin{aligned}P(A_3 \cup B_2) &= P(A_3) + P(B_2) - P(A_3 \cap B_2) \\&= \frac{79}{339} + \frac{60}{339} - \frac{21}{339} \\&= 0.233 + 0.177 - 0.06195 \\&= 0.3481\end{aligned}$$

- $\bar{A}_4$  = the selected physician is not 50 years or older.

$$= A_1 \cup A_2 \cup A_3$$

$$P(\bar{A}_4) = 1 - P(A_4) = 1 - \frac{n(A_4)}{n(\Omega)} = 1 - \frac{24}{339} = 0.9292$$

# Some Events

- $A_2 \cup A_3$  = the selected physician is aged (30 – 39) **or** is aged (40 – 49)  
= the selected physician is aged (30 – 49)

Since  $A_2 \cap A_3 = \emptyset$

$$P(A_2 \cup A_3) = \frac{n(A_2 \cup A_3)}{n(\Omega)} = \frac{189 + 79}{339} = \frac{268}{339} = 0.7906$$

Or

$$P(A_2 \cup A_3) = P(A_2) + P(A_3) = \frac{189}{339} + \frac{79}{339} = 0.7906$$

# Some Events

- What is the probability that the selected physician is not (40 – 49) years old and smokes occasionally?

$$\begin{aligned}P(\bar{A}_3 \cap B_2) &= P(A_1 \cap B_2) + P(A_2 \cap B_2) + P(A_4 \cap B_2) \\&= \frac{9}{339} + \frac{30}{339} + \frac{0}{339} = \frac{39}{339} \\&= 0.11504\end{aligned}$$

- What is the probability that the selected physician is (30 – 39) years old and is not a daily smoker?

$$\begin{aligned}P(A_2 \cap \bar{B}_1) &= P(A_2 \cap B_2) + P(A_2 \cap B_3) \\&= \frac{30}{339} + \frac{49}{339} = \frac{79}{339} \\&= 0.2330\end{aligned}$$

# Example

Suppose that there is a population of pregnant women with:

- 10% of the pregnant women delivered prematurely.
- 25% of the pregnant women used some sort of medication.
- 5% of the pregnant women delivered prematurely and used some sort of medication.

Experiment: Selecting a women randomly from this population.

Define the events:

- $D$  = The selected women delivered prematurely.
- $M$  = The selected women used medication.
- $D \cap M$  = The selected women delivered prematurely **and** used some sort of medication.

The complement events:

- $\bar{D}$  = The selected women did not deliver prematurely.
- $\bar{M}$  = The selected women did not use medication.

The probabilities of the given events are:

$$P(D) = 0.1, \quad P(M) = 0.25, \quad P(D \cap M) = 0.05$$

# A Two-way Table: (Percentage given by a two-way table):

Complete the table, then answer the questions:

	$M$	$\bar{M}$	Total
$D$	5	?	10
$\bar{D}$	?	?	?
Total	25	?	100



	$M$	$\bar{M}$	Total
$D$	5	5	10
$\bar{D}$	20	70	90
Total	25	75	100

# Calculating probabilities of some events:

- $D \cup M$  = the selected women delivered prematurely **or** used medication.

$$\begin{aligned} P(D \cup M) &= P(D) + P(M) - P(D \cap M) \\ &= 0.1 + 0.25 - 0.05 \\ &= 0.3 \end{aligned}$$

- $\bar{M}$  = The selected women did **not** use medication.

$$P(\bar{M}) = 1 - P(M) = 1 - 0.25 = 0.75 \quad (\text{By the rule})$$

$$P(\bar{M}) = \frac{75}{100} = 0.75 \quad (\text{From the table})$$

# Calculating probabilities of some events:

- $\bar{D}$  = The selected women did not deliver prematurely.

$$P(\bar{D}) = 1 - P(D) = 1 - 0.10 = 0.90 \quad (\text{By the rule})$$

$$P(\bar{D}) = \frac{90}{100} = 0.9 \quad (\text{From the table})$$

- $\bar{D} \cap \bar{M}$  = The selected women did not deliver prematurely **and** did not use medication

$$P(\bar{D} \cap \bar{M}) = \frac{70}{100} = 0.70 \quad (\text{From the table})$$

- $\bar{D} \cap M$  = The selected women did not deliver prematurely **and** used medication

$$P(\bar{D} \cap M) = \frac{20}{100} = 0.20 \quad (\text{From the table})$$

The selected woman delivered prematurely and did not use medication.

- $D \cap \bar{M}$  = ~~The selected women did not deliver prematurely and used medication~~

$$P(D \cap \bar{M}) = \frac{5}{100} = 0.05 \quad (\text{From the table})$$

The selected woman delivered prematurely or did not use medication.

- $D \cup \bar{M}$  = ~~The selected women did not deliver prematurely and used medication~~

$$\begin{aligned} P(D \cup \bar{M}) &= P(D) + P(\bar{M}) - P(D \cap \bar{M}) \\ &= 0.1 + 0.75 - 0.05 = 0.8 \quad (\text{By the rule}) \end{aligned}$$

- $\bar{D} \cup M$  = The selected women did not deliver prematurely **or** used medication

$$\begin{aligned} P(\bar{D} \cup M) &= P(\bar{D}) + P(M) - P(\bar{D} \cap M) \\ &= 0.9 + 0.25 - 0.20 = 0.95 \quad (\text{By the rule}) \end{aligned}$$

- $\bar{D} \cup \bar{M}$  = The selected women did not deliver prematurely **or** did not use medication

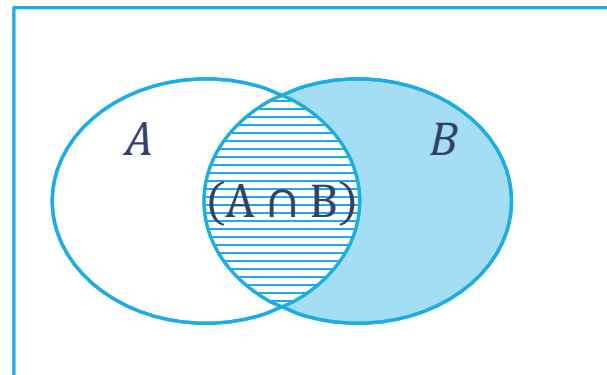
$$\begin{aligned} P(\bar{D} \cup \bar{M}) &= P(\bar{D}) + P(\bar{M}) - P(\bar{D} \cap \bar{M}) \\ &= 0.9 + 0.75 - 0.70 = 0.95 \quad (\text{By the rule}) \end{aligned}$$

# Conditional Probability:

The conditional probability of the event  $A$  when we know that the event  $B$  has already occurred is defined by:

$$P(A|B) = \frac{P(A \cap B)}{\underbrace{P(B)}_{\text{Given / Known}}} ; P(B) \neq 0$$

- $P(A|B)$  = The conditional probability of  $A$  given  $B$ .



Notes: For calculating  $P(A|B)$ , we may use any one of the following:

$$1) \quad P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{n(A \cap B)/n(\Omega)}{P(B)/n(\Omega)}$$

$$2) \quad P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{n(A \cap B)}{P(B)}$$

3) Using the restricted table directly.

### Multiplication Rules of Probability:

For any two events A and B, we have:

$$P(A \cap B) = P(B)P(A|B)$$

$$P(A \cap B) = P(A)P(B|A)$$

# Example

339 physicians are classified based on their ages and smoking habits as follows:

Age	Smoking Habit			Total
	Daily ( $B_1$ )	Occasionally ( $B_2$ )	Not at all ( $B_3$ )	
20 – 29 ( $A_1$ )	31	9	7	47
30 – 39 ( $A_2$ )	110	30	49	189
40 – 49 ( $A_3$ )	29	21	29	79
50+ ( $A_4$ )	6	0	18	24
Total	176	60	103	339

Experiment: Selecting a physician at random.

The number of elements of the sample space is  $n(\Omega) = 339$ .

The outcomes of the experiment are equally likely.

Consider the event  $P(B_1|A_2)$  = the selected physician smokes daily known/given that his age is between 30 and 39.

$$\bullet P(B_1) = \frac{n(B_1)}{n(\Omega)} = \frac{176}{339} = 0.519$$

$$\bullet P(B_1|A_2) = \frac{P(B_1 \cap A_2)}{P(A_2)} = \frac{0.324484}{0.557522} = 0.5820$$

$$\left\{ \begin{array}{l} P(B_1 \cap A_2) = \frac{n(B_1 \cap A_2)}{n(\Omega)} = \frac{110}{339} = 0.324484 \\ P(A_2) = \frac{n(A_2)}{n(\Omega)} = \frac{189}{339} = 0.557522 \end{array} \right.$$

Another solution:

$$P(B_1|A_2) = \frac{n(B_1 \cap A_2)}{n(A_2)} = \frac{110}{189} = 0.5820$$

Notice that:

$$P(B_1) = 0.519$$

$$P(B_1|A_2) = 0.5820$$

- $P(B_1|A_2) > P(B_1)$ ,  $P(B_1) \neq P(B_1|A_2)$ .

What does this mean?

We will answer this question after talking about the concept of independent events.

# Example (Multiplication Rule of Probability)

If we have  $P(A) = 0.9$ ,  $P(B|A) = 0.8$ . Find  $P(A \cap B)$ :

Solution:

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$
$$0.8 = \frac{P(A \cap B)}{0.9}$$
$$\Rightarrow P(A \cap B) = 0.8 \times 0.9 = 0.72$$

# Independent Events

Two events  $A$  and  $B$  are independent if one of the following conditions is satisfied:

- 1)  $P(A|B) = P(A)$
- 2)  $P(B|A) = P(B)$
- 3)  $P(A \cap B) = P(A)P(B)$

Note: The third condition is the multiplication rule of independent events.

# Example

Suppose that A and B are two events such that:

$$P(A) = 0.9, \quad P(B|A) = 0.8, \quad P(A \cap B) = 0.2.$$

Are A , B independent ? We need to satisfies one of the following to prove independent

$$P(A)P(B) = 0.5 \times 0.6 = 0.3 \neq P(A \cap B)$$

$$\text{or} \quad P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{0.2}{0.6} = 0.3333 \neq P(A)$$

$$\text{Or} \quad P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{0.2}{0.5} = 0.4 \neq P(B)$$

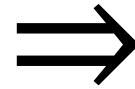
All the rules not satisfies ,then A,B are not independent( A,B are Dependent )

For this example, we may calculate probabilities of all events. We can use a two-way table of the probabilities.

# A Two-way Table:

Complete the table, then answer the questions:

	$B$	$\bar{B}$	Total
$A$	0.2	?	0.5
$\bar{A}$	?	?	?
Total	0.6	?	1



	$B$	$\bar{B}$	Total
$A$	0.2	0.3	0.5
$\bar{A}$	0.4	0.1	0.5
Total	0.6	0.4	1

Q1: Are  $A$  and  $B$  independent events?

Q2: Are  $A$  and  $B$  disjoint events?

Q3: Are  $A$  and  $B$  exhaustive events?

$$P(\bar{A}) = 0.5$$

$$P(\bar{B}) = 0.4$$

$$P(A \cap \bar{B}) = 0.3$$

$$P(\bar{A} \cap B) = 0.4$$

$$P(\bar{A} \cap \bar{B}) = 0.1$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) = 0.5 + 0.6 - 0.2 = 0.9$$

$$P(A \cup \bar{B}) = P(A) + P(\bar{B}) - P(A \cap \bar{B}) = 0.5 + 0.4 - 0.3 = 0.6$$

$$P(\bar{A} \cup B) = \textit{exersice}$$

$$P(\bar{A} \cup \bar{B}) = \textit{exersice}$$

Q1: Are A and B independent events?

**Rule of independent :**

$$P(A \cap B) = P(A)P(B)$$

$$P(A|B) = P(A)$$

$$P(B|A) = P(B)$$

$$P(A \cap B) = P(A)P(B)$$

$$0.2 \neq (0.5)(0.6)$$

$$0.2 \neq 0.3$$

So, A and B are not independent.

Q2: Are A and B disjoint events?

Rule of disjoint :  
 $P(A \cap B) = 0$

$$P(A \cup B) = P(A) + P(B) \quad \text{or} \quad P(A \cap B) = 0$$

$$P(A \cap B) = 0.2 \neq 0$$

So, A and B are not disjoint.

Are A and B exhaustive events?

$$P(A \cup B) = P(\Omega) = 1$$

$$P(A \cup B) = 0.9 \neq 1$$

So, A and B are not exhaustive.

Rule of exhaustive :

$$P(A \cup B) = 1$$

# Example: (Reading Assignment)

Suppose that a dental clinic has 12 nurses classified as follows:

Nurse	1	2	3	4	5	6	7	8	9	10	11	12
Has children	Yes	No	No	No	No	Yes	No	No	Yes	No	No	No
Works at night	No	No	Yes	Yes	Yes	Yes	No	No	Yes	Yes	Yes	Yes

The experiment is to randomly choose one of these nurses. Consider the following events:

$C$  = The chosen nurse has children.

$N$  = The chosen nurse works night shift.

- a) Find the probabilities of the following events:
1. The chosen nurse has children.
  2. The chosen nurse works night shift.
  3. The chosen nurse has children and works night shift.
  4. The chosen nurse has children and does not work night shift.
- b) Find the probability of choosing a nurse who works at night given that she has children.
- c) Are the events C and N independent? Why?
- d) Are the events C and N disjoint? Why?
- e) Sketch the events C and N with their probabilities using Venn diagram.

# Solution

We can classify the nurses as follows:

	$N$ (Night shift)	$\bar{N}$ (No night shift)	Total
$C$ (Has children)	2	1	3
$\bar{C}$ (No Children)	6	3	9
Total	8	4	12

a) The experiment has  $n(\Omega) = 12$  equally likely outcomes.

$$1. P(\text{The chosen nurse has children}) = P(C) = \frac{n(C)}{n(\Omega)} = \frac{3}{12} = 0.25$$

$$2. P(\text{The chosen nurse works night shift}) = P(N) = \frac{n(N)}{n(\Omega)} = \frac{8}{12} = 0.6667$$

$$3. P(\text{The chosen nurse has children and works night shift}) = P(C \cap N) \\ = \frac{n(C \cap N)}{n(\Omega)} = \frac{2}{12} = 0.16667$$

$$4. P(\text{The chosen nurse children and does not work night shift}) = P(C \cap \bar{N}) \\ = \frac{n(C \cap \bar{N})}{n(\Omega)} = \frac{1}{12} = 0.0833$$

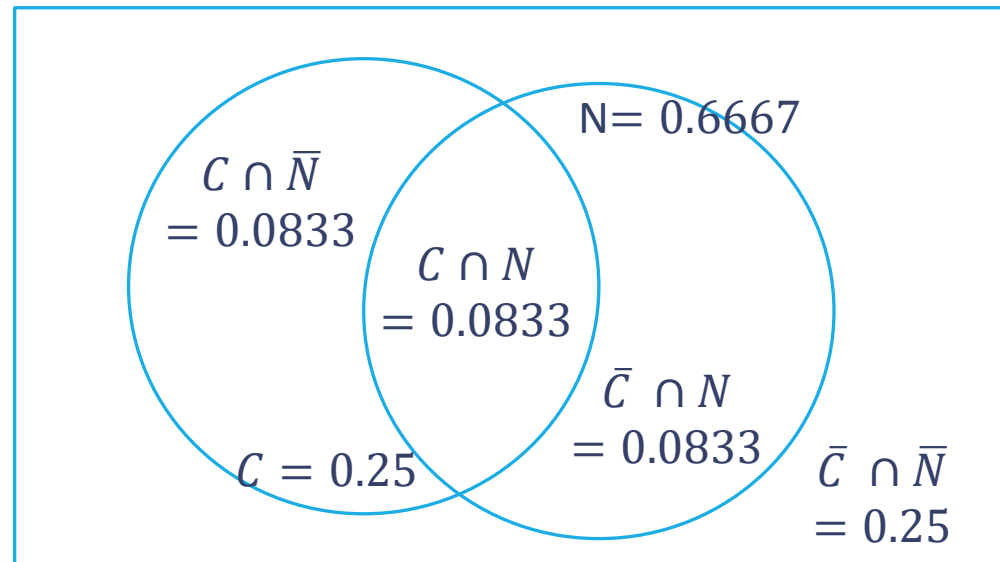
b) The probability of choosing a nurse who works at night given that she has children.

$$P(N|C) = \frac{P(C \cap N)}{P(C)} = \frac{2/12}{0.25} = 0.6667$$

c) The events C and N are independent because  $P(N|C) = P(N)$ .

d) The events C and N are not disjoint because  $n(C \cap N) \neq 0$ . (Note:  $n(C \cap N) = 2$ ).

e) Venn diagram



# Bayes' Theorem, Screening Tests, Sensitivity, Specificity, and Predictive Value Positive and Negative

There are two states regarding the disease and two state regarding the result of the screening test:

State of the Disease:  $\left\{ \begin{array}{l} \text{Present } (D) \\ \text{Absent } (\bar{D}) \end{array} \right.$

Result of the Test:  $\left\{ \begin{array}{l} \text{Positive } (T) \\ \text{Negative } (\bar{T}) \end{array} \right.$

We define the following events of interest:

$D$ : The individual has the disease (presence of the disease).

$\bar{D}$ : The individual does not have the disease (absence of the disease).

$T$ : The individual has a positive screening test result.

$\bar{T}$ : The individual has a negative screening test result.

There are possible situations:

		True status of the disease	
		+ve (D: Present)	-ve ( $\bar{D}$ : Absent)
Result of the test	+ve ( $T$ )	Correct diagnosing	False positive result
	-ve ( $\bar{T}$ )	False negative result	Correct diagnosing

# Definitions of False Results

There are two false results:

1. A false positive result:

This result happens when a test indicates a positive status, when the true status is negative.

$$P(T|\bar{D}) = P(\text{positive result} | \text{absence of the disease})$$

2. A false negative result:

This result happens when a test indicates a negative status, when the true status is positive.

$$P(\bar{T}|D) = P(\text{negative result} | \text{presence of the disease})$$

# Definitions of the Sensitivity and Specificity of the Test

## 1. The Sensitivity:

The sensitivity of a test is the probability of a positive test result given the presence of the disease.

$$P(T|D) = P(\text{positive result of the test} | \text{Presence of the disease})$$

## 2. The Specificity:

The specificity of a test is the probability of a negative test result given the absence of the disease.

$$P(\bar{T}|\bar{D}) = P(\text{negative result of the test} | \text{absence of the disease})$$

To clarify these concepts, suppose we have a sample of  $(n)$  subjects who are cross-classified according to disease status and screening test result as follows:

	Disease		
Test Result	Present ( $D$ )	Absent ( $\bar{D}$ )	Total
Positive ( $T$ )	$a$	$b$	$a + b = n(T)$
Negative ( $\bar{T}$ )	$c$	$d$	$c + d = n(\bar{T})$
Total	$a + c = n(D)$	$b + d = n(\bar{D})$	$n$

For example, there are subjects who have the disease and whose screening test result was positive.

From this table we may compute the following conditional probabilities:

1. The probability of false positive result:  $P(T|\bar{D}) = \frac{n(T \cap \bar{D})}{n(\bar{D})} = \frac{b}{b+d}$

2. The probability of false negative result:  $P(\bar{T}|D) = \frac{n(\bar{T} \cap D)}{n(D)} = \frac{c}{a+c}$

3. The sensitivity of the screening test:  $P(T|D) = \frac{n(T \cap D)}{n(D)} = \frac{a}{a+c}$

4. The specificity of false screening test:  $P(\bar{T}|\bar{D}) = \frac{n(\bar{T} \cap \bar{D})}{n(\bar{D})} = \frac{d}{b+d}$

# Definitions of the Predictive Value Positive and Predictive Value Negative of a Screening Test:

1) The predictive value positive of a screening test:

The predictive value positive is the probability that a subject has the disease, given that the subject has a positive screening test result:

$$\begin{aligned} P(D|T) &= P(\text{the subject has the disease}|\text{positive result}) \\ &= P(\text{presence of the disease}|\text{positive result}) \end{aligned}$$

2) The predictive value negative of a screening test:

The predictive value negative is the probability that a subject does not have the disease, given that the subject has a negative screening test result:

$$\begin{aligned} P(\bar{D}|\bar{T}) &= P(\text{the subject does not have the disease}|\text{negative result}) \\ &= P(\text{absence of the disease}|\text{negative result}) \end{aligned}$$

## Calculating the Predictive Value Positive and Predictive Value Negative

How to calculate  $P(D|T)$  and  $P(\bar{D}|\bar{T})$ :

We calculate these conditional probabilities using the knowledge if:

1) The sensitivity of the test =  $P(T|D)$

2) The specificity of the test =  $P(\bar{T}|\bar{D})$

3) The probability of the relevant disease in the general population,  $P(D)$ . (It is usually obtained from another independent study).

Calculating the Predictive Value Positive,  $P(D|T)$ :

$$P(D|T) = \frac{P(T \cap D)}{P(T)}$$

Therefore, we reach the following version of Bayes' theorem:

$$P(D|T) = \frac{P(T|D)P(D)}{P(T|D)P(D) + P(T|\bar{D})P(\bar{D})} \quad (1)$$

Note:

$P(T|D)$  = sensitivity.

$P(T|\bar{D}) = 1 - P(\bar{T}|\bar{D}) = 1 - \text{specificity}$ .

$P(D)$  = The probability of the relevant disease in the general population.

$P(\bar{D}) = 1 - P(D)$ .

Calculating the Predictive Value Negative,  $P(\bar{D}|\bar{T})$ :

To obtain the predictive value negative of a screening test, we use the following statement of Bays' theorem:

$$P(\bar{D}|\bar{T}) = \frac{P(\bar{T}|\bar{D})P(\bar{D})}{P(\bar{T}|\bar{D})P(\bar{D}) + P(\bar{T}|D)P(D)} \quad (2)$$

Note:

$$P(\bar{T}|\bar{D}) = \text{specificity.}$$

$$P(\bar{T}|D) = 1 - P(T|D) = 1 - \text{sensitivity.}$$

# Example

A medical research team wished to evaluate a proposed screening test for Alzheimer's disease. The test was given to a random sample of 450 patients with Alzheimer's disease and an independent random sample of 500 patients without symptoms of the disease. The two samples were drawn from populations of subjects who were 65 years of age or older. The result are as follows:

Test Result	Alzheimer Disease		Total
	Present ( $D$ )	Absent ( $\bar{D}$ )	
Positive ( $T$ )	436	5	441
Negative ( $\bar{T}$ )	14	495	509
Total	450	500	950

Based on another independent study, it is known that the percentage of patients with Alzheimer's disease (the rate of prevalence of the disease) is 11.3% out of all subjects who were 65 years of age or older.

# Solution:

Using these data we estimate the following quantities:

1) The sensitivity of the test:

$$P(T|D) = \frac{n(T \cap D)}{n(D)} = \frac{436}{450} = 0.9689$$

2) The specificity of the test:

$$P(\bar{T}|\bar{D}) = \frac{n(\bar{T} \cap \bar{D})}{n(\bar{D})} = \frac{495}{500} = 0.99$$

3) The probability of the disease in the general population,  $P(D)$ :

The rate of disease in the relevant general population,  $P(D)$ , cannot be computed from the sample data given in the table. However, it is given that the percentage of patients with Alzheimer's disease is 11.3% out of all subjects who were 65 years of age or older. Therefore  $P(D)$  can be computed to be:

$$P(D) = \frac{11.3\%}{100} = 0.113$$

4) The predictive value positive of the test:

We wish to estimate the probability that a subject who is positive on the test has Alzheimer disease. We use the Bayes' formula of Equation (1):

$$P(D|T) = \frac{P(T|D)P(D)}{P(T|D)P(D) + P(T|\bar{D})P(\bar{D})}$$

From the tabulated data we compute:

$$P(T|D) = \frac{436}{450} = 0.9689 \quad (\text{From part no. 1})$$

$$P(T|\bar{D}) = \frac{n(T \cap \bar{D})}{n(\bar{D})} = \frac{5}{500} = 0.01 = 1 - \text{specificity} = 1 - 0.99$$

Substituting of these result into Equation (1), we get:

$$\begin{aligned}
 P(D|T) &= \frac{(0.9689)P(D)}{(0.9689)P(D)+(0.01)P(\bar{D})} = \frac{\text{Sensitivity} * P(D)}{\text{Sensitivity} * P(D) + (1 - \text{Speicificity}) * P(\bar{D})} \\
 &= \frac{(0.9689)(0.113)}{(0.9689)(0.113) + (0.01)(1 - 0.113)} = 0.93
 \end{aligned}$$

As we see, in this case, the predictive value positive of the test is very high.

5) The predictive value negative of the test:

We wish to estimate the probability that a subject who is negative on the test does not have Alzheimer disease. We use the Bayes' formula of Equation (2):

$$P(\bar{D}|\bar{T}) = \frac{P(\bar{T}|\bar{D})P(\bar{D})}{P(\bar{T}|\bar{D})P(\bar{D}) + P(\bar{T}|D)P(D)}$$

To compute  $P(\bar{D}|\bar{T})$ , we first compute the following probabilities:

$$P(\bar{T}|\bar{D}) = \frac{495}{500} = 0.99 \quad (\text{From part no. 2})$$

$$P(\bar{D}) = 1 - P(D) = 1 - 0.113 = 0.887$$

$$P(\bar{T}|D) = \frac{n(\bar{T} \cap D)}{n(D)} = \frac{14}{45} = 0.0311 = 1 - \text{sensitivity} = 1 - 0.9689$$

Substituting in Equation (2), gives:

$$P(\bar{D}|\bar{T}) = \frac{P(\bar{T}|\bar{D})P(\bar{D})}{P(\bar{T}|\bar{D})P(\bar{D}) + P(\bar{T}|D)P(D)}$$

$$= \frac{\text{Specificity} * P(\bar{D})}{\text{Specificity} * P(\bar{D}) + (1 - \text{Sensitivity}) * P(D)} = \frac{(0.99)(0.887)}{(0.99)(0.887) + (0.0311)(0.113)}$$

$$= 0.996$$

As we see, the predictive value negative is also very high.

# Bayes Theorem pages(48-52)

The result of the test	Has the disease ( $D$ )	Does not have the disease ( $\bar{D}$ )	Total
Positive ( $T$ )	Correct decision $n(T \cap D)$	False decision $n(T \cap \bar{D})$	$n(T)$
	Sensitivity $P(T D) = \frac{n(T \cap D)}{n(D)}$	False positive result $P(T \bar{D}) = \frac{n(T \cap \bar{D})}{n(\bar{D})}$	
Negative ( $\bar{T}$ )	False decision $n(\bar{T} \cap D)$	Correct decision $n(\bar{T} \cap \bar{D})$	$n(\bar{T})$
	False negative result $P(\bar{T} D) = \frac{n(\bar{T} \cap D)}{n(D)}$	Specificity $P(\bar{T} \bar{D}) = \frac{n(\bar{T} \cap \bar{D})}{n(\bar{D})}$	
Total	$n(D)$	$n(\bar{D})$	$n(\Omega)$

$$P(\bar{T}|D) + P(T|D) = 1 \text{ and } P(\bar{T}|\bar{D}) + P(T|\bar{D}) = 1$$

Note that from the table:

$$P(\bar{T} \cap \bar{D}) + P(T|D) = 1, \text{ and } P(\bar{T} \cap \bar{D}) + P(T \cap \bar{D}) = 1.$$

i.e. False negative + Sensitivity = 1, and Specificity + False Positive = 1.

The probability of the relevant disease in the general population,  $P(D)$

[or  $P(\bar{D}) = 1 - P(D)$ ] which is obtained from another independent study.

## Predictive Value Positive

$$\begin{aligned}
 P(D|T) &= \frac{P(T|D) * P(D)}{\text{نفس البسط} + \text{نفس البسط} (D \rightarrow \bar{D})} \\
 &= \frac{P(T|D) * P(D)}{P(T|D) * P(D) + P(T|\bar{D}) * P(\bar{D})} \\
 &= \frac{\text{Sensitivity} * P(D)}{\text{Sensitivity} * P(D) + (1 - \text{Speicificity}) * P(\bar{D})}
 \end{aligned}$$

## Predictive Value Negative

$$\begin{aligned}
 P(\bar{D}|\bar{T}) &= \frac{P(\bar{T}|\bar{D}) * P(\bar{D})}{\text{نفس البسط} + \text{نفس البسط} (\bar{D} \rightarrow D)} \\
 &= \frac{P(\bar{T}|\bar{D}) * P(\bar{D})}{P(\bar{T}|\bar{D}) * P(\bar{D}) + P(\bar{T}|D) * P(D)} \\
 &= \frac{\text{Specificity} * P(\bar{D})}{\text{Specificity} * P(\bar{D}) + (1 - \text{Sensitivity}) * P(D)}
 \end{aligned}$$