

King Saud University

College of Sciences

Department of Statistics & OR

STAT 145

Biostatistics

CHAPTER 1: Getting Acquainted with Biostatistics

1.1 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)

1.2 Some Basic Concepts:

Data: Data is the raw material of statistics.

there are two types of 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 = 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 the weights of students enrolled in the college of engineering at KSU. If we are randomly select 50 students from engineering college at KSU and measure their weights. **Identify the population and the sample in the study?**

*The **population** consists of the weights of all of these students, and our variable of interest is the weight.*

*The weights of these 50 students forms a **sample**.*

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:

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.

Examples:

- (i) Family Size
- (ii) No. of patients
- (iii) Weight
- (iv) height

(a) Discrete Variables:

There are jumps or gaps between the values.

- Examples: - 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.

- Examples: - Height ($140 < x < 190$)
- Blood sugar level ($10 < x < 15$)
- hemoglobin level (g\dl)

- Blood type
- Nationality
- Students Grades
- Educational level

(a) Nominal Qualitative Variables:

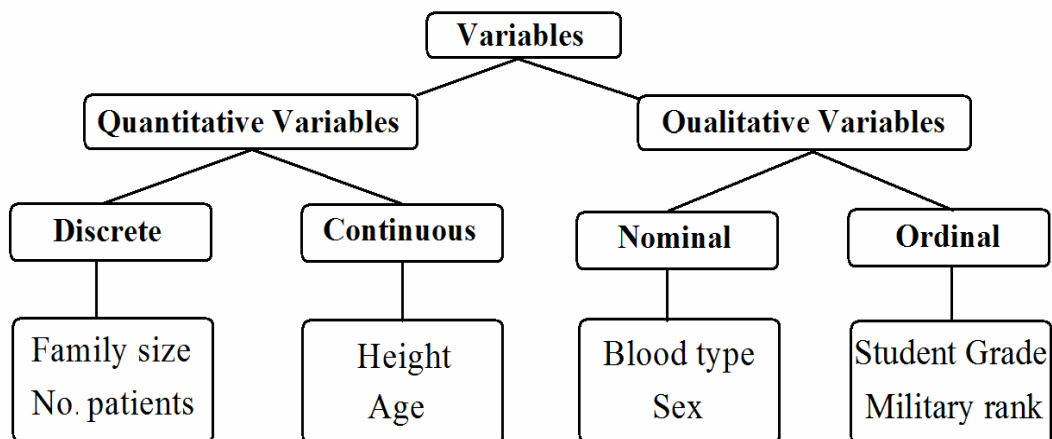
A nominal variable classifies the observations into various mutually exclusive and collectively non-ranked categories.

- Examples:-**
- 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**

- Examples:**
- Blood pressure level (high- normal- low)
 - Educational level (elementary, intermediate, ...)
 - Students grade (A, B, C, D, F)
 - Military rank



1.4 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.



Note: Explanation of (level) in the variables

- blood pressure level: *ordinal qualitative variable*
- blood sugar level($10 < x < 15$): *continuous quantitative variable*
- hemoglobin level (g\dl): *continuous quantitative variable*

CHAPTER 2: Strategies for Understanding the Meaning of Data:

2.1 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

2.2 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 smoking cessation:

55	46	58	54	52	69	40	65	53	58
----	----	----	----	----	----	----	----	----	----

The ordered array is:

40	46	52	53	54	55	58	58	65	69
----	----	----	----	----	----	----	----	----	----

2.3 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	13.5	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	18.3	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

Class Interval	Tally	Frequency
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

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	$n=50$

what is the variable?
 what is the type of variable ?

what is the sample size ?
 you will probably face missing value in the Frq ? if i remove 10
 how can i find missing value ?

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 = gap between class intervals
- d = lower limit – upper limit of the preceding class interval
- true upper limit = upper limit + $d/2$
- true lower limit = lower limit – $d/2$

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

$$(12.95+13.95)/2=13.45$$

For example: Mid-point of the 1st interval = $(13.0+13.9)/2 = 13.45$

Mid-point of the last interval = $(18.0+18.9)/2 = 18.45$

Note:

(1) Mid-point of a class interval is considered as a typical (approximated) value for all values in that class interval.

For example: approximately we may say that:

there are 3 observations with the value of 13.45

there are 5 observations with the value of 14.45

⋮

there are 1 observation with the value of 18.45

(2) There are no gaps between true class intervals. The end-point (true upper limit) of each true class interval equals to the start-point (true lower limit) of the following true class interval.

Cumulative frequency:

Cumulative frequency of the 1st 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 = $\text{Relative frequency} \times 100\%$

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%

$n=50$

total = 1

total = 100%

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

The number of people whose hemoglobin levels are less than or equal to 17.9 = 49

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 less than or equal to 16.9 = 78%

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

The percentage of people whose hemoglobin levels are more than or equal 16 = 54%

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-points
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	$n=100$		

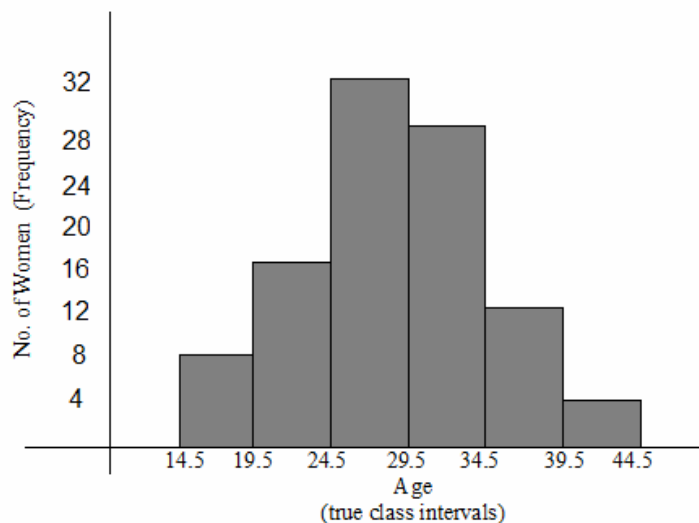
Width of the interval:

True Class Interval >> $W = \text{true upper limit} - \text{true lower limit} = 19.5 - 14.5 = 5$

Mid-points & Class interval & True Class $W = \text{lower limit} - \text{lower limit of the preceding interval}$

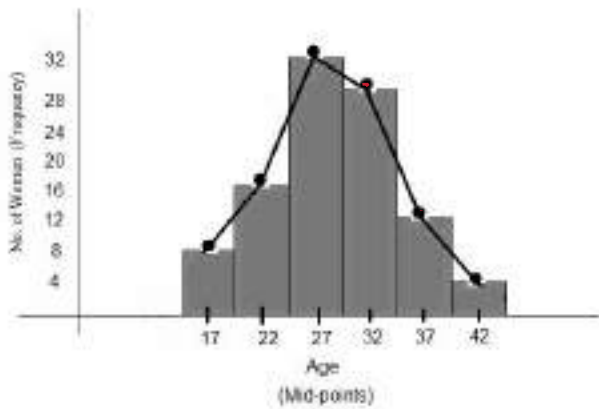
(1) Histogram: Organizing and Displaying Data using Histogram:

frequency
or
percentage
frequency
or
relative
frequency

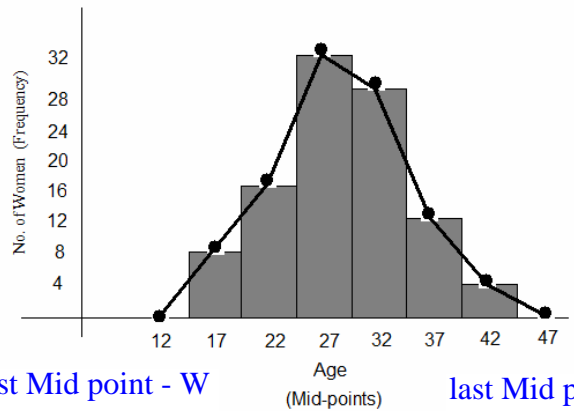


(2) Frequency Polygon: Organizing and Displaying Data using Polygon:

Polygon (Open)

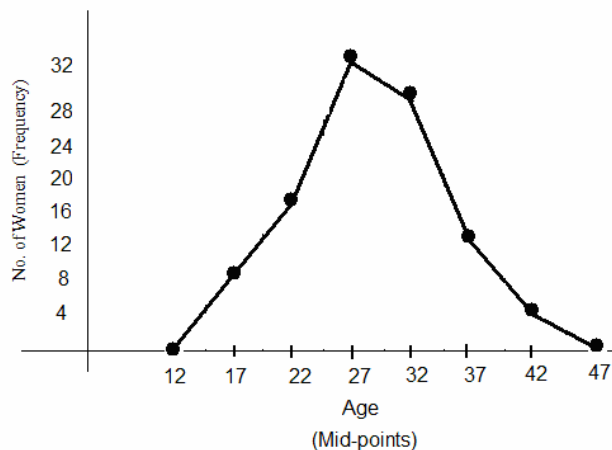


Polygon (Closed)



first Mid point - W
 $17 - 5 = 12$

last Mid point + W
 $42 + 5 = 47$



To calculate width :

First method:

width = $W = \text{lower limit} - \text{lower limit of the previous interval}$

width = $W = \text{Mid point} - \text{Mid point of the previous interval}$.

Second method:

$W = \text{True upper limit} - \text{True lower limit}$

2.4 Descriptive Statistics: Measures of Central Tendency:

- 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:

- (1) "N" population values.
- (2) "n" sample of values.

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

The population size = N

- Let x_1, x_2, \dots, x_n be the sample values (these values are known).

The sample size = n .

- (i) 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.

- (ii) A **statistic** is a measure (or number) obtained from the sample values: x_1, x_2, \dots, x_n .

- Values of statistics are known in general.
- Since parameters are unknown, statistics are used to approximate (estimate) parameters.

Measures of Central Tendency (or measures of location):

The most 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 considered to be a typical (or a representative) value of the set of data as a whole.

Mean:

(1) The Population mean (μ):

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 μ is a **parameter** (it is usually unknown and we are interested to know its value)

(2) 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 μ .

Example: Suppose that we have a population of 5 population values:

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

Suppose that we randomly select a sample of size :

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

The population mean is: $\mu = \frac{41+30+35+22+27}{5} = \frac{155}{5} = 31$ (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 takes into account all values of the data.

Disadvantages:

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

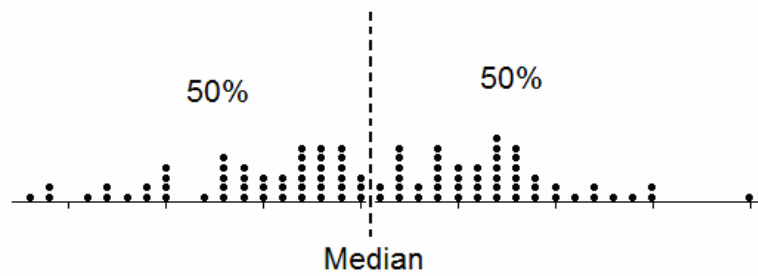
For example:

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:

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

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

$$\text{Median} = y_m.$$

Ordered set → (smallest to largest)	y_1	y_2	...	y_m middle value	...	y_n
Rank (or order) →	1	2	...	m	...	<i>n</i>

- 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 ($m+1$). The median is the mean (average) of the two middle values of the **ordered** observations:

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

Ordered set	→	y_1	y_2	...	y_m middle value	y_{m+1} middle value	...	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.

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

Ordered set	→	10	21	38 (middle value)	53	54
Rank (or order)	→	1	2	3 (= m)	4	5

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.

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

Therefore, the ranks of the middle values are:

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

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.

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

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

Advantages and disadvantages 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:

- The median does not take into account all values of the sample.
- 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**

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 and disadvantages 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	Mode
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:

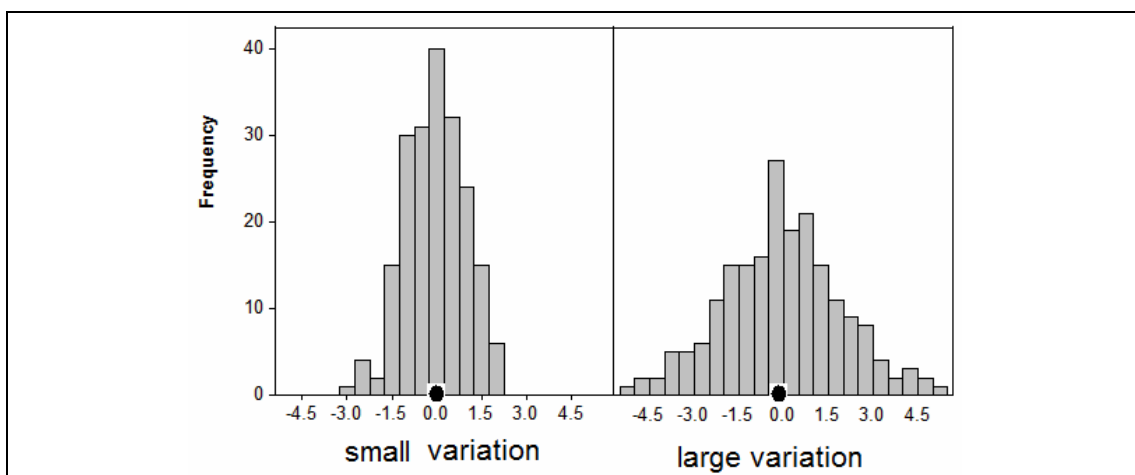
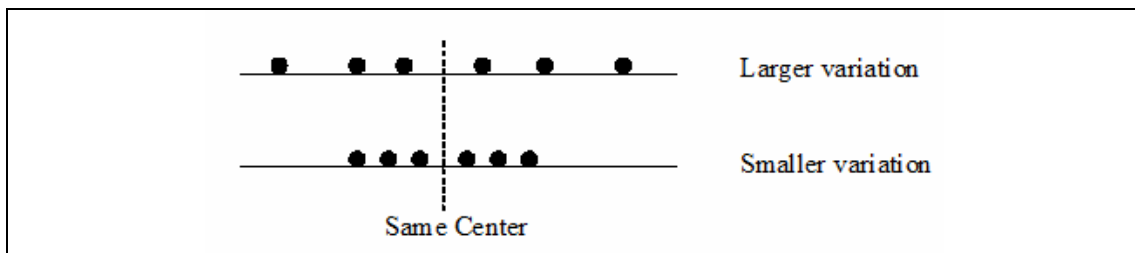
- The mode is not a “good” measure of location, because it depends on a few values of the data.
- The mode does not take into account all values of the sample.
- There might be no mode for a data set.
- There might be more than one mode for a data set.

2.6 Descriptive Statistics: 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).

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

Example:

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

Solution:

$$.max = 35$$

$$.min = 25$$

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

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 takes into account two of the values; however, it plays a significant role in many applications.

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 x_i from the sample mean \bar{x} is:

$$x_i - \bar{x}$$

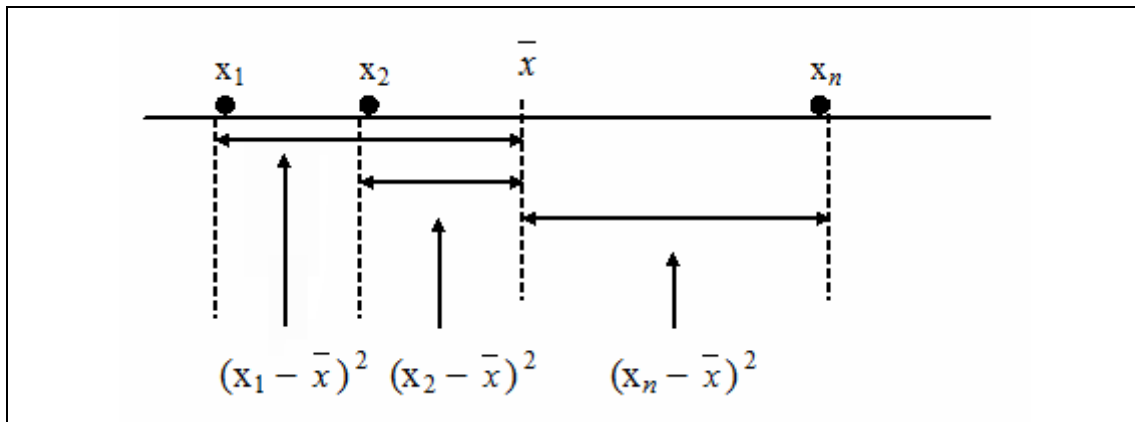
The squared deviation is:

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

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:



(1) The Population Variance σ^2 :

(Variance computed from the population)

Let X_1, X_2, \dots, X_N be the population values. The population variance (σ^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:
- σ^2 is a parameter because it is obtained from the population values (it is **unknown in general**).
 - $\sigma^2 \geq 0$

(2) 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:

$$\begin{aligned}
 S^2 &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \\
 &= \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n-1} \quad (\text{unit})^2 \\
 &= \frac{\sum_{i=1}^n x_i^2 - nx\bar{x}}{n-1} \quad \text{(practical formula)}
 \end{aligned}$$

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) σ^2 .
 - $S^2 \geq 0$
 - $S^2 = 0 \iff$ all observation have the same value
 \iff there is no dispersion (no variation)

Example:

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

Solution: $n = 5$

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

$$S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} = \frac{\sum_{i=1}^5 (x_i - 34.2)^2}{5-1}$$

$$\begin{aligned}
 S^2 &= \frac{(10 - 34.2)^2 + (21 - 34.2)^2 + (33 - 34.2)^2 + (53 - 34.2)^2 + (54 - 34.2)^2}{4} \\
 &= \frac{1506.8}{4} = 376.7 \quad (\text{unit})^2
 \end{aligned}$$

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

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

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

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: $S = \sqrt{S^2}$ (unit)

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

Example:

For the previous example, the sample standard deviation is

$$S = \sqrt{S^2} = \sqrt{376.7} = 19.41 \quad \text{(unit)}$$

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:

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

- 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 st data set	\bar{x}_1	S_1	$C.V_1 = \frac{S_1}{\bar{x}_1} 100\%$
2 nd 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:

1st data set: $\bar{x}_1 = 66$ kg, $S_1 = 4.5$ kg
 $\Rightarrow C.V_1 = \frac{4.5}{66} * 100\% = 6.8\%$

2nd data set: $\bar{x}_2 = 36$ kg, $S_2 = 4.5$ kg
 $\Rightarrow C.V_2 = \frac{4.5}{36} * 100\% = 12.5\%$

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

Descriptive statistics Measures of



Central tendency or (location)

- Mean (unit)
- Median (unit)
- Mode (unit)

Dispersion or (Variation)

- Range (unit)
- Variance = (standard deviation)² (unit)²
- Standard deviation = $\sqrt{\text{variance}}$ (unit)
- Coefficient of variation C.V. = $\frac{S}{\bar{x}} \times 100$ (unit-less)

Population:

X_1, X_2, \dots, X_N

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

Sample:

x_1, x_2, \dots, x_n

A descriptive measure computed from the values of a **sample** is called a "statistics"

	Population	Sample
Size	N	n
Mean	$\mu = \frac{\sum_{i=1}^N X_i}{N}$	$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$
Variance	$\sigma^2 = \frac{\sum_{i=1}^N (X_i - \mu)^2}{N}, \quad \sigma^2 \geq 0$	$S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}, \quad S^2 \geq 0$
Standard deviation	$\sigma = \sqrt{\sigma^2}, \quad \sigma \geq 0$	$S = \sqrt{S^2}, \quad S \geq 0$

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 Ω . called "Omega"

An Event:

Any subset of the sample space Ω is called an event.

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

Example: 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:

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

$E_2 =$ 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 = $\{\} = \phi \subseteq \Omega$

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

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

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

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

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

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

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

The outcomes are equally likely.

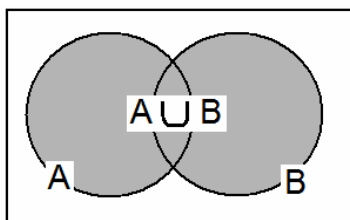
$$\therefore P(E_1) = \frac{3}{6}, \quad P(E_2) = \frac{3}{6}, \quad P(E_3) = \frac{2}{6}$$

Some Operations on Events:

Let A and B be two events defined on the sample space Ω .

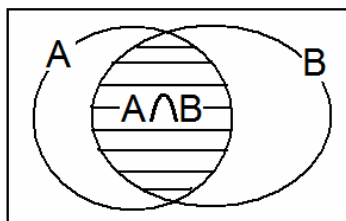
Union of Two events: $(A \cup B)$ or $(A + B)$

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.



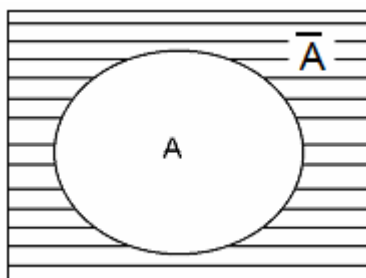
Intersection of Two Events: $(A \cap B)$

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 even A is denoted by \bar{A} . The even \bar{A} consists of all outcomes of Ω but are not in A . The even \bar{A} occurs if A does not.



Example:

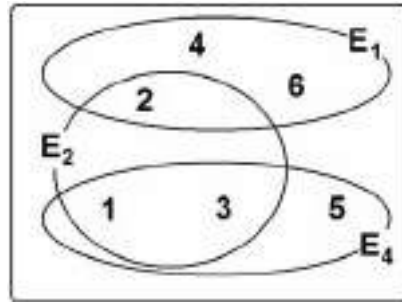
Experiment: Selecting a ball from a box containing 6 balls numbered 1, 2, 3, 4, 5, and 6 randomly.

Define the following events:

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

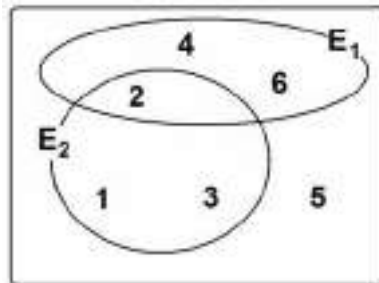
$E_2 = \{1, 2, 3\} =$ getting a number < 4 .

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



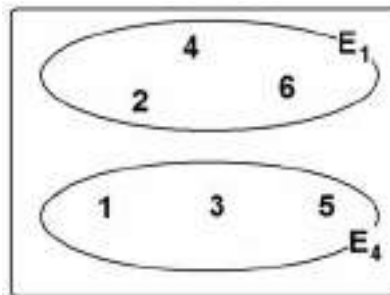
- (1) $E_1 \cup E_2 = \{1, 2, 3, 4, 6\}$
= getting an even number **or** a number less than 4.

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



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

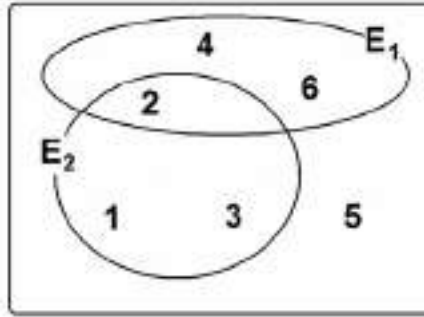
$$P(E_1 \cup E_4) = \frac{n(E_1 \cup E_4)}{n(\Omega)} = \frac{6}{6} = 1$$



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**.

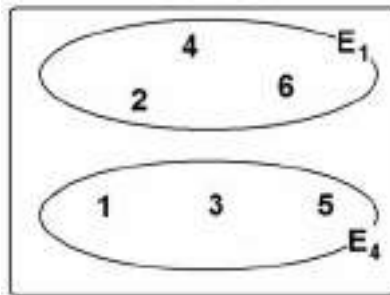
- (3) $E_1 \cap E_2 = \{2\} =$ getting an even number **and** a number less than 4.

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



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

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



Note: $E_1 \cap E_4 = \phi$. E_1 and E_4 are called disjoint (or mutually exclusive) events. These kinds of events can not occurred simultaneously (together in the same time).

(5) The complement of E_1

$$\begin{aligned} \bar{E}_1 &= \text{not getting an even number} = \overline{\{2, 4, 6\}} = \{1, 3, 5\} \\ &= \text{getting an odd number.} \\ &= E_4 \end{aligned}$$

Mutually exclusive (disjoint) Events:

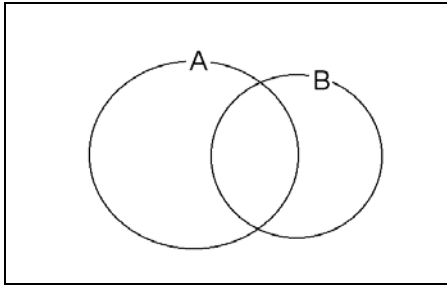
The events A and B are disjoint (or mutually exclusive) if:

$$A \cap B = \phi.$$

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

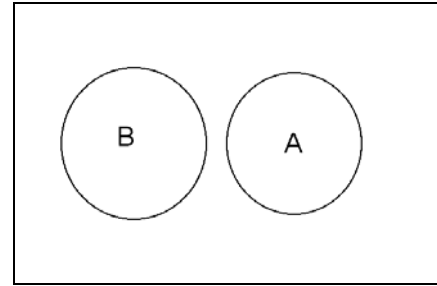
- (i) $P(A \cap B) = 0$
- (ii) $P(A \cup B) = P(A) + P(B)$

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



$$A \cap B \neq \phi$$

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



$$A \cap B = \phi$$

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

Exhaustive Events:

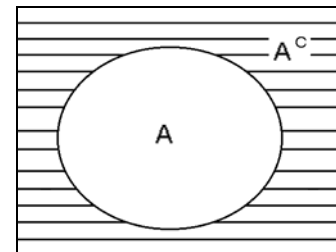
The events A_1, A_2, \dots, A_n are exhaustive events if:

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

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. $A \cap \bar{A} = \phi$ (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. $P(\phi) = 0$
4. $P(\bar{A}) = 1 - P(A)$

The Addition Rule:

For any two events A and B :

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

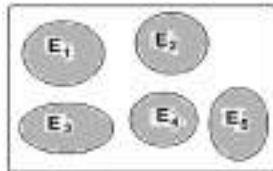
Special Cases:

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:

$$P(A_1 \cup A_2 \cup \dots \cup A_n) = P(A_1) + P(A_2) + \dots + P(A_n) = P(\Omega) = 1$$

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	$n(A_2 \cap B_1)$	$n(A_2 \cap B_2)$...	$n(A_2 \cap B_n)$	$n(A_2)$
.
.
.
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)

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:

$$P(A_i) = P(A_i \cap B_1) + P(A_i \cap B_2) + \dots + P(A_i \cap B_n) = \sum_{j=1}^n P(A_i \cap B_j)$$

For example,

$$P(A_2) = P(A_2 \cap B_1) + P(A_2 \cap B_2) + \dots + P(A_2 \cap B_n) = \sum_{j=1}^n P(A_2 \cap B_j)$$

We define the marginal probability $P(B_j)$ 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	$n(A_1)$
A_2	20	70	10	100	$n(A_2)$
A_3	30	100	120	250	$n(A_3)$
Total	100	200	200	500	

Table of **probabilities** of each event:

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



For example:
$$P(A_2) = P(A_2 \cap B_1) + P(A_2 \cap B_2) + P(A_2 \cap B_n)$$

$$= 0.04 + 0.14 + 0.02$$

$$= 0.2$$

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

Probabilities of the events:

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

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

Some operations on the events:

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

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

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

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

since $E_2 \cap E_4 = \phi$

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

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

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

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

another solution: $P(E_1^c) = 1 - P(E_1) = 1 - 0.4508 = 0.5492$

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

		Smoking Habit			Total
		Daily (B_1)	Occasionally (B_2)	Not at all (B_3)	
Age	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:

- A_3 = the selected physician is aged 40 - 49

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

Example: (Page 39)

Find: (سؤال اضافي)

$$\begin{aligned}P(A_3^C \cap B_2) &= P(A_1 \cap B_2) + P(A_2 \cap B_2) + P(A_4 \cap B_2) \\&= 9/339 + 30/339 + 0/339 \\&= 39/339 \\&= 0.11504\end{aligned}$$

$$\begin{aligned}P(A_2 \cap B_1^C) &= P(A_2 \cap B_2) + P(A_2 \cap B_3) \\&= 30/339 + 49/339 \\&= 79/339 \\&= 0.2330\end{aligned}$$

- 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.

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

- $A_3 \cup B_2$ = 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.
= $A_1 \cup A_2 \cup A_3$

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

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

Since $A_2 \cap A_3 = \phi$

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

or

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

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 woman randomly from this population.

Define the events:

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

The complement events:

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

~~Percentages: $\%(\bar{D}) = 10\%$ $\%(M) = 25\%$ $\%(D \cap M) = 5\%$~~

The probabilities of the given events are:

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

Q: Complete the table, then answer:

A Two-way table: (Percentages given by a two-way table):

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

Calculating probabilities of some events:

$D \cup M$ = the selected woman delivered prematurely or used medication.

$$P(D \cup M) = P(D) + P(M) - P(D \cap M) = 0.1 + 0.25 - 0.05 = 0.3$$

\bar{M} = The selected woman did not use medication

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

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

\bar{D} = The selected woman did not deliver prematurely

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

$$P(\bar{D}) = \frac{90}{100} = 0.90 \quad (\text{from the table})$$

$\bar{D} \cap \bar{M}$ = the selected woman 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 woman did not deliver prematurely and used medication.

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

$D \cap \bar{M}$ = the selected woman delivered prematurely and did not use medication.

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

$D \cup \bar{M}$ = the selected woman delivered prematurely or did not use medication.

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

$\bar{D} \cup M$ = the selected woman did not deliver prematurely or used medication.

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

$\bar{D} \cup \bar{M}$ = the selected woman did not deliver prematurely or did not use medication.

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

Conditional Probability:

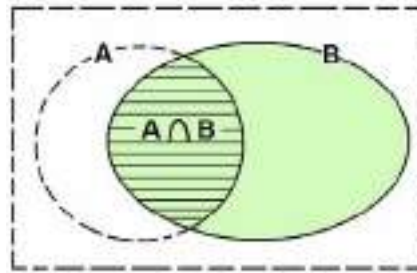
- The conditional probability of the event A when we know that the event B has already occurred is defined by:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad ; P(B) \neq 0$$

Known
given



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



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

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

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

(iii) 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:

		Smoking Habit			
		Daily (B_1)	Occasionally (B_2)	Not at all (B_3)	Total
Age	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

- Consider the event $(B_1 | A_2)$ knowing
 $(B_1 | A_2)$ = the selected physician smokes daily **given** that his age is between 30 and 39

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

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

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

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

Notice that:

$$P(B_1) = 0.519$$

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

$$P(B_1 | A_2) > P(B_1) , \quad P(B_1) \neq P(B_1 | A_2) \quad !$$

What does this mean?

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

Example: (Multiplication Rule of Probability)

Example : (Page 44)

If we have $P(A) = 0.9$, $P(B/A) = 0.8$ find $P(A \cap B) = ?$

Solution :

$$P(B/A) = \frac{P(A \cap B)}{P(A)}$$

$$0.8 = \frac{P(A \cap B)}{0.9}$$

$$\underline{P(A \cap B)} = 0.8 \times 0.9$$

$$= 0.72$$

~~B = the event of passing the second part~~

~~$A \cap B$ = the event of passing the first part and the second Part
= the event of passing both parts
= the event of passing the program~~

~~Therefore, the probability of passing the program is $P(A \cap B)$.~~

~~From the given information:~~

~~The probability of passing the first part is:~~

$$~~P(A) = 0.9 \quad \left(\frac{90\%}{100\%} = 0.9 \right)~~$$

~~The probability of passing the second part given that the trainee has already passed the first part is:~~

$$~~P(B|A) = 0.8 \quad \left(\frac{80\%}{100\%} = 0.8 \right)~~$$

~~Now, we use the multiplication rule to find $P(A \cap B)$ as follows:~~

$$~~P(A \cap B) = P(A) P(B|A) = (0.9)(0.8) = 0.72~~$$

~~We can conclude that 72% of the trainees pass the program.~~

Independent Events There are 3 cases:

- ~~• $P(A|B) > P(A)$ (knowing B increases the probability of occurrence of A)~~
- ~~• $P(A|B) < P(A)$ (knowing B decreases the probability of occurrence of A)~~
- ~~• $P(A|B) = P(A)$ (knowing B has no effect on the probability of occurrence of A)~~

~~In this case A is independent of B .~~

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

$$(i) P(A|B) = P(A) \Leftrightarrow (ii) P(B|A) = P(B) \Leftrightarrow (iii) P(B \cap A) = P(A)P(B)$$

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

Example: Suppose that A and B are two events such that:

$$P(A) = 0.5, P(B) = 0.6, P(A \cap B) = 0.2.$$

These two events are not independent (they are dependent)

because: $P(A)P(B) = 0.5 \times 0.6 = 0.3 \neq P(A \cap B) = 0.2.$

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

$$\text{Also, } P(A) = 0.5 \neq P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{0.2}{0.6} = 0.3333.$$

$$\text{Also, } P(B) = 0.6 \neq P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{0.2}{0.5} = 0.4.$$

For this example, we may calculate probabilities of all events.

We can use a two-way table of the probabilities as follows:

Complete the table,
then answer:

	B	\bar{B}	Total
A	0.2	?	0.5
\bar{A}	?	?	?
Total	0.6	?	1.00

We complete the table:

	B	\bar{B}	Total
A	0.2	0.3	0.5
\bar{A}	0.4	0.1	0.5
Total	0.6	0.4	1.00

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

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

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

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

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

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

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

$$P(\bar{A} \cup B) = \text{exercise}$$

$$P(\bar{A} \cup \bar{B}) = \text{exercise}$$

Add some questions

Q1: Are A and B independent events?

Q2: Are A and B disjoint events?

Q3: Are A and B Exhaustive events?

Note: The Addition Rule for Independent Events: If the events A and B are independent, then by the addition rule,

$$\begin{aligned} P(A \cup B) &= P(A) + P(B) - P(A \cap B) \\ &= P(A) + P(B) - P(A)P(B) \end{aligned}$$

Example : (Page 46)

Q1: Are A and B independent events?

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

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

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

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

$$0.2 \neq 0.3$$

So A,B are not independent.

Q2: Are A and B disjoint events?

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

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

So A,B are not disjoint.

Q3: Are A and B Exhaustive events?

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

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

So A,B are not Exhaustive.

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.

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

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

$P(\text{The chosen nurse has children and works night shift})$

$$= P(C \cap N) = \frac{n(C \cap N)}{n(\Omega)} = \frac{2}{12} = 0.16667$$

P(The chosen nurse has children and does not work night shift)

$$= P(C \cap \bar{N}) = \frac{n(C \cap \bar{N})}{n(\Omega)} = \frac{1}{12} = 0.0833$$

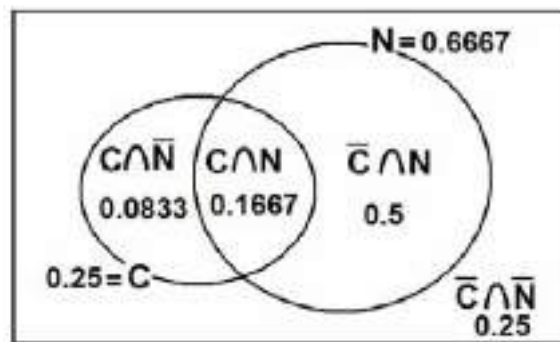
b) The probability of choosing a nurse who works at night given that she has children:

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

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

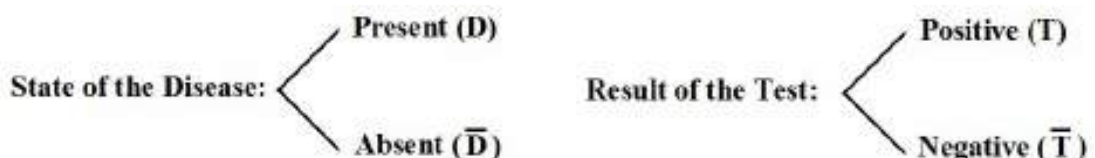
d) The events C and N are not disjoint because $C \cap N \neq \emptyset$. (Note: $n(C \cap N) = 2$)

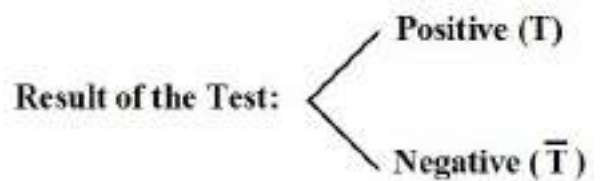
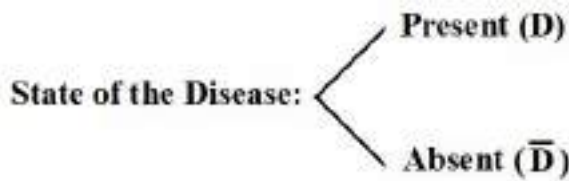
e) Venn diagram



3.5 Bayes' Theorem, Screening Tests, Sensitivity, Specificity, and Predictive Value Positive and Negative: (pp.79-83)

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





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 4 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. Its probability is:

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

2. **A false negative result:** This result happens when a test indicates a negative status when the true status is positive. Its probability is:

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

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:

Test Result	Disease		Total
	Present (D)	Absent (\bar{D})	
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 (a) subjects who have the disease and whose screening test result was positive.

From this table we may compute the following conditional probabilities:

- The probability of false positive result: $P(T | \bar{D}) = \frac{n(T \cap \bar{D})}{n(\bar{D})} = \frac{b}{b + d}$
- The probability of false negative result: $P(\bar{T} | D) = \frac{n(\bar{T} \cap D)}{n(D)} = \frac{c}{a + c}$
- The sensitivity of the screening test: $P(T | D) = \frac{n(T \cap D)}{n(D)} = \frac{a}{a + c}$
- The specificity of the 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 of:

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, $P(D | T)$:

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

Proof:

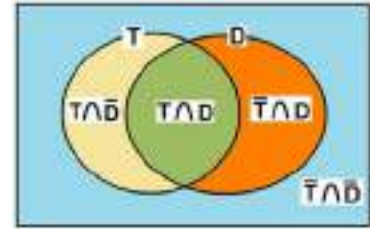
But we know that:

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

$$P(T \cap D) = P(T|D)P(D) \quad (\text{multiplication rule})$$

$$P(T \cap \bar{D}) = P(T|\bar{D})P(\bar{D}) \quad (\text{multiplication rule})$$

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



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

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

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 Bayes' theorem:

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

Note:

$P(\bar{T}|\bar{D})$ = 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 results 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 results into Equation (1), we get:

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

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:

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

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

$$P(\bar{T} | D) = \frac{n(\bar{T} \cap D)}{n(D)} = \frac{14}{450} = 0.0311 = 1 - \text{Sensitivity} = 1 - 0.9689$$

Substitution in Equation (2) gives:

$$\begin{aligned} P(\bar{D} | \bar{T}) &= \frac{P(\bar{T} | \bar{D}) P(\bar{D})}{P(\bar{T} | \bar{D}) P(\bar{D}) + P(\bar{T} | D) P(D)} \\ &= \frac{(0.99)(0.887)}{(0.99)(0.887) + (0.0311)(0.113)} \\ &= 0.996 \end{aligned}$$

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

Bayes Theorem pages 48-52

		Has the disease (D)	Dose not have the disease (\bar{D})	Total
The result of the test	Positive (T)	<p>Correct decision $n(T \cap D)$</p> <p>Sensitivity $P(T D) = \frac{n(T \cap D)}{n(D)}$</p>	<p>False decision $n(T \cap \bar{D})$</p> <p>false positive result $P(T \bar{D}) = \frac{n(T \cap \bar{D})}{n(\bar{D})}$</p>	n(T)
	Negative (\bar{T})	<p>False decision $n(\bar{T} \cap D)$</p> <p>false negative result $P(\bar{T} D) = \frac{n(\bar{T} \cap D)}{n(D)}$</p>	<p>Correct decision $n(\bar{T} \cap \bar{D})$</p> <p>Specificity $P(\bar{T} \bar{D}) = \frac{n(\bar{T} \cap \bar{D})}{n(\bar{D})}$</p>	n(\bar{T})
Total		n(D)	n(\bar{D})	n(Ω)

Note that from the table:

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

i.e. false negative + Sensitivity = 1 and Specificity + false positive = 1

The probability of the relevant disease in the general population, $P(D)$ [or $P(D') = 1 - P(D)$] which is obtained from another independent study.

Predictive value Positive:

$$P(D|T) = \frac{P(T|D) * P(D)}{\text{نفس البسط } (D \rightarrow \bar{D}) + \text{نفس البسط}}$$

$$= \frac{P(T|D) * P(D)}{P(T|D) * P(D) + P(T|\bar{D}) * P(\bar{D})} = \frac{\text{Sensitivity} * P(D)}{\text{Sensitivity} * P(D) + (1 - \text{Specificity}) * P(\bar{D})}$$

Predictive value Negative:

$$P(\bar{D}|\bar{T}) = \frac{P(\bar{T}|\bar{D}) * P(\bar{D})}{\text{نفس البسط } (\bar{D} \rightarrow D) + \text{نفس البسط}}$$

$$= \frac{P(\bar{T}|\bar{D}) * P(\bar{D})}{P(\bar{T}|\bar{D}) * P(\bar{D}) + P(\bar{T}|D) * P(D)} = \frac{\text{Specificity} * P(\bar{D})}{\text{Specificity} * P(\bar{D}) + (1 - \text{Sensitivity}) * P(D)}$$