

# ملاخص STAT 101

هذا الملاحض لا يعني

نحن الكتاب او التجميعات

هو فقط ترتيب وربط

للافكار لتسهيل الفهم

Rayhan

# CHAPTER 1

Descriptive statistics

1.0 - 1.1 - 1.2 - 1.3 - 1.4 - 1.5



**Data :-** a collection of information collected by means of experiments, observations, or real life events and stored in a proper format (data derived from dating)

**Statistics :-** a branch of science deals with collection, organization, presentation, analysis, interpretation of data and take the appropriate decisions. 1.0

## Types of statistics 1.1

### Descriptive

consists of methods and techniques which are used for presenting and summarizing data in tables or graph forms and provide some numerical measures for it.

\* Examples:-


The study of some human characteristics includes coverage height or weight based on data taken from a group of people is a descriptive study.


\* Forms of tables:-


- Frequency table


- Frequency distribution table

\* Forms of graphs:-

- Pie chart 

- Bar chart 

- Histogram 

- Polygon 

\* Forms of numerical measures:-

- averages

- Percentiles

- dispersions

and other measures:-

descriptive	inferential
<ul style="list-style-type: none"><li>- does not make conclusions</li><li>- organizes data and describes them</li></ul>	<ul style="list-style-type: none"><li>- does make conclusions</li><li>- Takes the data we organized and draws conclusions</li></ul>

### Inferential

some methods and techniques that can be used to make conclusions about the entire population using the observations from the samples taken from the population.

\* Methods based on

probability and distribution theory to make inferences:-

- point estimation

- interval estimation

- hypothesis testing

- statistical modeling

- clustering

- and many other methods

### Population :-

set of all things - have at least one common feature - that will be subjected to a study to obtain inferences for a specific problem.

\* the number of elements is called the size  $(N)$  or  $(n)$

\* The elements of a population is called an individuals

Examples:-

- Set of all student in a country

- Set of all palm trees in Ksar

### Samples :-

a subset of population which is used to collect information and to make inferences about the population.

\* the number of elements is called the size and denoted by  $(n)$  or  $(m)$

Examples:-

- the set of selecting Qasim trees from Ksar Palm trees.

\* since using the whole population to make inferences is hard, we use a subset of it (sample)

### Relationship between sample and population

Population

we take a portion of the population to study our hypothesis

Sample

based on the information we took from the sample we make inferences about the entire population



Parameter	Statistic
<ul style="list-style-type: none"> <li>- certain quantity or quality for describing a characteristic in a given <u>population</u> that <del>the</del> summarizes the data for the entire population</li> <li>- in a statistical study the measure is unknown</li> </ul>	<ul style="list-style-type: none"> <li>- certain quantity or quality for describing a characteristic in a given <u>sample</u> that summarizes the data for the entire sample</li> <li>- in a statistical study the measure is known</li> </ul>
one uses a statistic to make inferences about a parameter	

## # Variable :-

a ~~map~~ map or a function  $X$  defined on the population or sample and takes value in an arbitrary set  $M$   
 $X: \text{population or sample} \rightarrow M$

\* the variable measures a characteristic, feature, or factor in the population

\* not all individuals have the same ~~characteristics~~ exact all characteristics but they may share some of them

## Types of variables

### according to type

#### Qualitative

variables that takes non numerical ~~numbers~~ values or numerical values (can't perform mathematical operations on the values) which indicate an attribute or property.

- Examples:-
- marital status
  - eye color
  - ID number of Kisan Students

#### Quantitative

variables that takes numerical values and we can perform mathematical operations on them

\* can be ordered increasing or decreasing.

Examples:-

- Height
- weight
- temperature
- revenue

#### Discrete

variables takes finite or infinite countable numbers

Examples:-

- Quant Discrete:
  - \* number of accidents in a city
  - \* number of laptops sold a day
- Qual Discrete:
  - \* blood groups of people
  - \* types of cars

### according to numbers

#### continuous

variable takes uncountable number of values

Examples:-

- Quant continuous:
  - \* weight of person
  - \* distance between cities
- Qual continuous:
  - \* all spectrum colors resulting from white light



**Raw data :-** data that is not organized well.  
(ungrouped data)

1.2

Examples:- 1.2.1

Qualitative

June	July	January	December
March	March	April	September
August	May	May	February
July	February	June	June
April	February	November	August
July	April	June	January
December	the birth of 25 members of a community		

1.2.2

Discrete Quantitative

4	1	2	0	2	0	1	2
0	3	0	4	0	1	1	2
3	1	2	4	0	1	0	2
4	0	1	1	2	3	0	4
0	2	0	5	2	3	1	0
Number of children in 40 family in a society							

1.2.3

Continuous Quantitative

170	180	175	176	172
173	183	171	169	174
180	190	186	189	192
167	175	170	178	191
165	177	183	181	179
The height of 25 adult men in cm				

## Types of Organization of raw data

### 1- Frequency Table :- (Qual or discrete Quan)

class	frequency	Relative frequency	Percent frequency
all the classes of the variable are mentioned	The number of times the data occurs	$\text{The Relative Frequency (R.F)} = \frac{\text{The frequency of each}}{\text{The sum of frequencies}}$	$\text{Percent frequency} = R.F \times 100\%$
<b>sum</b>	$\Sigma \text{frequency}$	$\Sigma \text{Relative frequency} = 1$	$\Sigma \text{percent frequency} = 100\%$

Examples:-

- look at example 1.2.4 page 12 for qualitative data frequency table.

- look at example 1.2.6 page 13 for quantitative discrete frequency table.



# 2- frequency distribution table

(continuous Quan Data)

\* unwanted to use frequency dist table for data less than 32 ( $n \geq 32$ )

1- Range :  $x_L - x_S$

2- K 5-20 classes, if not provided:  $K = [3.322 \log n]$   
 $K = [x] = \text{integer number} \leq x$

3- class boundary length  $C = \frac{R + \text{one measuring unit}}{K}$

4- class limit length =  $C - 1$

5- Lower Limit = -0.5 from lower limit of class boundary | Higher Limit = +0.5 from higher limit of class boundary

6- class midpoint =  $\frac{\text{upper limit} + \text{lower limit}}{2}$

7- find the frequency

9- find ACF

8- find the Relative frequency

10- find DCF

Example :- look at data page 16

- 1- start with 1st frequency
- 2- add 1st and 2nd frequency
- 3- add 3rd frequency to the previous
- 1- start with SF
- 2- subtract 1st frequency
- 3- subtract 2nd frequency from previous DCF

class limit	class boundary	Midpoint	frequency	RF	ACF	DCF
1-5	0.5 → 5.5	3	7	0.14	<del>7</del> 7	50
6-10	5.5 → 10.5	8	9	0.18	7+9=16	50-7=43
11-15	10.5 → 15.5	13	14	0.28	7+9+14=30	50-7-9=34
16-20	15.5 → 20.5	18	12	0.24	7+9+14+12=42	50-7-9-14=20
21-25	20.5 → 25.5	23	8	0.16	7+9+14+12+8=50	50-7-9-14-12=8
total	-	-	50	1	-	-



# Types of graphs

1.3

## 1- Pie chart :-

a simple way of representing the portions of each class or category of data on a circular disk so that each category is allocated a circular sector representing it.

the measure of angle =  $RF_{ij} \times 360$



Example:-

data from 1.2.6 in the book

category (AB) the Angle is:  $0.05 \times 360 = 18$  degree

" (B) " " " :  $0.1 \times 360 = 36$  degree

" (A) " " " :  $0.45 \times 360 = 162$  degree

" (O) " " " :  $0.40 \times 360 = 144$  degree

Angles drawn are not accurate!

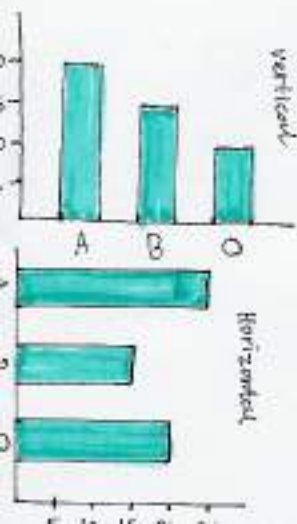
## 2- Bar chart (class + frequency)

a representation of data of discrete variable with finite values (quant or qual) this is done through vertical or horizontal bars so that each statement has a height equal to the frequency of it.

- \* the bar chart is useless if the values are large
- \* the width is uniform for all bars
- \* the width doesn't matter
- \* the bars must be separated

### Types of bar charts

#### 1- normal



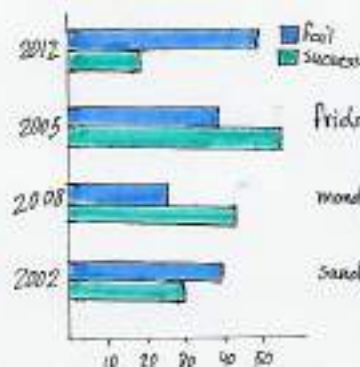
#### 2- Two direction bar chart

a bar chart where we can represent both positive and negative values. check pg 20 for details



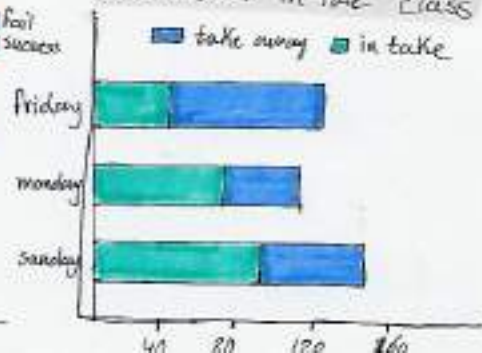
#### 3- Multiple bar chart

a bar chart where we can use it to represent multiple inter related values by clustering



#### 4- component or stacked bar chart

a bar chart where we can represent each component by a section in the bar and the size is proportional to its contribution in the class





# 3 - Histogram :-

## class boundary + frequency (continuous)

a graph in which class boundaries are marked on a vertical axis and is constructed by drawing a rectangular column above each actual category so that its height equals the frequency.

\* Similar to a bar chart but the bars must touch and the classes are close to each other

## Types of histograms

based on shape

### Symmetric

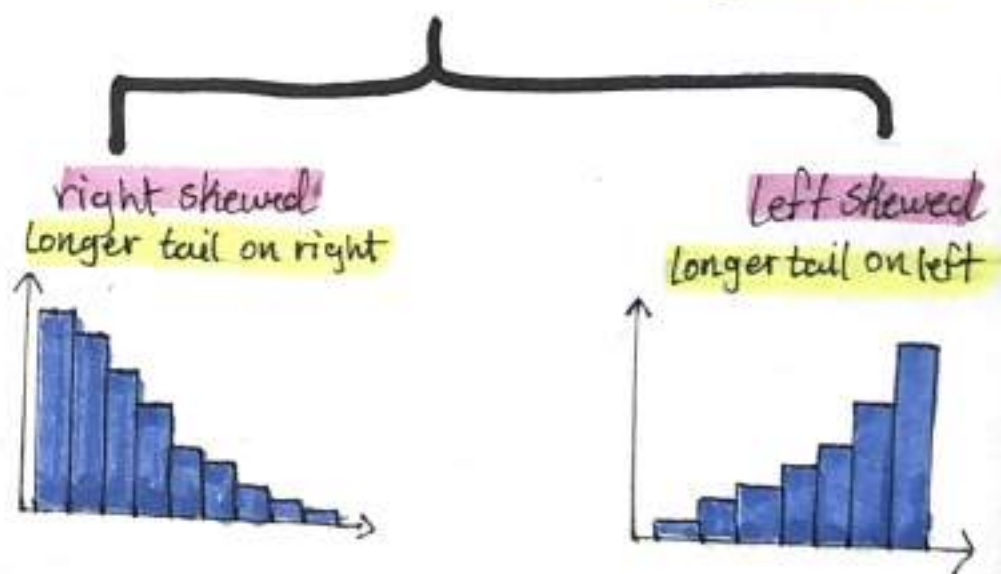
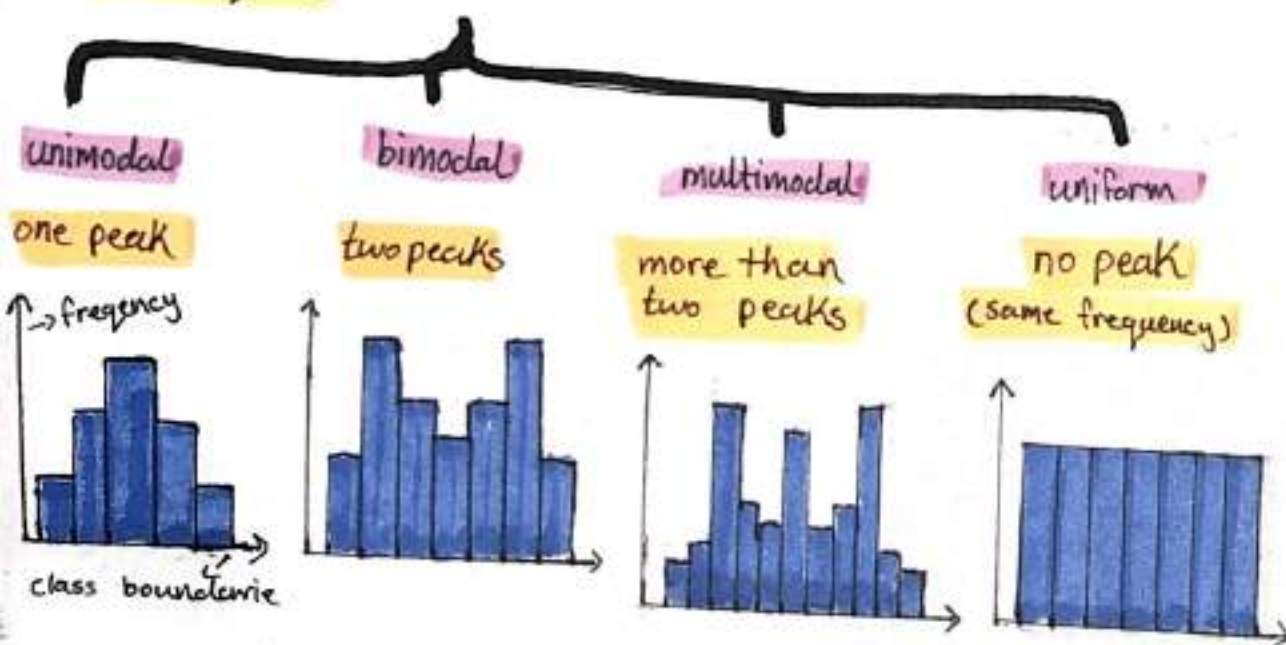
divided into 2 identical sections with respect to a column in the middle

- \* the highest frequency is in the middle then it decreases from both left and right boundary
- \* some have the bell shape
- \* the tails from left and right are equal

### (non-symmetric) Skewed

**Skewedness:-** bins on one side have ~~one~~ high frequency which ~~decrease~~ decrease as we move to the other side

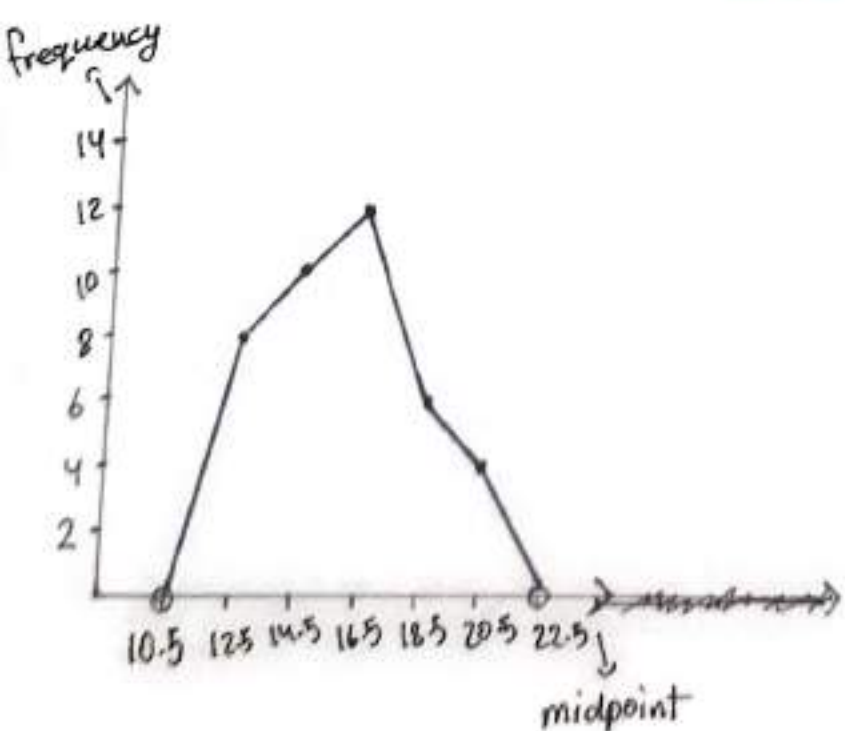
\* the side with the lower frequency is said to have a longer tail



# 4 - Polygon

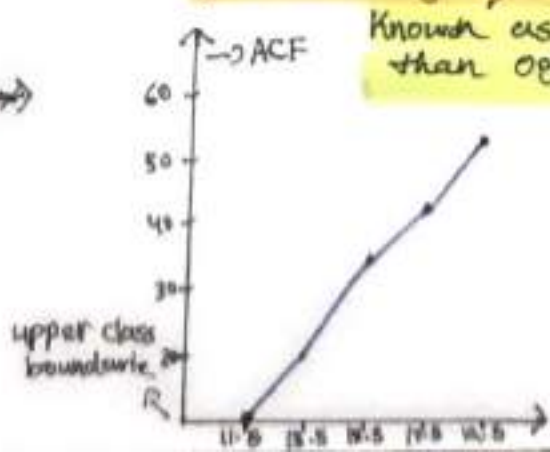
## (continuous)

class midpoint + frequency



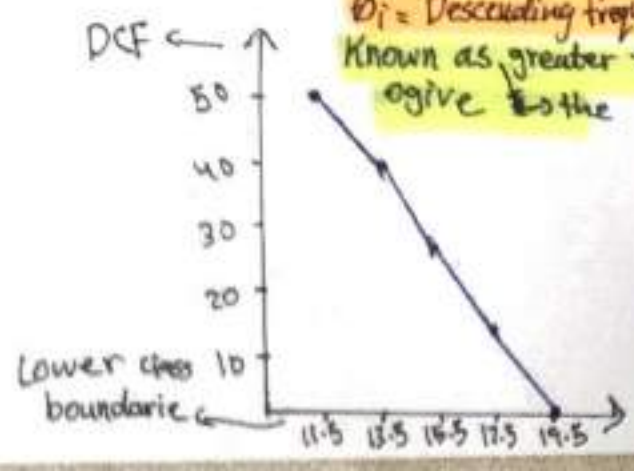
### Descending cumulative frequency table (DCFP)

Polygon which connects with a straight line the point  $(b_i, F_i)$   
 $b_i$  = upper bound of the class  
 $F_i$  = Ascending frequency of a class



### Ascending cumulative frequency table (ACFP)

Polygon which connects with a straight line  $(b_i, \phi_i)$  that being a point.  $b_i$  = lower bound of a class  
 $\phi_i$  = Descending frequency. Known as greater than ogive to the





1.4

# Measure of central Tendencies

## 1 - Mean $\bar{x}$ only for quantitative data

### Raw data Mean

$$\bar{x} = \frac{\sum x_i}{n}$$

Example:-

calculate the mean for this data set

20-18-15-15-14-4  
12-11-9-7-6-1?

$$\bar{x} = \frac{\sum x_i}{n}$$

$$= \frac{20+18+15+15+14+4+12+11+9+7+6+1}{12}$$

$$= \frac{132}{12} = 11$$

### frequency table Mean

$$\bar{x} = \frac{\sum f_i x_i}{\sum f_i} = \frac{\sum f_i x_i}{n}$$

Example:-

No. Subject	Frequency
0	8
1	18
2	12
3	2
-	40

$$\bar{x} = \frac{\sum f_i x_i}{\sum f_i}$$

$$= \frac{(0 \times 8) + (1 \times 18) + (2 \times 12) + (3 \times 2)}{40}$$

$$= \frac{48}{40} = 1.2$$

### frequency dist table Mean

$$\bar{x} = \frac{\sum f_i x_m}{\sum f_i}$$

Example:-

class boundary	class midpoint	frequency
11.5-13.5	12.5	8
13.5-15.5	14.5	10
15.5-17.5	16.5	12
17.5-19.5	18.5	6
19.5-21.5	20.5	4
-	-	40

$$\bar{x} = \frac{\sum f_i x_m}{\sum f_i}$$

$$= \frac{(12.5 \times 8) + (14.5 \times 10) + (16.5 \times 12) + (18.5 \times 6) + (20.5 \times 4)}{40}$$

$$= \frac{636}{40} = 15.9$$

### weighted Mean

$$\bar{x} = \frac{\sum w_i x_i}{\sum w_i}$$

when weights are equal then ordinary  $\bar{x} = \text{weight } \bar{x}$

Examples:-

course	grade	points	hours	$\sum w_i x_i$
math	B	4	4	16
stat	A	5	3	15
Eng	C	3	3	9
Physics	C	3	4	12
total	-	-	14	52

$$\bar{x} = \frac{\sum w_i x_i}{\sum w_i}$$

$$= \frac{52}{14}$$

$$= 3.71$$

## 2 - Median $\tilde{x}$ only for Quan data

divides the data in 2 halves after ordering them

### Raw data

$$\tilde{x} = \frac{x_{\frac{n}{2}} + x_{\frac{n}{2}+1}}{2}$$

n odd

Example:-

21-22-23-24-26-28-29  
 $x_1, x_2, x_3, x_4, x_5, x_6, x_7$

$$\tilde{x} = \frac{x_{\frac{n}{2}} + x_{\frac{n}{2}+1}}{2}$$

$$= \frac{24}{2} = 24$$

$$\tilde{x} = \frac{x_{\frac{n}{2}} + x_{\frac{n}{2}+1}}{2}$$

n even

Example:-

21-22-23-24-26-28-29-35  
 $x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8$

$$\tilde{x} = \frac{x_{\frac{n}{2}} + x_{\frac{n}{2}+1}}{2}$$

$$= \frac{x_4 + x_5}{2} = \frac{24+26}{2} = 25$$

### frequency table

$$\frac{\sum f_i + 1}{2} \leq x_i$$

 $\sum f_i$  odd

value  $\tilde{x}$  is the class  $x_i$  with the smallest cumulative frequency greater or equal to  $\frac{\sum f_i + 1}{2}$

Example:- see page 33

$$\sum f_i \text{ even}$$

cumulative frequency

$$\text{if } x_i \leq \frac{\sum f_i}{2} \text{ and } x_{i+1} < \frac{\sum f_i}{2}$$

$$\tilde{x} = \frac{x_i + x_{i+1}}{2}$$

cumulative frequency

### frequency dist table

1- Find median class (smallest cumulative frequency  $\geq \frac{\sum f_i}{2}$ )

$$2- \tilde{x} = \tilde{L} + \frac{\frac{\sum f_i}{2} - (\tilde{F} - \tilde{f})}{\tilde{f}} \times C$$

$\tilde{L}$  = Lower limit of median class

$\tilde{F}$  = cumulative frequency of median class

$\tilde{f}$  = frequency of median class

$C$  = class length of median class

Example:- see page 34



# 3-Mode $\hat{x}$

The highest frequency

## Quan and Qual data

- \* if the highest frequency is constant for all data then the data has no mode.
- \* multiple values have the highest frequency then the data has more than one mode.

### Raw data

Take the value with the highest frequency

Example:-

\* AA, C, A, D, A, B, B, C, D, A, B

The mode is A

\* 12, 11, 13, 14, 13, 12, 11, 12, 13, 12, 12, 13, 14, 13

The modes are 12, 13

\* Symmetric distribution  
 $\bar{x} = \hat{x} = \tilde{x}$



### Relationship between Mean, Median, Mode

\* Left skewed distribution  
 $\hat{x} > \tilde{x} > \bar{x}$



\* Right skewed distribution  
 $\bar{x} > \tilde{x} > \hat{x}$



### Frequency table

Take the value with the highest frequency

Example:-

Frequency	18	4	16	2
Blood group	A	B	O	AB

The Mode here is A

### Frequency dist table

Find the modal class (frequency is greater than all the others) not extremity

$$2 - \hat{x} = \hat{L} + \frac{d_1}{d_1 + d_2} \times C$$

$\hat{L}$  = Lower class limit of modal

$d_1 = f_{\text{modal}} - f_{\text{previous class directly}}$

$d_2 = f_{\text{modal}} - f_{\text{next class directly}}$

$C$  = class length

Example:-

Class boundary	50.5-55.5	55.5-60.5	60.5-65.5	65.5-70.5	70.5-75.5
Frequency	2	7	8	4	3

$$\hat{x} = \hat{L} + \frac{d_1}{d_1 + d_2} \times C = 60.5 + \frac{8-7}{(8-7)+(4-3)} \times 5 = 61.5$$

	Mean $\bar{x}$	Median $\tilde{x}$	Mode $\hat{x}$
Advantage	<ul style="list-style-type: none"> <li>quick and easy to compute</li> <li>all values are considered</li> <li>one value for a set of data</li> </ul>	<ul style="list-style-type: none"> <li>easy to compute and understand</li> <li>not affected by outliers or extreme values</li> <li>can be used if there is missing data not in the middle</li> </ul>	<ul style="list-style-type: none"> <li>quick and easy to compute</li> <li>for quan and qual</li> <li>not affected by extreme values</li> </ul>
Disadvantage	<ul style="list-style-type: none"> <li>not for Qual data</li> <li>highly affected by outliers since it considers all values</li> <li>not applicable if data is lost</li> </ul>	<ul style="list-style-type: none"> <li>doesn't take all values into account</li> <li>not used in many statistical studies</li> <li>can't be used for qual data</li> </ul>	<ul style="list-style-type: none"> <li>There may be two or more modes</li> <li>sometimes there is no mode at all</li> <li>may not reflect central tendencies accurately</li> </ul>



# dispersions

**1- Percentiles :-** ( $P_1, P_2, \dots, P_{99}$ ) divid the data on 100 ordered

$P_{50}$  = divid data to 2 halves so  $P_1 = 1\%$   $P_{99} = 99\%$   $P_{25} = Q_1$   $P_{50} = Q_2$   $P_{75} = Q_3$

**\* How to calculate?**

$$1 - P_r = \frac{r(n+1)}{100} = \frac{K}{S}, \dots$$

$$2 - P_r = x_k + S(x_{k+1} - x_k)$$

\* It's useful when the data is large ( $n \geq 99$ )

**Example :-**

10 20 36 36 36 38 39 40 40 50  
51 53 56 56 60 60 60 60 70 70  
70 70 72 80 88 88 92 92 92 95

calculate  $P_{35}$  for this data?

$$P_{35} = \frac{35(30+1)}{100} = 10.85 \quad \left| \quad P_{35} = x_{10} + 0.85(x_{11} - x_{10}) \right. \\ = 50 + 0.85(51 - 50) = 50.85$$

**2- Deciles :-** ( $D_1, D_2, \dots, D_9$ ) divid the data on 10  $D_5 = Q_2 = P_{50}$

**\* How to calculate?**

$$1 - d_r = \frac{r(n+1)}{10} = \frac{K}{S}, \dots$$

$$2 - D_r = x_k + S(x_{k+1} - x_k)$$

**Example :-**

calculate the decile 6th for the above data?

$$d_6 = \frac{6(30+1)}{10} = 18.6 \quad \left| \quad D_6 = x_{18} + 0.6(x_{19} - x_{18}) \right. \\ = 60 + 0.6(70 - 60) = 66$$

**3- Quartiles :-** ( $Q_1, Q_2, Q_3$ ) divid the data on 4

$Q_1 = P_{25}$   $Q_2 = P_{50} = D_5$   $Q_3 = P_{75}$

**Example :-**

calculate the  $Q_1$  for this data

$Q_2$   
 $Q_3$

21 22 23 24 26 28 29

$$* Q_1 = \frac{1(7+1)}{4} = 2.0 \quad \left| \quad Q_1 = x_2 + 0(x_3 - x_2) = x_2 = 22 \right.$$

$$* Q_2 = \frac{2(7+1)}{4} = 4.0 \quad \left| \quad Q_2 = x_4 + 0(x_5 - x_4) = x_4 = 24 \right.$$

$$* Q_3 = \frac{3(7+1)}{4} = 6.0 \quad \left| \quad Q_3 = x_6 + 0(x_7 - x_6) = x_6 = 28 \right.$$

\* used to determine whether a value  $x$  of a given data is an extreme values.

**Extreme values**

$$* LF = Q_1 - 1.5(Q_3 - Q_1)$$

$$* HF = Q_3 + 1.5(Q_3 - Q_1)$$

**Five numbers**

Smallest value,  $Q_1$ ,  $Q_2$ ,  $Q_3$   
and largest value



**Box plot:-** the graphical representation of the 5 numbers

## \* How to construct a box plot?

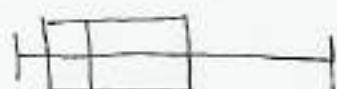
- 1- calculate the five numbers
  - 2- draw an axis
  - 3- draw vertical lines on all 3 Quantiles and join them to make a box
  - 4- calculate the end of the right whisker:
    - if data does not have a great value then it ends with highest value ( $x_n$ )
    - if data has great value then it ends at HF
  - 5 calculate the end of the left whisker:
    - if data does not have a small value then it ends with smallest value ( $x_1$ )
    - if data has a small value then it ends at LF
  - 6- extreme values are represented by ~~at~~ (\*) or (.)
- \* for continuous data we can use a box plot to determine if the distribution is symmetric or skewed



left skewed



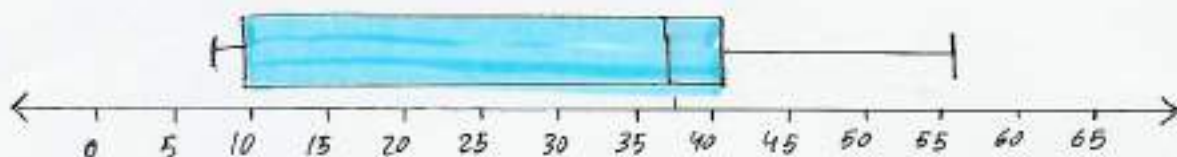
symmetric



Right skewed

Example:-

draw the box plot for these five numbers  
7, 9.75, 37.5, 40.75, 56



it does not have extreme values



# Measures of variation dispersions

1.5

## 1 - Variance $S^2$ $S^2 \geq 0$ always

### Raw data

- \* if  $n=1$  then  $S^2=0$
- \* expressed in square ( $m^2$ )

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

$$S^2 = \frac{\sum x_i^2 - \frac{(\sum x_i)^2}{n}}{n-1}$$

### Frequency table

- \* if  $\sum f_i = 1$  then  $S^2=0$
- \* expressed in square ( $m^2$ )

$$S^2 = \frac{\sum f_i (x_i - \bar{x})^2}{\sum f_i - 1}$$

### Frequency dist

- \* if  $\sum f_i = 1$  then  $S^2=0$
- \* expressed in square ( $m^2$ )

$$S^2 = \frac{\sum f_i (x_m - \bar{x})^2}{\sum f_i - 1}$$

## 2 - Standard deviation $S$ $S \geq 0$ always

$$S = +\sqrt{S^2}$$



\* Standard deviation is the best measure of dispersion

- \* we use measures of dispersion when the mean is used as a central tendency measure
- small value of standard deviation = variables close to mean
- large value of standard deviation = variables far from mean

- \* highly affected by extreme values
- \* not applicable if data is lost

## 3 - coefficient of variation $CV$

$$CV = \frac{S}{\bar{x}} \times 100\%$$

$$\bar{x} \neq 0$$

- \* a useful measure of variation to compare between sets of data with different units (measures)



## 4- Range R

Raw data

$$R = x_L - x_S$$

Frequency

$$R = x_m - x_d$$

grouped data  
(distribution table)

$$R = x_K - x_1$$

Interquartile

$$IQR = Q_3 - Q_1$$

\* it is also called as mid-spread because it gives us the middle 50% of the values

\* we can use the 5 numbers to find:-  
1- Ranges  
2- IQR  
3- Median

## 5- Z-SCORE (Standard Score) $Z_{x_i}$

\* converts data to make its mean = 0 and S = 1

$$Z_{x_i} = \frac{x_i - \bar{x}}{S}$$

\* we use the numerical value to assign a degree of data to tell us its position compare to the rest of the data

\*

مقياس لادرجة ان قيمة  $x_i$  عن المتوسط

$Z_{x_i} > 0$  يعني ان  $x_i$  اكبر من المتوسط

$Z_{x_i} < 0$  يعني ان  $x_i$  اصغر من المتوسط

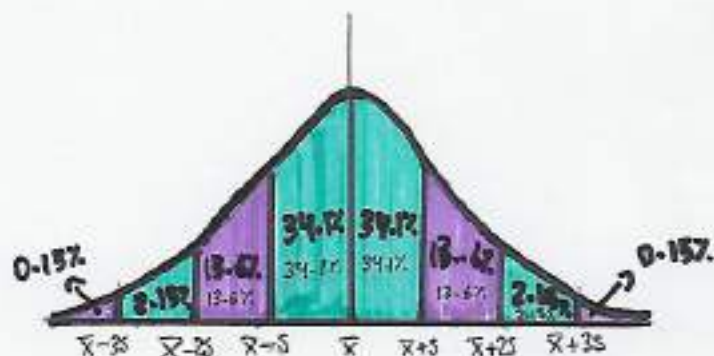
## The empirical Rule

\* if a data set has an approximate bell shape RF histogram

- 68.2% of data lies in  $\bar{x} \pm S$

- 95.4% of data lies in  $\bar{x} \pm 2S$

- 99.7% of data lies in  $\bar{x} \pm 3S$





# chapter 2

## Probability

2.1 - 2.2 - 2.3 - 2.4



# The fundamental principle of counting

2.1

has 2 Rules

## Multiplicative Rule

If we have  $A$  ways of doing something and have  $B$  ways of doing another then there are  $A \times B$  ways of doing both

\* "and" usually requires us to multiply

\* independent events

## Addition Rule

If we have  $A$  ways of doing something and have  $B$  ways of doing another and we can't do both at the same time then there are  $A+B$  to choose 1 action

\* "or" usually requires us to add

\* mutually exclusive events

for shortcuts we use factorial notation which is

$$n! = n \times (n-1) \times (n-2) \times (n-3) \times \dots \times 2 \times 1$$

$$1-0! = 1$$

$$2-1! = 1$$

some special cases:-

$$3 - n \geq 2 \rightarrow n! \approx n^n e^{-n} \sqrt{2\pi n}$$

Permutation	Combination
- Any arrangement of $r$ distinct objects, from a set of $n \in \mathbb{N}$ (natural numbers) different objects.	- Any unordered group of $r$ distinct objects from a set of $n$ different objects.
$nPr = \frac{n!}{(n-r)!}$ ; $0 \leq r \leq n$   $nPn = n!$	$nCr = \frac{n!}{r!(n-r)!} = \frac{nPr}{r!} = \binom{n}{r}$ ; $0 \leq r \leq n$
care about the order	don't care about the order
عناية بترتيب الأشياء في الترتيب	لا عناية بترتيب الأشياء unordered

## cardinal number of a set $|\Omega|$

Let  $\Omega$  be a given set, then the number of all elements in  $\Omega$  is  $|\Omega|$

\*  $\Omega$  is infinite and it's countable (elements  $\in \mathbb{N}$ ) then  $|\Omega| = \infty$ .

\*  $\Omega$  is infinite and it's not countable (elements  $\notin \mathbb{N}$ ) then  $|\Omega| = \mathfrak{P}$

( $\Omega$  has continuous capacity)



# Types of Experiments

2.2

## Regular (systematic)

we know its outcomes in advance

## Random (stochastically)

we don't know its outcomes in advance but we can determine the set of results

## Probability science :-

a branch of mathematics that deals with theoretical mathematical models of random experiments.

- \* each possible outcome is called an ~~out~~ elementary event
- \* 2 or more outcomes are said to have the same chance of appearing

## theoretical mathematical model (probability space) $[\Omega, \mathcal{F}, P]$

$\Omega$

The set of all possible results of a random experiments (space of elementary events)



### Discrete

Finite or countable  
Infinite events

### Continuous

uncountable number of outcomes

$\mathcal{F}$

(Algebra of events)  
a collection of events of interest with these condition :-

1-  $\Omega \in \mathcal{F}$  or  $\mathcal{F} \subseteq \Omega$

2- \* any elements A and B  $\in \mathcal{F}$

then  $A \cup B \in \mathcal{F}$

3- \* any element A  $\in \mathcal{F}$

then  $A^c \in \mathcal{F}$

\* if this condition is true :-

-  $A_1, A_2, A_3, \dots, A_n \in \mathcal{F}$  and  $A_1 \cup A_2 \cup A_3 \dots$

then  $\mathcal{F}$  is  $\sigma$ -algebra

\* elements of  $\mathcal{F}$  (algebra or  $\sigma$ -algebra) are called events

\*  $|\mathcal{F}| = 2^n$

P

(probability function)  
a measure of the likelihood a specific event will occur

RF

if the experiment A repeating N times and the event has n trials

$$P(A) = \frac{n(A)}{N}$$

\* Relative frequency is not a probability but they're approximate probability

\* But when repeated again and again they become actual probability

### classical

if all the events have the same chance in appearing (fair)

$$P(A) = \frac{|A|}{|\Omega|}$$

↑ compound

$$P(A) = \frac{1}{|\Omega|}$$

↑ simple

### theorem

if the experiment is not fair and is finite

$$P(A) = \frac{|A|}{|\Omega|}$$

## Types of events



### Simple

if it contains one elementary event (outcome)

### compound

if it contains two or more elementary events



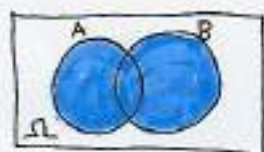
# Some operations on events 2.2-2.3-2.4

## 1- UNION of two events:-

$A \cup B$

an event containing all elementary events that belong to A or B or to both.

$$A \cup B = \{w: w \in A \text{ or } w \in B\}$$



Example:-  $A = \{1, 3\}$   
 $A \cup B = \{1, 2, 3, 4\}$   $B = \{2, 4\}$

Additive law in probability:-

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

$$P(A \cup B) = P(A) + P(B) \Rightarrow \text{if mutually exclusive}$$

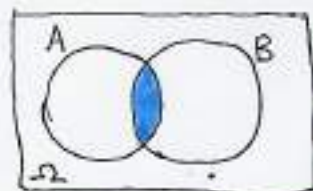
$$P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap B) - P(A \cap C) - P(B \cap C) + P(A \cap B \cap C)$$

## 2- Intersection of 2 events:-

$A \cap B$

an event containing all elementary events that belong to A and B or both occurring together

$$A \cap B = \{w: w \in A \text{ and } w \in B\}$$



Example:-  
 $A = \{TT, HT\}$   $B = \{TT, TH\}$   
 $A \cap B = \{TT\}$

Multiplicative law in probability:-

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

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

$$P(A \cap B) = P(A) \times P(B) \Rightarrow \text{if independent}$$

## 3- Complement of an event:-

$\bar{A}$

an event that occurs when A doesn't

$$\bar{A} = \{w: w \in \Omega, w \notin A\}$$



$\Omega = \{1, 2, 3, 4, 5, 6\}$   $A = \{1, 2\}$   
 $\bar{A} = \{3, 4, 5, 6\}$

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

De Morgan's Law:-

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

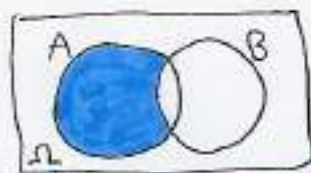
$$P(\bar{A} \cap \bar{B}) = P(\overline{A \cup B})$$

## 4- Difference Between 2 events:-

$A \setminus B$

the event of the occurrence of A but not B

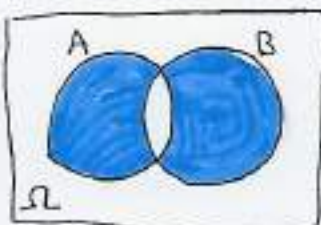
$$A \setminus B = A \cap \bar{B} = \{w: w \in A \text{ and } w \notin B\}$$



## 5- exactly one of event $A \Delta B$

the event of A and B occur but not together  $A \Delta B = \{x: x \in A \cap \bar{B} \text{ or } x \in \bar{A} \cap B\}$   
 $= \{x: x \in A \cup B \text{ and } x \notin A \cap B\}$

$$A \Delta B = (A \setminus B) \cup (B \setminus A)$$



## 6- impossible event $\emptyset$

The event of an occurrence being impossible

Example:-

$$A \cap \bar{A} = \emptyset$$



## certain event $\Omega$

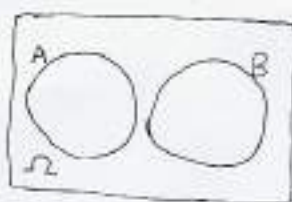
An event that is sure to happen

Example:-  $* A \cup \bar{A} = \Omega$

\* Flipping a coin once and getting a head and tail

## Mutually Exclusive event

they can't occur at the same time



that means  $A \cap B = \emptyset$

Example:-

getting on a plane or not getting on a plane

## conditional probability

the probability of A occurring given that B also occurs

$$P(A|B) \quad 2.4$$

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

independent  $\leftarrow P(A|B) = P(A)$

## independent events

The occurrence of B doesn't effect the occurrence of A

\* باختيار حدث ما لم يمتد بحدوثه من اتي  
اذا لم او تنجح يسي. ما يترون الى بعد

\* Mutually exclusive events are not the same as independent

\* A and B  $\rightarrow$  independent

\*  $\bar{A}$  and B  $\rightarrow$  independent

\*  $\bar{A}$  and  $\bar{B} \rightarrow$  independent

## conditions for partitioning $\Omega$ :

1-  $Z_i \neq \emptyset \rightarrow$  no impossible events

2-  $Z_i \cap Z_j = \emptyset$  (mutually exclusive)

3-  $\cup Z_i = \Omega \rightarrow$  union gives  $\Omega$

## Total probability theory

$$P(B) = \sum P(Z_k) P(B|Z_k)$$



$$P(B) = P(B \cap Z_1) + P(B \cap Z_2) + P(B \cap Z_3) + \dots + P(B \cap Z_k)$$

$$= P(Z_1)P(B|Z_1) + P(Z_2)P(B|Z_2) + P(Z_3)P(B|Z_3) + \dots + P(Z_k)P(B|Z_k)$$

## Bayes' Theorem

an equation that relates two conditional probabilities of the form  $P(B|A)$  and  $P(A|B)$

$$P(Z_i|B) = \frac{P(Z_i) P(B|Z_i)}{\sum P(Z_k) P(B|Z_k)}$$



# CHAPTER 3

Random variables and  
Probability distribution

3.1

3.2

3.3



# Random variable

3.1 3.2 3.3

\* denoted by  $X, Z, Y$  and the values by  $x, z, y$

\*  $\{W \in \Omega; X(W) \leq x\} \in \mathcal{F} \quad \forall x \in \mathbb{R}$   
 $(X \leq x) \uparrow_{x \rightarrow \infty} 1$

## Discrete

finite or infinite countable

Proving  $X$  is a discrete random variable has 2 conditions:-

1-  $P(X=x) \geq 0$

2-  $\sum P(X=x) = 1$



## Continuous

uncountable set of numbers

Probability density function (P.d.F)

Properties and conditions:-

1-  $P(a \leq x < b) \geq 0$

2-  $\int_{-\infty}^{+\infty} f(x) dx = 1$

\*  $P(a \leq x \leq b) = \int_a^b f(x) dx$

\*  $f(x) = \frac{d}{dx} F(x)$

distribution function

$F(x) = \int_{-\infty}^x f(x) dx$

Mean

$\mu = E(x) = \int_{-\infty}^{+\infty} x f(x) dx$

Variance

$\sigma^2 = \int_{-\infty}^{+\infty} x^2 f(x) dx - \mu^2$

some special distribution

## Exponential

times between events occurring

\*  $\lambda > 0$

P.d.F

$f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{for } x \geq 0 \\ 0 & \text{for } x < 0 \end{cases}$

$F(x) = \begin{cases} 1 - e^{-\lambda x} & \text{for } x \geq 0 \\ 0 & \text{for } x < 0 \end{cases}$

distribution function

mean

$\mu = \frac{1}{\lambda}$

Variance

$\sigma^2 = \frac{1}{\lambda^2}$

Standard deviation

$\sigma = \frac{1}{\lambda}$

## Normal

natural phenomena in Real life

P.d.F

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

distribution function

$F(x) = \frac{1}{\sigma \sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt$

\* In a graph it has a bell shaped

\* one peak in the center = mean = mode = median

## Standardizing Normal distribution :-

to help compute areas under P.d.F curve

$Z = \frac{x - \mu}{\sigma}$

Probability mass function (P.M.F)  $f_x(x)$

Properties and conditions:-

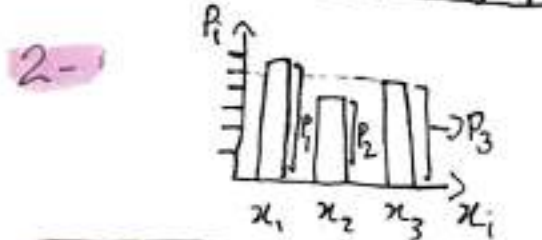
1-  $P(X=x) \geq 0$

2-  $\sum P(X=x) = 1$

Can be represented in a table or graph

1-

X Random Variable	$x_1$	$x_2$	...	$x_n$	$\Sigma$
$P(X=x) = f_x(x)$	$P_1$	$P_2$	...	$P_n$	1



mean:-  $\mu = E(x) = \sum x_i P(X=x_i)$

Variance:-  $\sigma^2 = E(x^2) - \mu^2 = \sum x_i^2 P(X=x_i) - \mu^2$

Standard deviation:-  $\sigma = \sqrt{\sigma^2}$

Distribution function  $F_x(x)$

$F_x(x) = \sum P(X=x)$

\*  $0 \leq F_x(x) \leq 1$

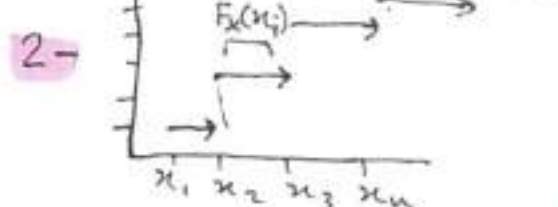
\*  $\lim_{x \rightarrow -\infty} F_x(x) = 0, \lim_{x \rightarrow +\infty} F_x(x) = 1$

$F_x(x) = \begin{cases} 0 & x < x_1 \\ P_1 & x_1 \leq x < x_2 \\ P_1 + P_2 & x_2 \leq x < x_3 \\ \vdots & \vdots \\ P_1 + P_2 + \dots + P_n = 1 & x \geq x_n \end{cases}$

Can be represented in a table or graph

1-

For	$x < x_1$	$x_1 \leq x < x_2$	...
$F_x(x)$	0	$P_1$	...



some Special distributions

## Binomial

2 possible outcomes success (p) and failure (q) repeated n times

\*  $1 > p > 0$  \*  $1 - p = q$

\* trials are independent

P.M.F

$P(X=x) = \binom{n}{x} p^x q^{n-x}$

mean

$\mu = E(x) = np$

Variance

$\sigma^2 = npq$

## Poisson

events occurring in an amount of specific time

\*  $\lambda > 0$

P.M.F

$P(X=x) = \frac{\lambda^x e^{-\lambda}}{x!}$

Distribution function

$F_x(x) = \sum P(X=x)$

Variance

$\sigma^2 = \lambda$

mean

$\mu = E(x) = \lambda$



# CHAPTER 4

Introduction to  
statistical inference

4.1 4.2 4.3 4.4 4.5 4.6



we have something called an

## Estimator:-

a statistic (function of the random sample)  
whose value depends on the particular sample  
is drawn on the population

in a sample space  $\Omega$  we take a simple random sample of  $\Omega$  denoted by  $\mathcal{X}$

$\Omega$  = population  
 $\mathcal{X}$  = sample

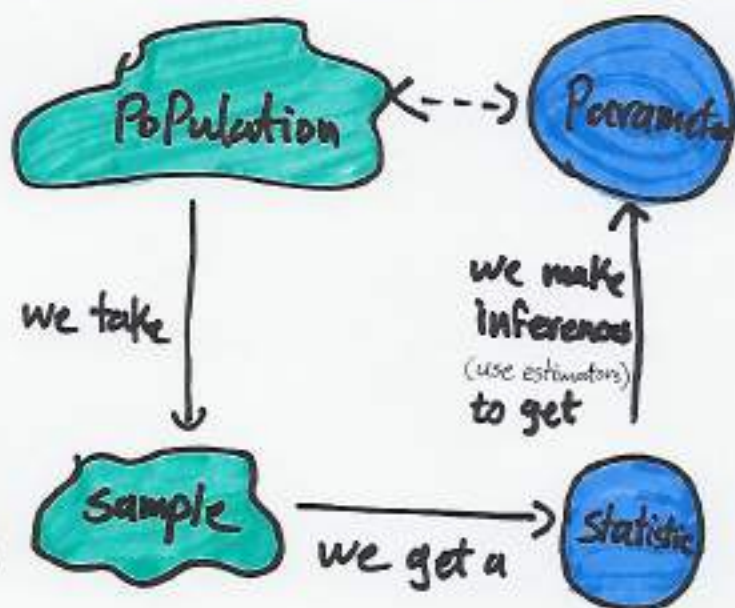
the mean of  $\mathcal{X}$  is  $\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$  which is an estimator

the variance of  $\mathcal{X}$  is  $S^2 = \frac{\sum (X_i - \bar{X})^2}{n-1}$  which is an estimator

\* We use the estimators to make predictions about the value of a Population. explained in this graph

Since the Population is HUGE we can't get an accurate value of parameter

so we use samples



(( شرح المفردات estimators وكيف نستعملها للإيجاد Parameter ))

والمعاني  $\bar{X}$  و  $S^2$  هي هذه المقاييس التي نستخدمها كثيرًا



# Standard Normal Probabilities $Z_\alpha$

The value on the real axis for which is the area between the left part of the curve of  $p(x)$  and the straight  $x = Z_\alpha$  equal to  $\alpha$



الفديو  
هذا يشرح  
الفكرة افضل

استحبوا على كل ذا الكلام. ففي تشابير  $\chi^2$  عن Normal distribution قلنا بأن عشان نخليه قياسي نستعمل ذا القانون  $Z = \frac{x - \mu}{\frac{\sigma}{\sqrt{n}}}$  ونطلع لنا قيمة. يعني الحالة يعطيك القيمة الاخرى ويبينك توجد القيمة الادوية

Confidence level (1- $\alpha$ ) %	Confidence coefficient $\alpha$	Reliability coefficient	
		$Z_{1-(\frac{\alpha}{2})}$	$Z_{1-(\frac{\alpha}{2})}$
0.90	0.10	$Z_{0.95}$	1.645
0.95	0.05	$Z_{0.975}$	1.96
0.98	0.02	$Z_{0.99}$	2.33
0.99	0.01	$Z_{0.995}$	2.575

## the central limit theorem

When the sample size is large ( $n \geq 30$ ) then the sample distribution follows a normal distribution

### Sample distribution of mean

mean:  $\mu_{\bar{x}} = \mu$

standard error:  $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$

Standardizing the distribution:

$$Z = \frac{\bar{x} - \mu_{\bar{x}}}{\frac{\sigma}{\sqrt{n}}} = \frac{\bar{x} - \mu_{\bar{x}}}{\sigma_{\bar{x}}}$$

### Sample distribution of proportions

mean:  $\mu_p = p$  standard error:  $\sigma_p = \sqrt{\frac{p(1-p)}{n}}$

Standardizing the distribution:

$$Z = \frac{\bar{x} - \mu_{\bar{x}}}{\sigma_p}$$

Point estimation: - an estimate of the Population Parameter by a single number

## Types of estimations

### Estimation of the Population mean

Interval estimation: - Predicted to have the value of Population Parameter with a certain Probability

#### Interval estimation

\* confidence interval: - a range of values (level of confidence) that have the true value of population Parameter with specific probability

#### Point estimation

Parameter = statistic  
 $\mu = \bar{x}$

case	Population	Sample size	Standard deviation S.D	confidence interval
1	Normal	any	Known	$\bar{x} \pm Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$
2	any	large $n \geq 30$	Known	$\bar{x} \pm Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$
3	any	large $n \geq 30$	unknown	$\bar{x} \pm Z_{1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}}$

cases and their confidence interval

But replaced with Sample S.D

\* margin error: - the maximum error of using the value of statistics instead of Parameter

$$\delta_{\mu} = Z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \text{ or } \delta_{\mu} = Z_{1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}}$$

to find sample size:

$$n = \left[ \frac{Z_{1-\frac{\alpha}{2}} \sigma}{\delta_{\mu}} \right]^2$$

### Estimation of the Population Proportion

\* confidence interval for Population Proportion: -

$$P = \hat{p} \pm Z_{1-\frac{\alpha}{2}} \left( \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right)$$

$$P = \frac{N(A)}{N}, \hat{p} = \frac{n(A)}{n}$$

\* Margin Error: -

$$\delta_p = \hat{p} \pm Z_{1-\frac{\alpha}{2}} \left( \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right)$$

\* Sample size: -

1-  $\hat{p}$  is given:

$$n = \left[ \frac{Z_{1-\frac{\alpha}{2}}}{\delta_p} \right]^2 \hat{p}(1-\hat{p})$$

2-  $\hat{p}$  not given:

$$n = 0.25 \left[ \frac{Z_{1-\frac{\alpha}{2}}}{\delta_p} \right]^2$$

3- L is given:

$$L = \text{maximum length}$$

$$n = \left[ \frac{Z_{1-\frac{\alpha}{2}}}{L} \right]^2$$



# Statistical Hypothesis

4.4 4.5 4.6

an argument about a specific statistical question, it can be true or wrong

Null Hypothesis  $H_0$

a statement under an investigation or testing

- \* we usually use =
- \* represent a statement of "no effect" or "no difference"

باختیار null نپذیریم  
شرکت ادعت ان حلبیا  
وزنه 1 kg است  
null وان اینها عینه  
وطلح 0.99 kg است  
alternate  
لفرضیه‌ای غیر از طریق

Alternate Hypothesis  $H_1$

a statement we will adopt in the situation in which the evidence is so strong that we will reject the null

- \* help us know if the null hypothesis is true or not
- \* we usually use  $<, \neq, >$

Types of errors:-

actual situation

	$H_0$ true	$H_0$ is false
reject $H_0$	correct decision	Type II error $P(\text{Type II error}) = \beta$
don't reject $H_0$	Type I error $P(\text{Type I error}) = \alpha$	correct decision

$\alpha$  = significance error

$1 - \beta$  = the power of test

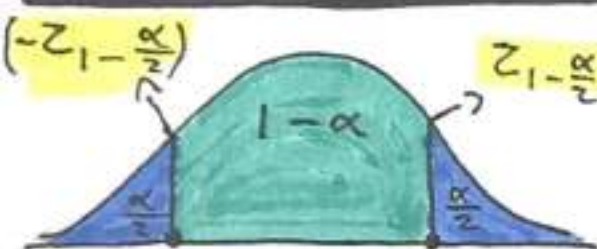
test statistic  
و نشون ای فرضیه درست

The critical regions and values:-

\* a region that produced by the values that corresponds to the rejection of the null hypothesis at some chosen level of significance

\* the value(s) that determine the critical region(s) is called critical value

$-(Z_{1-\alpha}), (Z_{1-\frac{\alpha}{2}}), (-Z_{1-\alpha}), (-Z_{1-\frac{\alpha}{2}})$



Two tailed region

$\mu \neq \mu_0$   
( $H_1 \neq H_0$ )



Left tailed region

$\mu < \mu_0$   
( $H_1 < H_0$ )



Right tailed region

$\mu > \mu_0$   
( $H_1 > H_0$ )

Accepted region  
قبول فيه فرضیه  $H_0$

rejected region  
نرفتن فيه فرضیه  $H_0$

\* نذر ای critical value تست  
عن طریق العلاقة بين  $H_0$  و  $H_1$   
وايضاً عن طريق التوزيع الطبيعي  
القياسي (الرسمه)

For mean:-

$$Z_0 = \begin{cases} \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} & \sigma \text{ know} \\ \frac{\bar{X} - \mu_0}{\frac{s}{\sqrt{n}}} & \sigma \text{ unknown and } n \geq 30 \end{cases}$$

Test statistic

is a statistic whose calculated value from the sample is used to make a decision on Hypothesis test

Hypothesis test

For proportion:-

$$Z_0 = \frac{\hat{P} - P_0}{\sqrt{\frac{P_0(1-P_0)}{n}}}$$

Population mean

How to:-

- 1- determine  $H_0, H_1$
- 2- determine test statistic
- 3- determine critical region
- 4- make a decision whether  $H_0$  is true or not

Power

$H_0$  is true

P value:- Assuming  $H_0$  is true, the probability that the test statistic will take values as extreme as or more extreme than the observed test statistic

check example: 4.5.1 & 4.6.1

\* the smaller p value the stronger the evidence against  $H_0$

Population Proportion

- 1- determine  $H_0, H_1$
- 2- determine test statistic
- 3- determine critical region
- 4- make a decision



# CHAPTER 5

Correlation and  
Regression

5.1 ~~5.1.1~~ ~~5.1.2~~ 5.2



# Pearson's correlation coefficient $r$ :- 5.1

a measure of linear correlation between  $X$  &  $Y$

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \cdot \sqrt{\sum (y_i - \bar{y})^2}}$$

or

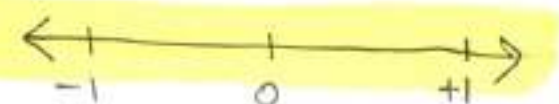
$$r = \frac{n \cdot \sum x_i y_i - \sum x_i \cdot \sum y_i}{\sqrt{n \cdot \sum y_i^2 - (\sum y_i)^2} \cdot \sqrt{n \cdot \sum x_i^2 - (\sum x_i)^2}}$$

Relationship between $X$ & $Y$	Range of $r$
No linear $S_y = 0$ $S_x = 0$	$r = 0$
very weak	$0 < \pm r \leq 0.30$
weak	$0.3 < \pm r \leq 0.5$
moderately strong linear	$0.5 < \pm r \leq 0.7$
Strong	$0.7 < \pm r \leq 0.86$
Very strong	$0.86 < \pm r \leq 1$
complete	$\pm r = 1$

\* The values ranges between  $-1$  and  $+1$

\*  $+r$  means  $X \propto Y$  and it's called Positive monotone

\*  $-r$  means  $X \propto \frac{1}{Y}$  and it's called negative monotone

\*  كلما كانت  $r$  قريبة من  $\pm 1$  كلما قوت العلاقة وإذا  $r = 0$  فلا توجد علاقة

## \* Scatter plot :-

a graph of data give in the form of  $(x_i, y_i)$



\* check page 205 For Examples

## Equation for linear regression:

$$\hat{y} = \hat{a} + \hat{b}x \quad 5.2$$

\* constant  $\hat{b}$  : the slope \* constant  $\hat{a}$  : Y-intercept

$$\hat{b} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \hat{a} = \bar{y} - \hat{b}\bar{x}$$

$\bar{y} = \frac{\sum y_i}{n_y} \quad \bar{x} = \frac{\sum x_i}{n_x}$

## coefficient of determination $r^2$ :

$$r^2 = \frac{SSR}{SS_{tot}} = 1 - \frac{SSE}{SS_{tot}} = \%$$

$$* SS_{tot} = \sum (y_i - \bar{y})^2 = SSR + SSE$$

total sum of squared deviation

$$* SSR = \sum (\bar{y} - \hat{y}_i)^2$$

sum of squared regression error (Explained variation)

\* the values range  $0 \leq r^2 \leq 1$

\*  $r^2 = 0$  no explanatory value

\*  $r^2 = 1$  explains 100%

$$* SSE = \sum (y_i - \hat{y}_i)^2$$

## تمثيل معادلة الانحدار الخطي

نحتاج نقطتين فقط

فنختار واسم من ذولا

$$1 - (X_1, Y_1) \text{ نقطة}, (X_2, Y_2) \text{ نقطة}$$

نأخذ قيمة  $x$  ونعوضها

$$2 - (\bar{x}, \bar{y}) \text{ نقطة}, (X_1, \hat{Y}_1) \text{ نقطة}$$

نأخذ  $\bar{x}$  ونعوضها في معادلة  $\bar{y}$

ونأخذ قيمة  $x$  ونعوضها في معادلة  $y$

\* scatter plot يعطينا تمثيل النقطتين على محورين ونأخذ قيمتين