بسم الله الرحمن الرحيم

السلام عليكم ورحمة الله وبركاته

# Biostatistics
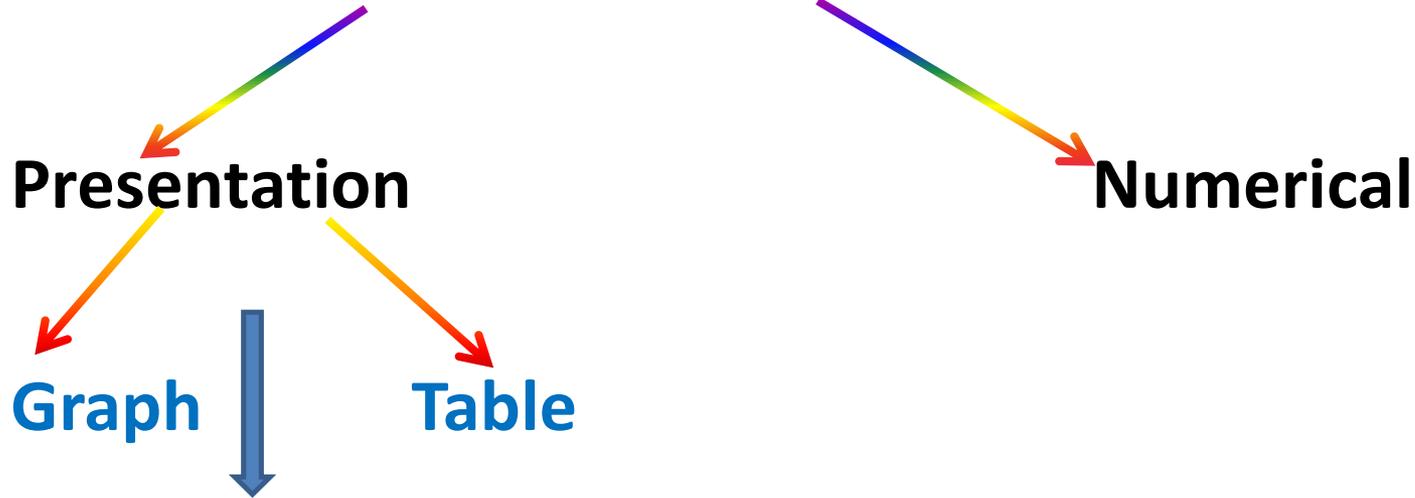
## L IV
## 12 –July 2023

# PROF.  DR. WAQAR   AL-KUBAISY

# Description statistics summarization

**Presentation**                                **Numerical**

**Graph**            **Table**

-*this approach might not be enough,*

-*comparisons between one set of data & another*

-*summarize data by one more step further .*

-*presenting a set of data by a*

-            *single Numerical value*

# The central value as representative value in a set of data,

## 1-Measures of central tendencies (Location) .
A value around which the data has a tendency to congregate (come together )or cluster

## 2-Measures of Dispersion, scatter around average
A value which measures
the degree to which the data are  or  are not, spread out

**The central value as**

1-Measures of central tendencies (Location) .
  A value around which the data has a tendency
to congregate (come together )or cluster
2-Measures of Dispersion, scatter around average
  A value which measures
the degree to which the data are  or  are not ,
spread out

**1-Measures of central tendencies (Location)**

75, 75, 75, 75, 75, 75,          Mean =  ????

75, 70, 75. 80, 85.              Mean =  ????

60, 65, 55, 70, 75, 75, ,70, 80, Mean=  ????

$$\overline{X} = \frac{\Sigma X}{N}$$

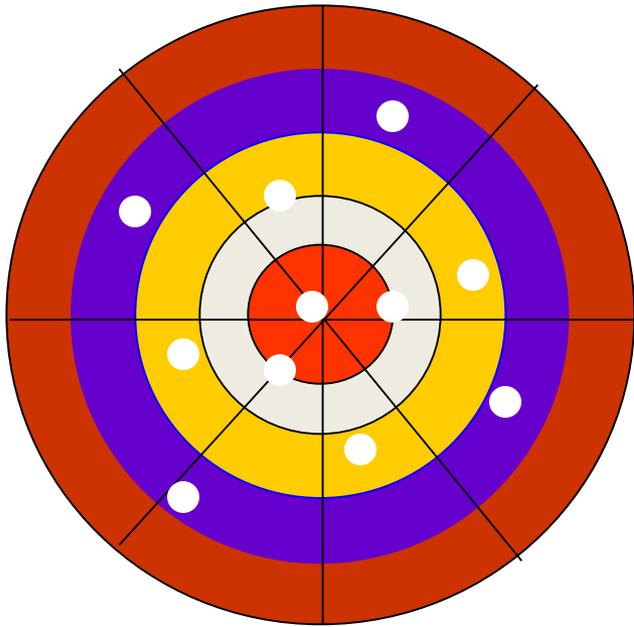**2-Measures of Dispersion,**

7/12/2023

5

The central value as
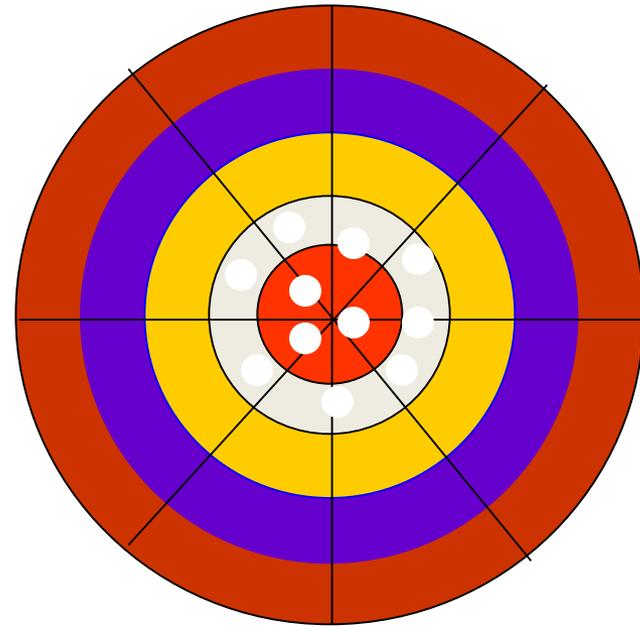1-Measures of central tendencies
2-Measures of Dispersion,

# Measures of Dispersion
## (Measures of Variation)
## (Measures of Scattering)
## Measures of spread

# Measures of Dispersion



**SHOOTER A**

**SHOOTER B**

*Both shooters are hitting around the "centre" but shooter B is more "accurate"*

# Measures of Dispersion

1- **Range**

**2-Interquartile range**

**3- Variance**

**4- Stander Deviation**

**5- Coefficient of variance**

> **the choice of the most appropriate measure depends crucially on the type of data involved**

# Measures of spread

**Measuring of spread are very useful.**

**There are <u>three main</u> measures in common use .**

**once again the type of data influence the choice of an appropriate measure**

> **the choice of the most appropriate measure depends crucially on the type of data involved**

# The Range
**simplest**

**most obvious one of dispersion.**

**It is the distance from the smallest to the largest**

**It *Obtained by***

**subtracting lowest value from the highest value in a set of data .**

**Pulse rate   70   76   74   78   72    74    76**

**Range = 78 – 70    =**

The range is **best written**
like rang of data (from- to)   70-78
rather than single-valued difference  which is much less informative

**The range is not affected by skewness**

    70 72   74   76   76   78   78         78-70  70-78

**sensitive to the addition or removal of an outlier value**

   66 70  74    90, 100  120 124      124-66  66-124
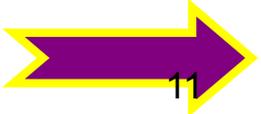
   66 70  74    90, 100  120 124  **545**   66-545

   I**ts disadvantage**

  it is based on **only two observations**

          (the lowest and highest value) and

❖ give no idea about others,

❖  not take into consideration other values in data

❖sensitive to an **outlier value**         **Therefore**

❖  **It is not very useful measures of variation,**

✓      **because it does not use  other observation**

*Therefore* ;

*Therefore ;*

**Sensitive an outlier value**

**Interquartile rang  (I q r).**

✓ *measure the variation of one observation from the other*
✓ **Standard deviation**

**Interquartile rang  (I q r).** →

7/12/2023

12

# Percentile

**A percentile provides information about how the data are spread over the interval from the smallest value to the largest value.**

The pth percentile (25%) (30%):
is a value such that at least p percent of the observations are less than or equal to this value   **and**
**at least (100 - p)** (75% ) (70%) percent of the observations **are greater than or equal** to this value.

The pth percentile is a value so that **roughly p% of the data are smaller and (100-p)% of the data are larger**.

Three Steps

**Three Steps for computing a percentile**.
1. Sort the data from low to high;
2. Count the number of values (n);
3. Select the p*(n+1) observation

**Examples**

The following data represents cotinine levels in saliva (ng/ml) after smoking. We want to compute the 50th percentile.

73, 58, 67, 93, 33, 18, 147

Sorted data: 18, 33, 58, 67, 73, 93, 147

There are n=7 observations.

**Select 0.50*(7+1)=4th observation**.

Therefore, the **50th percentile** equals 67.

**Notice that there are**

**three observations larger than 67 and**

**three observations smaller than 67.**

**Examples**

The following data represents cotinine levels in saliva (ng/ml) after smoking. We want to compute the 20th percentile.

73, 58, 67, 93, 33, 18, 147

Sorted data: 18, 33, 58, 67, 73, 93, 147

Suppose we want to compute the 20th percentile.
Notice that $p*(n+1)$ = **0.20*(7+1)=1.6.** This is not a whole number so we select halfway **between 1st and 2nd** observation
they have to go **six tenths of the way to the** second value.

# Calculation of percentile value

> **The pth percentile is**
> **the value in the p/100 (n+1) th  position.**

**For example**
**the 20th percentile**

Calculation of percentile value
**the birth weight(grm) of 30 infants which  we put in ascending order.**

 2860   2994   3193   3266   3287   3303   3388

3399   3400   3421   3447   3508   3541    3594

3613   3615   3650   3666  3710  3798

3800  3886   3896   4006   4010    4090   4094

 4200   4206   4490

# Calculation of percentile value

The pth percentile is
 the <u>value</u> in the p/100 (n+1) th  position.

the 20th percentile is the
  20/100(n+1)   with the BW values
   20/100 (30 +1)
 0.2x31 observations= 6.2observation

the birth weight of 30 infants which we put in ascending order.
 2860   2994  3193   3266   3287   3303    3388  3399  3400
3421   3447   3508  3541   3594   3613   3615   3650   3666
3710   3798  3800 3886   3896   4006   4010    4090   4094
4200   4206   4490

**The 6th value is 3303 g
the 7th value is 3388 g**

**a difference of** <span style="color:red">**85 g**</span>

**the 20th percentile is**

**3303 + 0.2 of 85 g**
  **which is**
**3303g + 0.2x 85 g   =**
**=3303g+17g**
**= 3320  g**

the birth weight of 30 infants which we put in ascending order.
 2860   2994  3193   3266   3287
 3303   3388  3399  3400  3421   3447
 3508  3541  3594  3613  3615   3650
 3666    3710   3798  3800 3886
 3896   4006   4010   4090   4094
 4200  4206   4490

**The pth percentile is**
 **the** <u>**value**</u> **in the p/100 (n+1) th  position.**

**Similarly we could calculate**

cont. ……Calculation of percentile value

## the deciles

**which subdivide the data values**

**into 10 (not 100 )equal division,**

**and**

## Quintiles

which sub-divide the values into

five equal –sized groups

**Collectively we call**

❖ **percentiles,**

❖ **deciles** divide the sorted data into ten equal parts, so that each part

represents 1/10 of the sample or population. **and**

❖**quintiles**

the birth weight of 30 infants which we put in ascending order.

| 2860 | 2994 | 3193 | 3266 | 3287 | 3303 |
|------|------|------|------|------|------|
| 3388 | 3399 | 3400 | 3421 | 3447 | 3508 |
| 3541 | 3594 | 3613 | 3615 | 3650 | 3666 |
| 3710 | 3798 | 3800 | 3886 | 3896 | 4006 |
| 4010 | 4090 | 4094 | 4200 | 4206 | 4490 |

**The pth percentile is
the value in the p/100 (n+1) th position.**

A quartile is :

 a division of observations into four defined     25%   50%

**Interquartile rang  (i q r).**
**One solution to the problem of the sensitivity**
 **to extreme value  (outlier) is to**


✓**chop the quarter(25 percent) of the values of both ends**
**of the distribution**
    **(which removes any troublesome outliers)**


**then measure the range of the remaining values**


❑ **this distance is called**
❑**interquartile range or i q r .**

# Calculation of iqr

## To calculate iqr we need to determine two values

| first quarantile ( Q1) The value which cuts off the **bottom** **25** percent of values | third quarantile (Q3) **The value which** **cuts off the top 25** percent of **values,** |

The interquartile range is then written as (Q1 to Q3)

**31X 0.25 = 7.75**

31X .75 = 23.25

the birth weight of 30 infants which we put in ascending order.

2860   2994   3193   3266   3287   3303   3388   3399

3400   3421   3447   3508   3541   3594   3613   3615

3650   3666   3710   3798   3800   3886   3896   4006

4010   4090   4094   4200   4206   4490

**The pth percentile is the value in the p/100 (n+1) th position.**

with the BW data
Q1= 3396.25g and
Q3 = 3923.50 g

**7.75**th  3399-3388=11x.75=8.25+3388= 3396.25

0.75x 31=**23.25**th

**4006-3896=110x.25=27.5+3896=3923.5**

**the birth weight of 30 infants which we put in ascending order.**

**2860   2994  3193   3266   3287   3303   3388   3399  3400**
**3421   3447   3508 3541  3594  3613 3615 3650   3666**
**3710   3798   3800 3886  3896      4006   4010   4090   4094**
**4200  4206   4490**

**Therefore iqr =   3369. 25 to 3923.50)g**

**the middle 50 percent**

23

# Calculation of iqr

the middle **50 percent** of infant weighed
   between **3396.25 and 3923.50 g**
✓**The interquartile range**
**indicate**
❖  the spread of the middle 50%of the distribution,

❖   **together with the median** **is  useful** adjunct
(accessory) to the range

❖ it is **less sensitive** to the **size of the sample**
providing that  this is not too small

**The interquartile range is not affected either by** Outlier skewness

**BUT**

**it** does not use all of the information in the data set since it ignores the bottom and top quarter of values.

✓*measure the variation of one observation from the other*
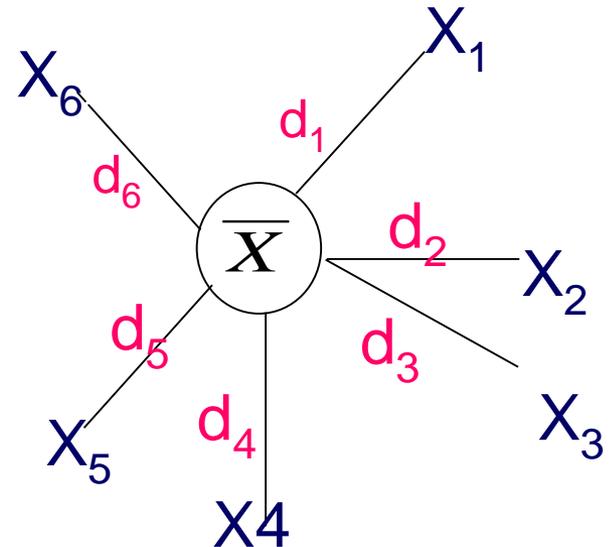
✓**Standard deviation**

**75, 70, 75. 80, 85.**          Mean =  ????

**60, 65, 55, 70, 75, 75, ,70, 80, Mean=  ????**

$$\overline{X} \quad = \quad \frac{\sum X}{N}$$

the **mean** (average) distance of all data values from the **over all mean** of all values.

$X_1$
$X_6$
$d_1$
$d_6$
$\overline{X}$
$d_2$
$X_2$
$d_5$
$d_3$
$X_3$
$d_4$
$X_5$
X4

## Standard deviation (SD)

 **The limitation of iqr it does not use all of the information in the data since it omits the top and bottom quarter of values.**

**An alternative approach use the idea of summarizing spread   by    measuring**
the **mean** (average) **distance of all data values from the over all mean of all values**.

▪**The smaller the mean distance is**
✓  **the narrower the spread of values must be**
        **and visa versa**
this is known as **standard deviation**