

Chapter 3 Descriptive Statistics

Numerical Measures

↳ Sample statistics & measures are computed for data from a sample
لإحصاءات العينة - يتم حسابها لقياس البيانات بأشرف من عينته

↳ Population parameters & measures are computed for data from a population

لإحصاءات المجتمع - يتم حسابها لقياس البيانات بأشرف من المجتمع

↳ A sample statistics is ^{تقدير نقطي} point estimator of the corresponding population parameter

3.1 Measures of location

[1] Mean - الوسط الحسابي

[2] Median - الوسيط

[3] Mode

[4] Percentiles

[5] Quartiles

[1] Mean (average) - الوسط الحسابي

↳ a measure of Central location for the data
مقياس للموقع المركزي للبيانات

Mean = $\frac{\text{sum of values}}{\text{number of values}}$ الوسط الحسابي

↳ for sample → $\bar{x} = \frac{\sum x_i}{n}$

\bar{x}
المتوسط

n = sample size

↳ for population → $\mu = \frac{\sum x_i}{N}$

μ
متوسط

n = population size

\bar{x} is a point estimator for μ

Example → Give the sample data 45, 95, 60, 77, 80 and find the mean

sample mean = $\frac{\sum x_i}{n}$

$\bar{x} = \frac{45 + 95 + 80 + 60 + 77}{5} = 71.4$

How to use Calculator (تطبيق على الآلة الحاسبة)

خطوة

① Mode → 2 (SD appears on screen)

- ② 45 → M+ ⇒ n = 1
- 95 → M+ ⇒ n = 2
- 80 → M+ ⇒ n = 3
- 60 → M+ ⇒ n = 4
- 77 → M+ ⇒ n = 5

هكذا حتى نصل إلى القيمة

نلاحظ على الشاشة عدد 5

تجربنا تطبيقاً
الحاسبة قبل البداية
Shift + mode + 3 + = + =

③ shift [2] → \bar{x}

\bar{x}	$x \sigma_n$	$x \sigma_{n-1}$
1	2	3

الجواب $\bar{x} = 71.4$ = answer [1]

* الآلة الحاسبة دقيقة ويتم إنتاجها الوقت في الأوقات
بناءً على سرعة استقار الآلة الحاسبة

* أهمية اختيار وضع الأرقام في البداية "خطوة 1"

[2] Median الوسيط

↳ is the value in the middle approximately 50% of data values are less than or equal median and approximately 50% of the value are greater than or equal the median

لـ الوسيط ٥٠٪ من البيانات هي أقل منه "أقل منه"
 ٥٠٪ من البيانات هي أكبر منه "أكثر منه"

How to find median

[1] arrange the data in ascending order
 لـ ترتيب البيانات ترتيباً

[2] Two Cases are possible

① if n is odd "فردية"
 ↳ then median is the middle value

② if n is even "زوجية"
 ↳ then median is the average of the two middle values

15, 20, 30, 23, 70, 80 ← مثال
 ↳ Median = $\frac{30 + 23}{2} = \underline{26}$
 ← تقرب العدد من الأعداد

Example → Find the median → 22, 10, 18, 22, 16, 17, 20

① arrange → 10, 16, 17, 18, 20, 22, 22

n = 7 median = 18

[3] Mode الكنوال

↳ the value(s) that has highest frequency
 "occure the most"
 له يعني العلى تكراراً

ex → find the mode → 100, 140, 120, 123, 137
 → No mode

ex → find the mode → 99, 140, 123, 140, 100, 140, 100

→ Mode = 140
 □ unimodal data → one mode

ex → find the mode → 8, 12, 15, 8, 12, 18, 17, 10

→ Mode = 8, 12

□ Bimodal data → two mode

Ex → find the mode → A, AB, O, O, B, O, AB, A, O

→ Mode = O

Note → when ever the data set contains extreme values
 the median is after the preferred measure

→ □ Trimodal data → 3 mode

Outlier

extreme values
5000 - 2000
20,000
outlier

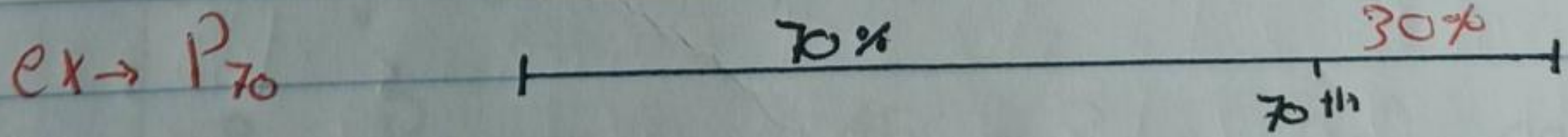
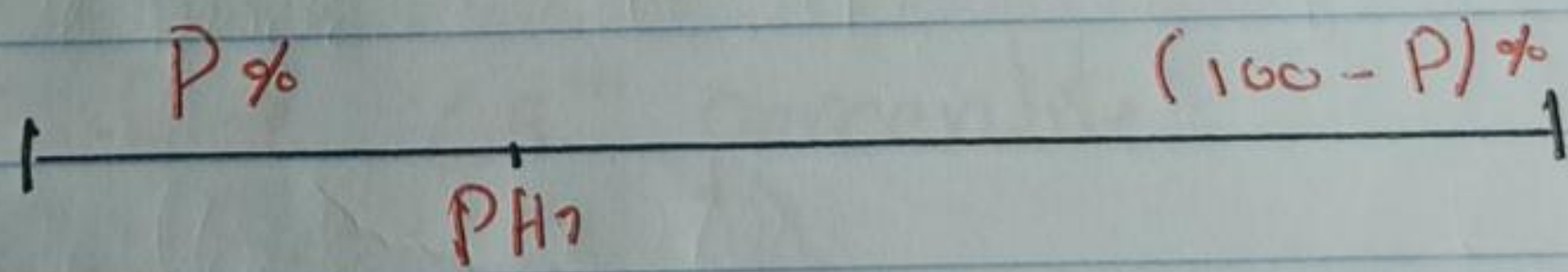
outlier 200

Outlier

Percentiles

Pth percentile is the value that approximately P% of the values are less than or equal this value and (100-p)% of the values are greater or equal to this value

ex -> P30 30 or more 30 or less



How to find a Percentile

① arrange data in a ascending order
الترتيب لبياناتنا تصاعدياً

② Compute $i = \frac{P}{100} \times n$

i = index (position)
موقع القيمة بالترتيب

n = حجم العينة
Sample size

[Case 1] if i is not integer, then the P^{th} percentile is the value in the next integer position
غير صحيح
لأنه أقرب إلى التالي

[Case 2] if i is an integer, then the P^{th} percentile is the average of the two values in the i and $(i+1)$ position
لأنه يكون إذا كان عدد صحيح، نأخذ وسطه

Example → Given the data 12, 18, 15, 25, 10, 8, 12, 9

Find the 65th percentile

arrange → 8, 9, 10, 12, 12, 15, 18, 25

$$P = \frac{P}{100} \times n = \frac{65}{100} \times 8 = 5.2 \text{ not integer}$$

غير صحيح

P_{65} = The value of the 6th position

$P_{65} = 15$

② Find the P_{25} (25th percentile)

arrange \rightarrow 8, 9, 10, 12, 12, 15, 18, 25

$$P = \frac{25}{100} \times 8 = \frac{2}{1} \text{ integer}$$

↳ الثاني، الثاني

$P_{25} \rightarrow$ The average of the values in the second and third position

$$P_{25} = \frac{9 + 10}{2} = 9.5$$

~~~~~  $\rightarrow$  Note المعدل

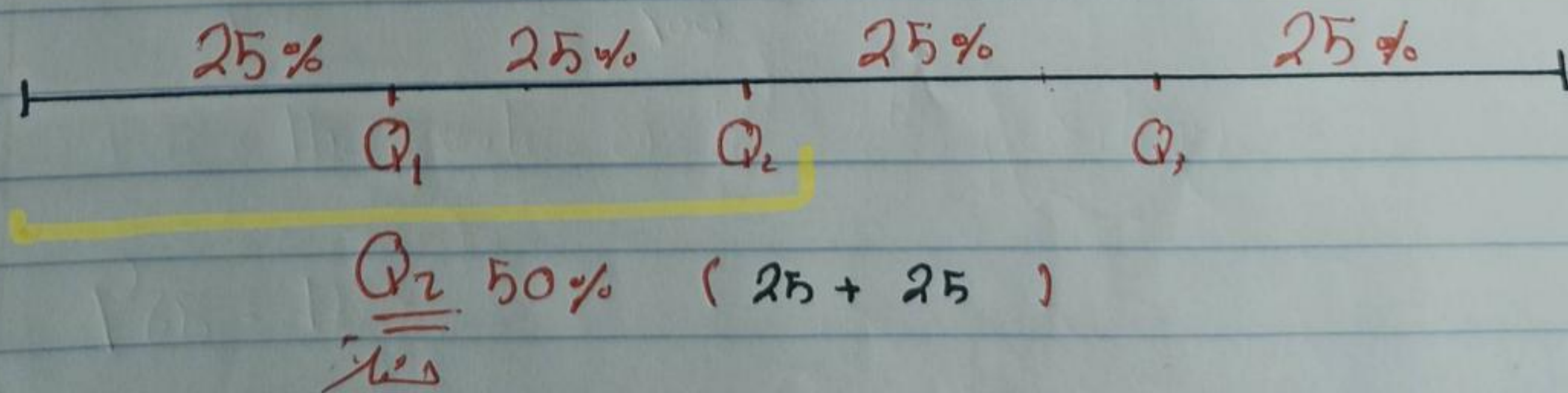
$P_{50}$  (50<sup>th</sup> percentile) = Median

[5] Quartiles الكوارتيل

①  $\hookrightarrow$  First quartile  $Q_1 = P_{25}$

②  $\hookrightarrow$  Second quartile  $Q_2 = P_{50} \rightsquigarrow$  Median

③  $\hookrightarrow$  Third quartile  $Q_3 = P_{75}$





Example → Given the data → 12, 18, 16, 15, 20, 10

Find quartiles →

Step 1 arrange → 10, 12, 15, 16, 18, 20

1 → Q<sub>1</sub> → P<sub>25</sub> =  $i = \frac{25}{100} \times 6 = 1.5$  not integer  
 2 ← نقطة ← Round up

Q<sub>1</sub> = 12

2 → Q<sub>2</sub> = Median "P<sub>50</sub>" =  $i = \frac{50}{100} \times 6 = 3$  integer

Q<sub>2</sub> =  $\frac{15 + 16}{2} = 15.5$

3 → Q<sub>3</sub> = P<sub>75</sub> →  $i = \frac{75}{100} \times 6 = 4.5$  not integer  
نقطة → 5

Q<sub>3</sub> = 18

يعني ثلاثة ارباع البيانات اقل او تساوي 18 or  
 75% من البيانات اقل او تساوي 18

3:2 Measures of variability "مقياس لبيان"

- 1] Range
- 2] Interquartile range "IQR"
- 3] Variance
- 4] Standard deviation
- 5] Coefficient of variation "CV"

These numbers measure the spread "dispersion"  
هذه الأرقام تقيس التشتت

لأنه يعني مثل هذا تكون علاقاته الكلاسيكية قريبة من بعضها

[1] Range "المدى"

↳ The simplest measure of variability.  
لأنه أبسط مقياس للتغير

Range = Largest value - Smallest value  
or  
Maximum - Minimum

هم كبيراً جداً بمقياس عشوائى  
لذا البيانات قريبة جداً  
من بعضها

Example → 73, 78, 81, 81, 85, 90, 90, 93, 94

\* Range is highly influenced by outliers

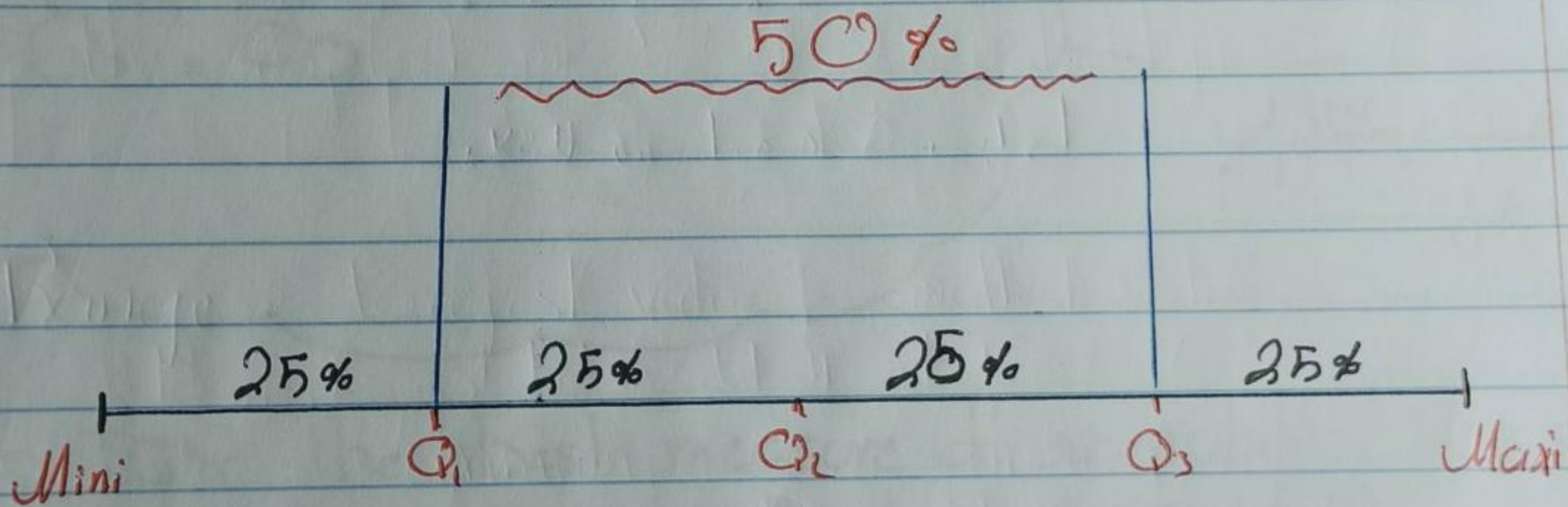
له مدى بنائي كبيراً بال Outlier ← القيم الشاذة

$$\begin{aligned} \text{Range} &= \text{Max} - \text{Mini} \\ &= 94 - 73 \\ &= \underline{\underline{21}} \end{aligned}$$

[2] Interquartile range "IQR" مدى ربيعي

↳ The range of the middle [50%] of the data  
له مدى ربيعي منتصف البيانات

$$\text{IQR} = Q_3 - Q_1$$



لا توجد مدى لهذا البيانات كما هو موضح بالرسم

لا Range اعلى من IQR كفضية

لا Range, IQR لا يمكن ان يكون سالبه  
ويجوز ان يكون كغيره فالاشارة لزيادة ج.أ. وقليلة

Example  $\rightarrow$   $Q_1 = 81$   
 $Q_3 = 90$

Find IQR  $\rightarrow$   $IQR = Q_3 - Q_1$   
 $= 90 - 81$   
 $= 9$

توزيع

\* هذا مثال الـ Range



$Q_1 = P_{25} = \frac{25}{100} \times 9$   
 $= 2.25 \rightsquigarrow 3$

$Q_1 = 81$

$Q_3 = P_{75} = \frac{75}{100} \times 9$   
 $= 6.75 \rightsquigarrow 7$

$Q_3 = 90$

من أجل المثال،  
 قبل

□ IQR The best measure of variability to use when extreme value "outlier" exist

يعني له "IQR" هو أفضل مقياس عند وجود قيم خارجة "outlier" لأن الـ IQR أقل تأثرًا بالقيم الخارجة

له ذلك فإن الـ IQR يدل على قلة الخرج عن المجتمع

[3] Variance التباين

↳ a measure that utilize all the data.  
 له وهو مقياس يستخدم في جميع البيانات

↳ Variance is the average of square deviation  
 له التباين هو متوسط الانحراف التربيعي

Population variance →  $\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$   
 التباين للمجتمع  $\mu = \text{pop size}$

Sample variance →  $s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$   
 التباين للعينة  $\bar{x} = \text{pop mean}$   
 $n = \text{samp size}$

\*we divided by  $n-1$  to have the best estimation that is unbiased "تقدير" for the population

له لأنه اقوائيه على ان يجادها وتطبيقها على  
 الالة اكاينة ا رى ا

↳  $s^2$  is a point estimator to  $\sigma^2$

note →  $\sum_{i=1}^n (x_i - \bar{x}) = 0$   
 (deviation about the mean  $x_i - \bar{x}$ )

Example → Given the sample "مثال معطى"  $15, 13, 17, 20, 22, 42$

$x_1, x_2, x_3, x_4, x_5, x_6$   
 $15, 13, 17, 20, 22, 42$

Find the variance " $s^2$ " " $n=6$ "

$\bar{x} = \frac{\text{مجموع القيم}}{\text{عدد الن}} = 21.5$

| $x_i$ | $x_i - \bar{x}$    | $(x_i - \bar{x})^2$ |
|-------|--------------------|---------------------|
| 13    | $13 - 21.5 = -8.5$ | 72.25               |
| 15    | " " = -6.5         | 42.25               |
| 17    | " " = -4.5         | 20.25               |
| 20    | " " = -1.5         | 2.25                |
| 22    | " " = 0.5          | 0.25                |
| 42    | " " = 20.5         | 420.25              |
| Total | 0                  | 557.5               |

$\sum x_i - n\bar{x} = 0 \leftrightarrow$  دائماً

هكذا نتأكد  
 بالتحقق على صيغة  
 وركابته

$$s^2 = \frac{\sum_{i=1}^6 (x_i - \bar{x})^2}{6 - 1}$$

$$s^2 = \frac{557.5}{5}$$

$$= 111.5$$

#### [4] Standard deviation "التباين المعياري"

↳ is the positive square root of the variance  
 له جذره هو الجذر التربيعي للتباين

$$\text{variance} = (\text{Standard devi.})^2$$

$$\text{Standard deviation} = \sqrt{\text{variance}}$$

↳ سبب الجذر التربيعي قبل

$$\text{variance} = 111.5$$

$$\text{S. devi} = \sqrt{111.5} \approx 10.56$$

↳ population standard deviation  $\sigma$

↳ sample standard deviation  $s$

$s$  is a point estimator for  $\sigma$

↳ Standard deviation and variance measure whether the data values cluster around the mean

له قيم للبيانات منفصلة حول المتوسط

↳ The variance is useful in comparing the variability of two or more variables

له التباين مفيد في مقارنة تباين متغيرين أو أكثر

□ we can use the SD mode to find standard deviation and variance

Step

[1] Reset all نحل الذاكرة

كيف تستخدم  
الالة في  
حسب الجان

[2] Mode → 2 "SD"

[3] Data نقل البيانات

15 → m+

13 → m+

17 → m+

20 → m+

22 → m+

42 → m+ → n = 6

[4] press shift + 2 → "Σ<sub>x</sub>"

|           |          |            |
|-----------|----------|------------|
| $\bar{x}$ | $\sum x$ | $\sum x^2$ |
| 1         | 2        | 3          |

[5] press "3" → then press "=" 10.5593 ≈ 10.56

في حالة طلب بعد ما نود  $\bar{x}$  (المتوسط) من الة نفور

[1] press shift + 2 → then press [1]

[2] press "=" → 21.5



## Notes

\* IQR ← أقل انحراف تأثر بال Outlier

\* Variance, s. deviation ← تأثر قليل بال Outlier

\* Range ← الأكثر تأثر بال Outlier

\* كلما زاد Standard Deviation في البيانات ثبتت عن ولا فرق والعكس صحيح

↑ البعد →  $s^2$  ↑

↓ البعد →  $s^2$  ↓

[5] Coefficient of variation معامل الاختلاف

↳ it ~~mean~~ measure the standard deviation relative to the mean

↳ يقاس الاختلاف المعياري بالنسبة للمتوسط

$$CV = \left( \frac{s}{\bar{x}} \times 100 \right) \%$$

Note → CV is used to Compare variability of two groups of data

↳ معامل الاختلاف يستخدم لمقارنة التباين بين مجموعتين

Example →

sample 1 → 10, 13, 12, 17, 18

sample 2 → 11, 16, 12, 9, 14

Sample 1

$$\text{mean} = \frac{\sum x_i}{n} = \frac{70}{5} = 14$$

Sample 2

$$= \frac{62}{5} = 12.4$$

\* إذا كان الـ mean متشابه للمجموعتين  
 نأخذ المجموعة ذات الـ s الأعلى

[2]  $s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$  "معامل التباين" = 3.39

= 2.7

[3] CV =  $\left( \frac{3.39}{14} \times 100 \right) \%$   
 = 24.2%

=  $\left( \frac{2.7}{12.4} \times 100 \right) \%$   
 = 21.8%

↳ sample 1 has higher variability

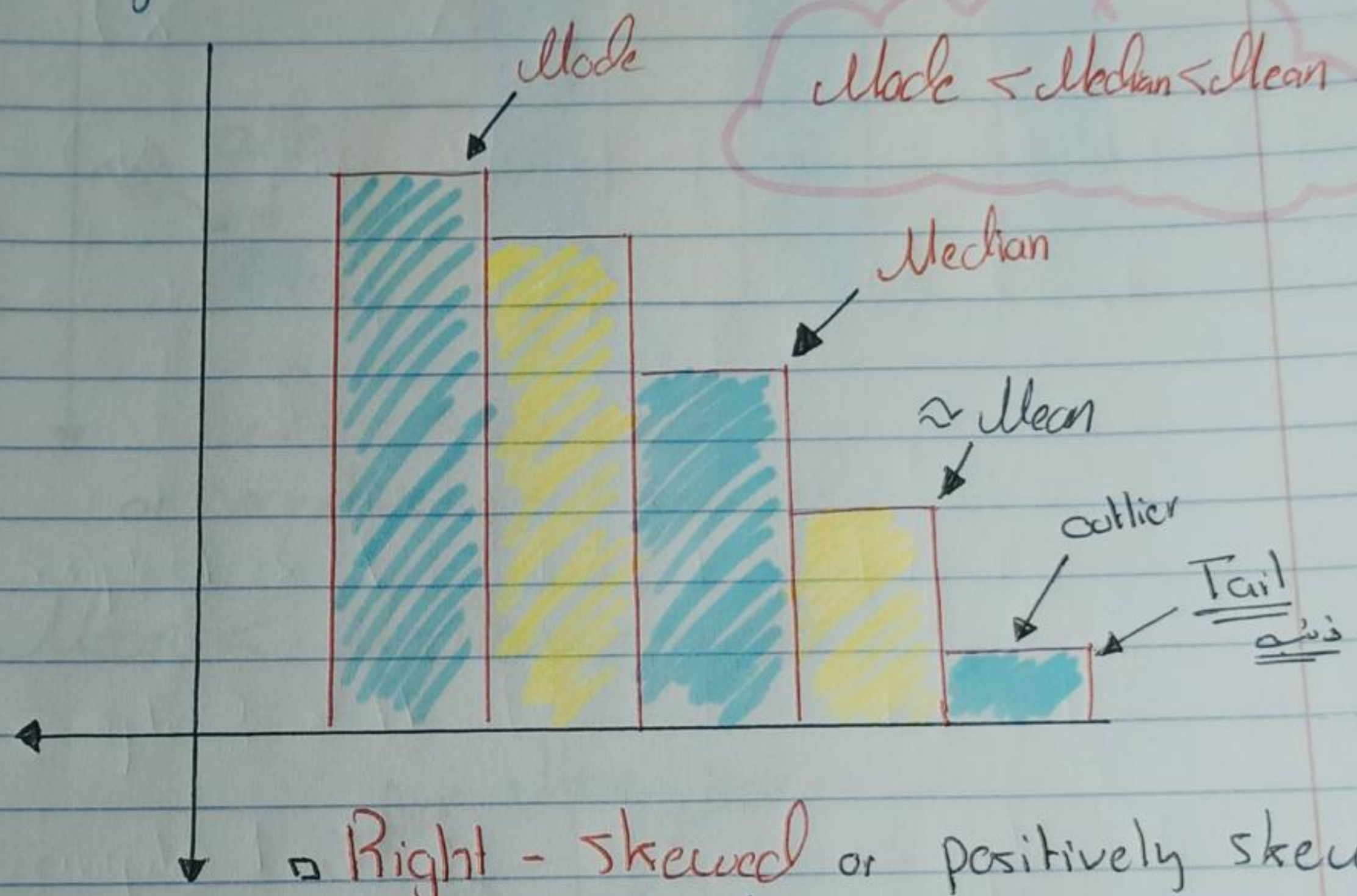
\* فنأخذ القيمة الأعلى من بين المعيارين "يعني ذات التباين الأعلى"

### 3.3 Measure of Distribution shape Relative location

□ Distribution shape

↳ Histogram

①



□ Right - skewed or positively skewed  
مانگ لایهون

# موقع اور "Tail" کے رسم ال Histogram

# Outliers ← کہناک شرطین ① تکرار قلیل

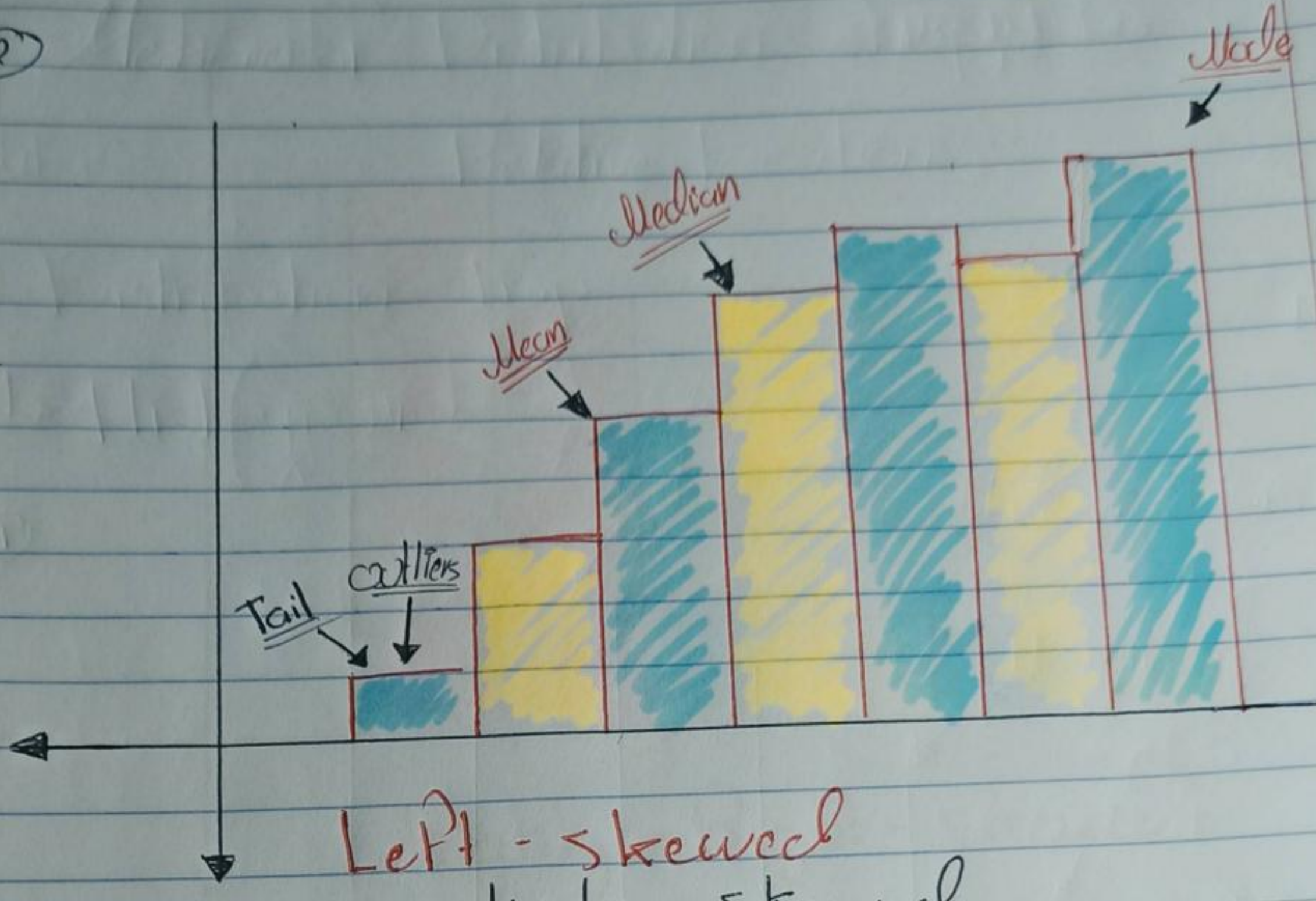
② قیمت عالیہ کبراً  
اوندافلت کبراً

# Mode ← کہو اقبیة الاعلی تکراراً

# Mean ← بتاثر بال Outliers

# Outliers ← Mean بانجا مہا

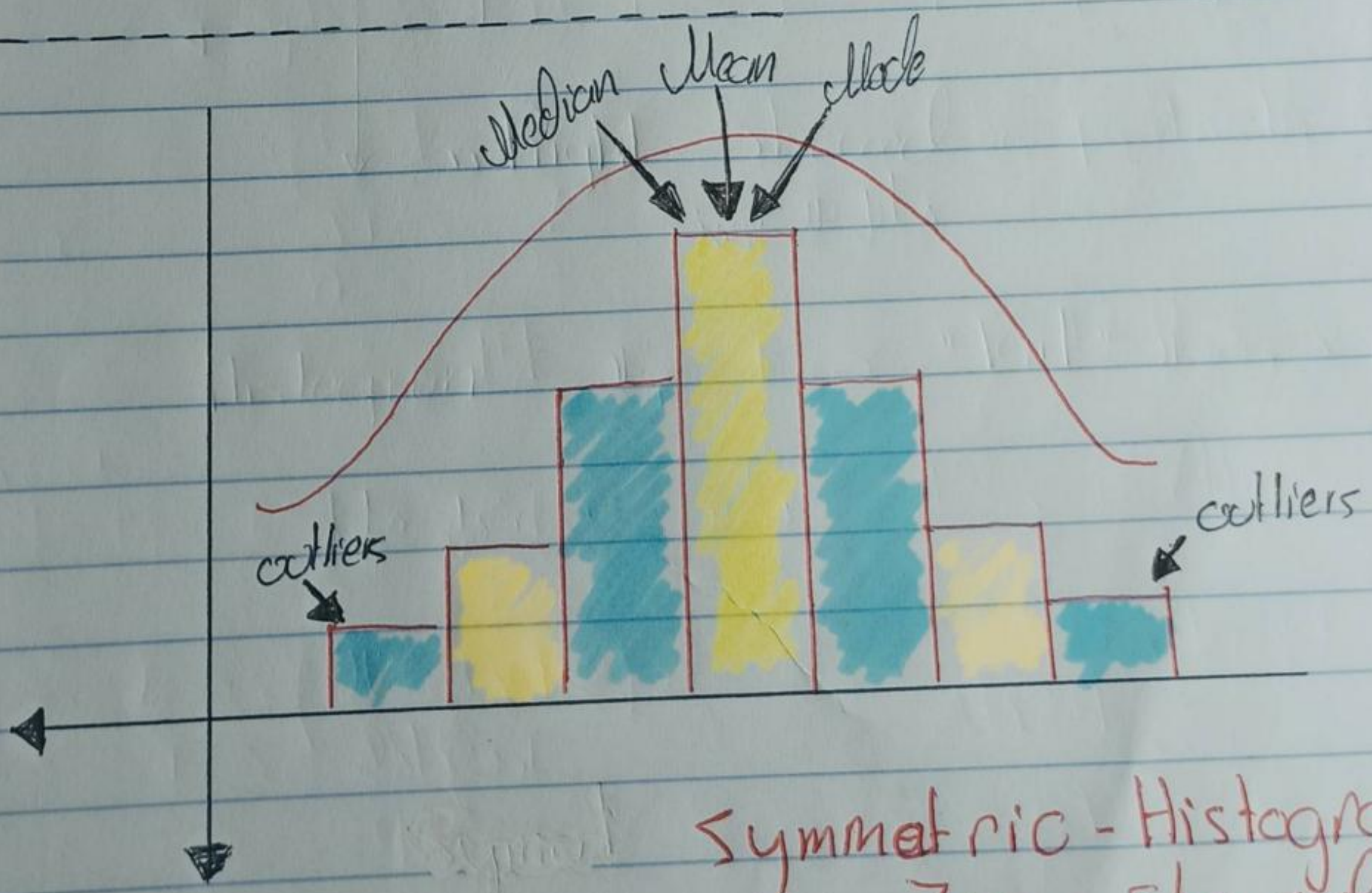
2



Left-skewed  
or negatively-skewed

$$\text{Mean} < \text{Median} < \text{Mode}$$

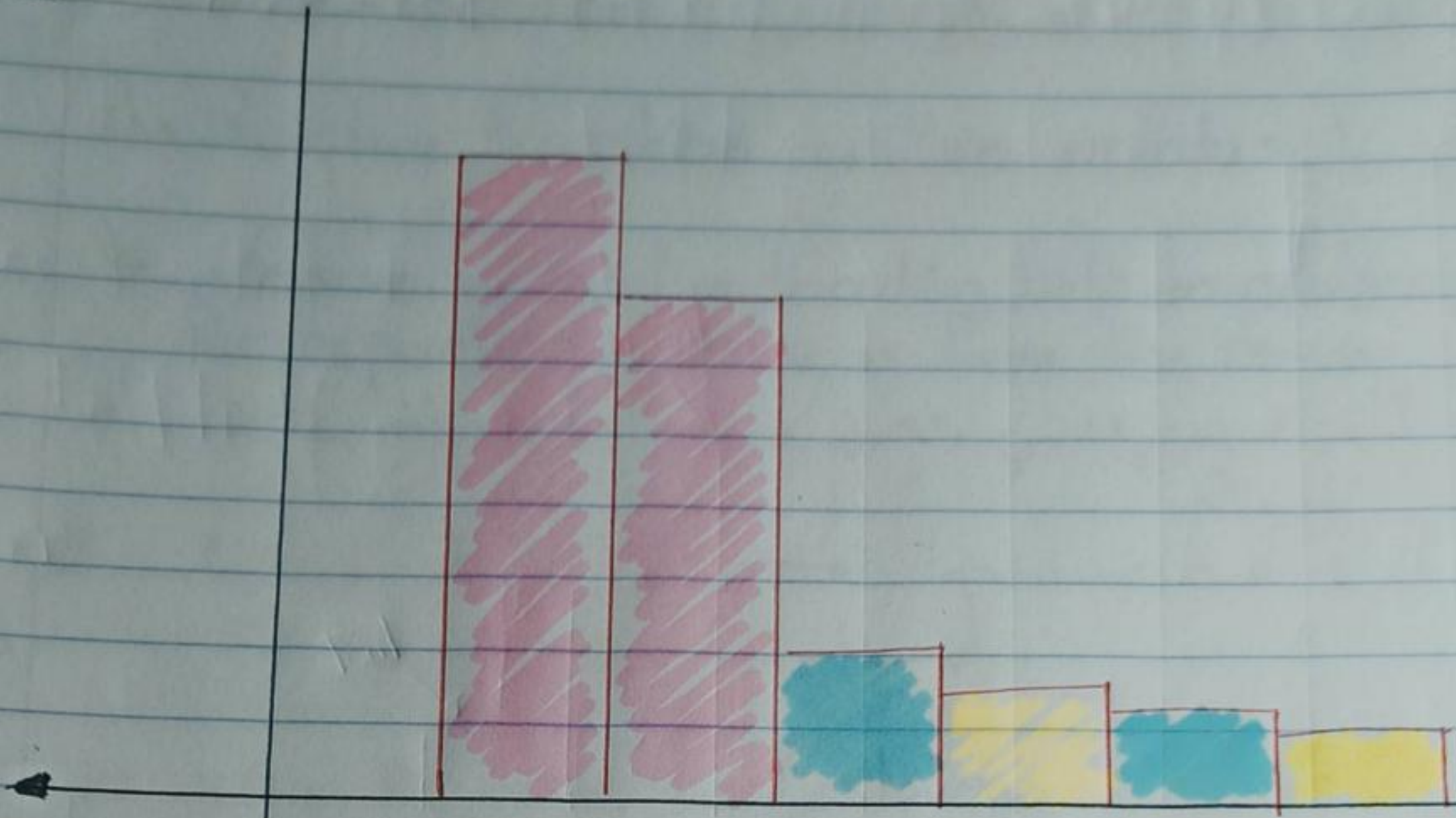
3



Symmetric - Histogram  
Zero-skewed

$$\text{Median} = \text{Mean} = \text{Mode}$$

4



Highly - Right - skewed

الفرقائه كبيرة جداً

□ Relative - location موقع نسبي

↳ Relative location within a dataset

↳ Measure of relative location help to determine how far a particular value is from the mean

لأنه يحدد الموقع النسبي بتحديد مدى بعد قيمة معينة عن المتوسط

بواسطة المقادير  $\rightarrow$  Z - SCORE  $\rightarrow$  قيمة معيارية

لأنه يقيس الموقع النسبي

□ using the mean and standard deviation we can determine the relative location of any observation

لأنه باستخدام المتوسط والانحراف المعياري يمكننا تحديد الموقع النسبي لأي ملاحظة

□ if we have values  $x_1, x_2, x_3, \dots, x_n$

□ and assume that the sample mean is  $\bar{x}$

□ and the sample standard deviation is  $s$

Remember that  $s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$

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

↳ We can compute a new value called **Z-SCORE** for each value we can compute Z-score

□ Z-score

$$Z_i = \frac{x_i - \bar{x}}{s}$$

$Z_i$  = z-score of  $x_i$

$\bar{x}$  = The sample mean

$s$  = sample standard deviation

↳ The z-score is called the standardized value Bo

Explanation of z-score

↳ is the number of standard deviation,  $x_i$  is from the mean  $\bar{x}$

EX  $\rightarrow$   $s = 10$   
 $\bar{x} = 80$

$x_1 = 70$     $x_2 = 85$     $x_3 = 95$

$Z_1 = \frac{x_1 - \bar{x}}{s} = \frac{70 - 80}{10} = -1 \rightarrow x_1 = 70$  is -1 standard deviation less than the mean  $= \bar{x} = 80$

$Z_2 = \frac{85 - 80}{10} = 0.5 \rightarrow x_2 = 85$  is 0.5 standard deviation more than the mean

$Z_3 = \frac{95 - 80}{10} = 1.5 \rightarrow x_3 = 95$  is 1.5 standard deviation more than mean

mean is said to be sign position  $\leftarrow$  z-score \*

Ex  $\rightarrow$  when Z-score = 0

this means that the value of  $x$  is the same as value of mean

$$\bar{x} = 80 \quad s = 10 \quad x_4 = 80$$

$$Z_4 = \frac{80 - 80}{10} = \underline{0}$$

also, Z-score for any observation can be interpreted as a measure of the relative location of the observation in the data set

أيضاً يمكن تفسير ~~هذه~~ مقياس لمدى انحراف كل ملاحظة عن المتوسط  
مقياس الموقع النسبي للملاحظة في مجموعة البيانات

|                       |                |                |
|-----------------------|----------------|----------------|
| Example $\rightarrow$ | $\bar{x} = 80$ | $\bar{y} = 20$ |
|                       | $s = 10$       | $s = 2$        |
|                       | $x_1 = 100$    | $y_1 = 24$     |

$x_1 = 100$  and  $y_1 = 24$  have same relative location

$$Z_{x_1} = \frac{100 - 80}{10} = 2$$

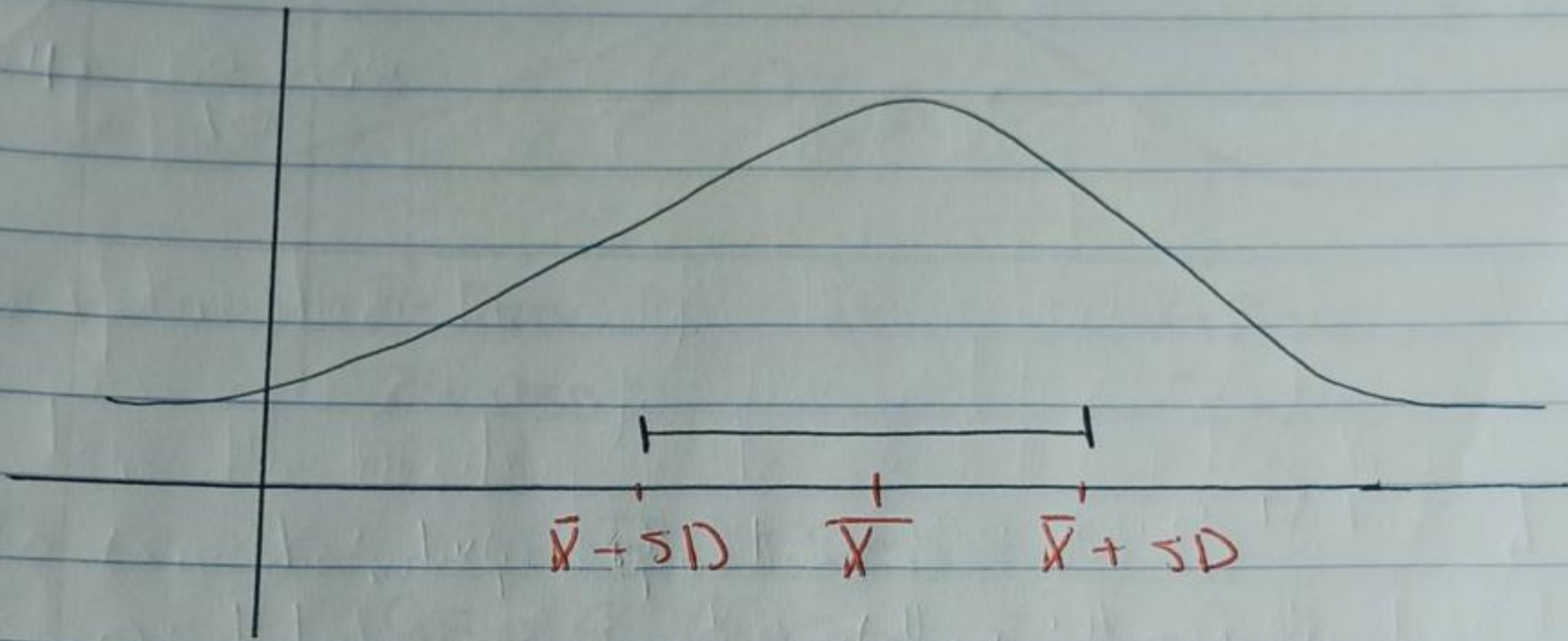
$$Z_{y_1} = \frac{24 - 20}{2} = 2$$



### Empirical Rule قاعدة التجريبية

↳ if the data set exhibits, a symmetric shaped or bell-shaped distribution like

↳ اننا عرفت في وحدة لبيانها شكل متماثل او توزيع على شكل جرس متماثل

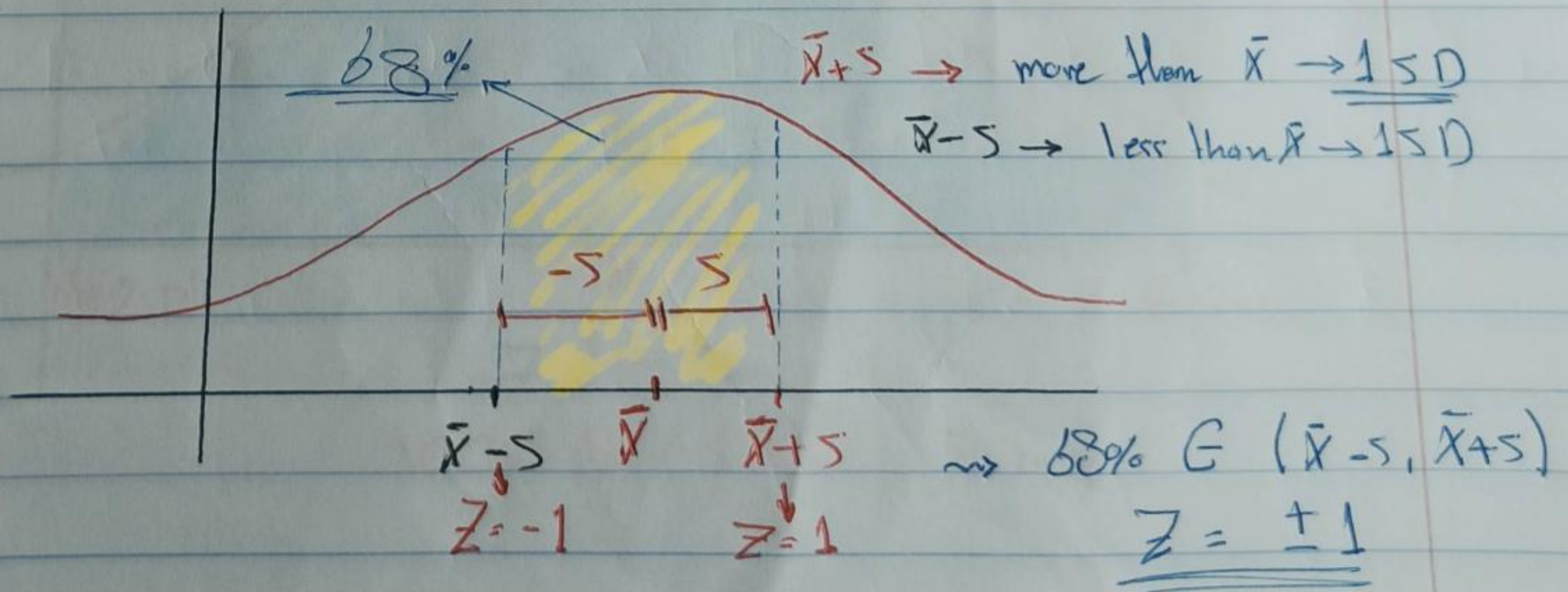


Then we can apply the empirical Rule.

↳ Empirical rule for data having a Bell-shape distribution

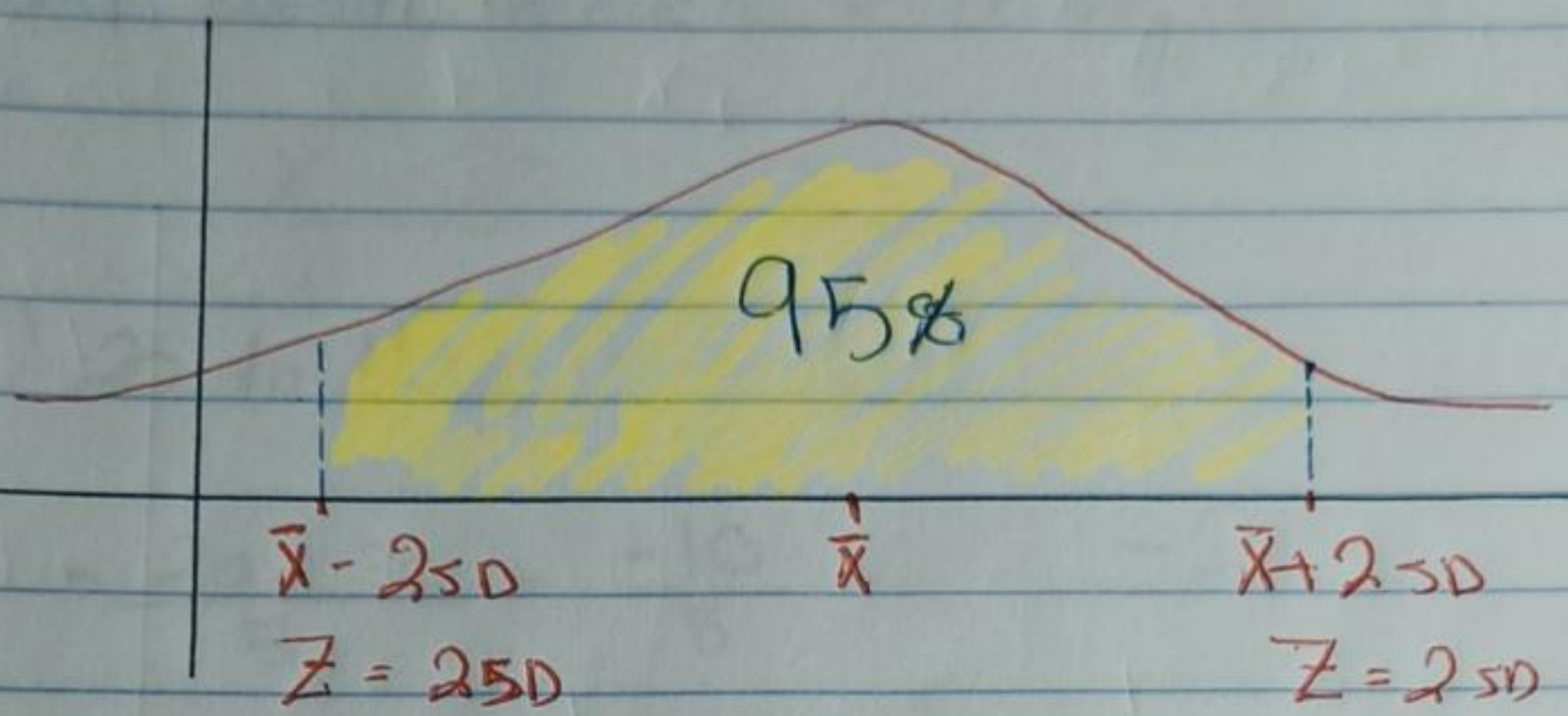
↳ approximately 68% of the data values will be within 1 standard deviation of the mean

↳ تقريباً 68% من البيانات من  $\pm 1$  انحراف معياري  
↳ بين  $z = -1$  and  $z = 1$  من البيانات



□ Approximately 95% of the data values will be within 2 standard deviation of the mean

له 95% من البيانات في 2 انحراف معياري

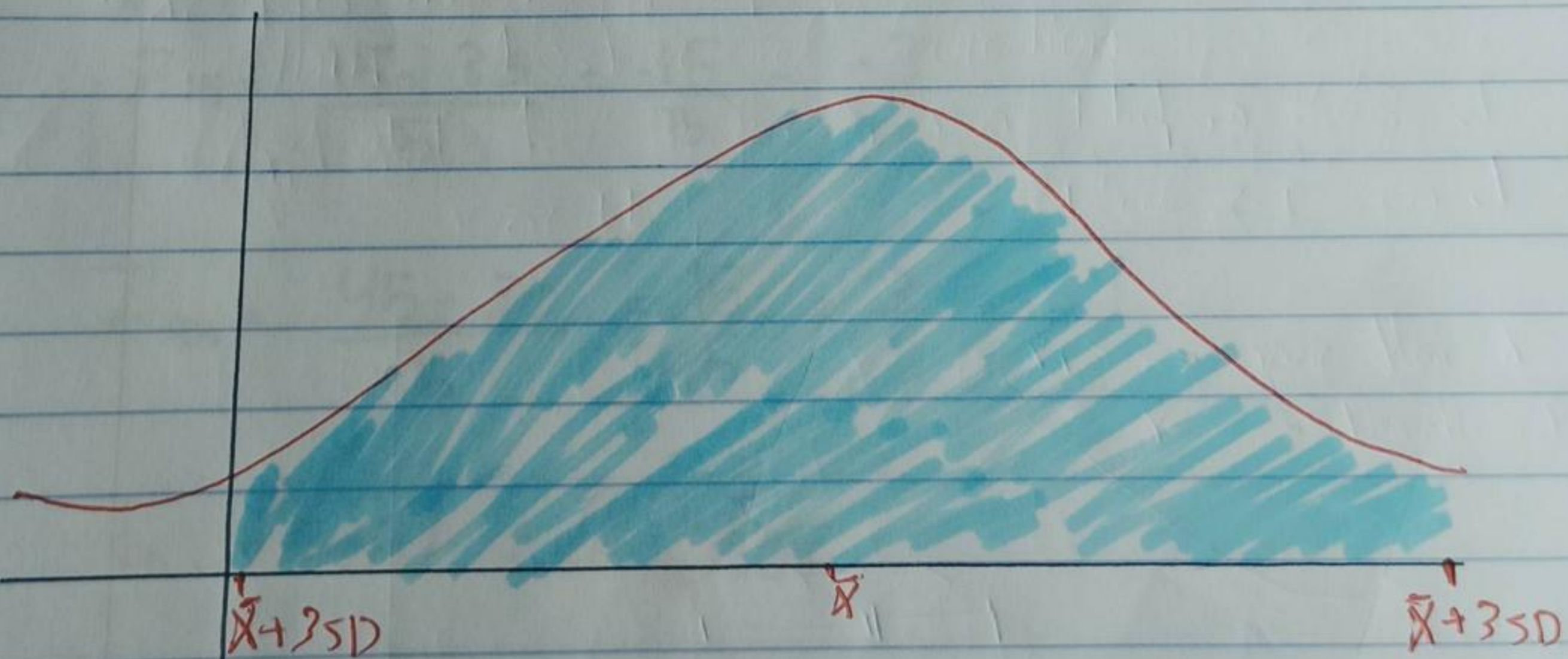


$$Z = \pm 2.50$$

$$95\% \in (\bar{x} + 2SD, \bar{x} - 2SD)$$

□ Almost all of the data values will be within 3 standard deviations of the mean

له تقريبا جميع البيانات موجودة في 3 انحراف معياري



$$Z = \pm 3SD$$

$$100\% \in (\bar{x} + 3SD, \bar{x} - 3SD)$$

Ex suppose that the data have a bell-shaped distribution with a mean of 30 and standard deviation of 5

use the empirical Rule to determine the percentage of data within of each following

$$\bar{x} = 30 \quad s = 5$$

A) 20 to 40

$$Z_{20} = \frac{20 - 30}{5} = \frac{-10}{5} = -2$$

$$Z_{40} = \frac{40 - 30}{5} = \frac{10}{5} = 2$$

↳ within 2 standard deviation, we have 95% of data

B) 15 to 45

$$Z_{15} = \frac{15 - 30}{5} = \frac{-15}{5} = -3$$

$$Z_{45} = \frac{45 - 30}{5} = \frac{15}{5} = 3$$

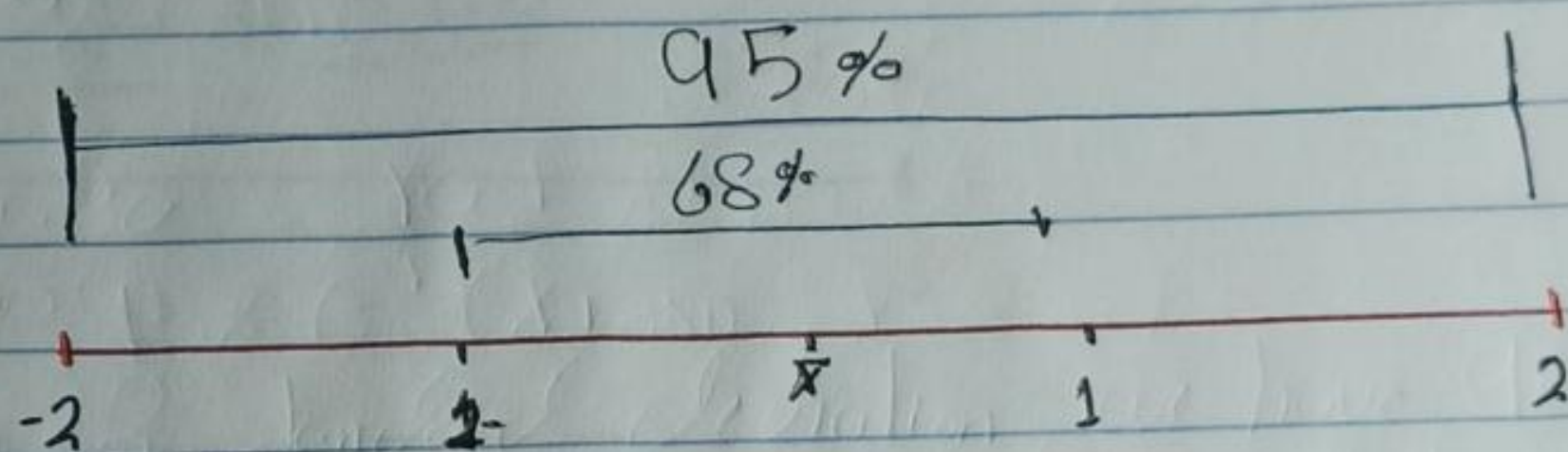
↳ Almost all data

Ex → Assume that the grades of stat 236 are Bell-shaped with  
mean = 60 and standard deviation = 10  
 $\mu = 60$        $\sigma = 10$

① Find percentage of data that are between 50 and 80

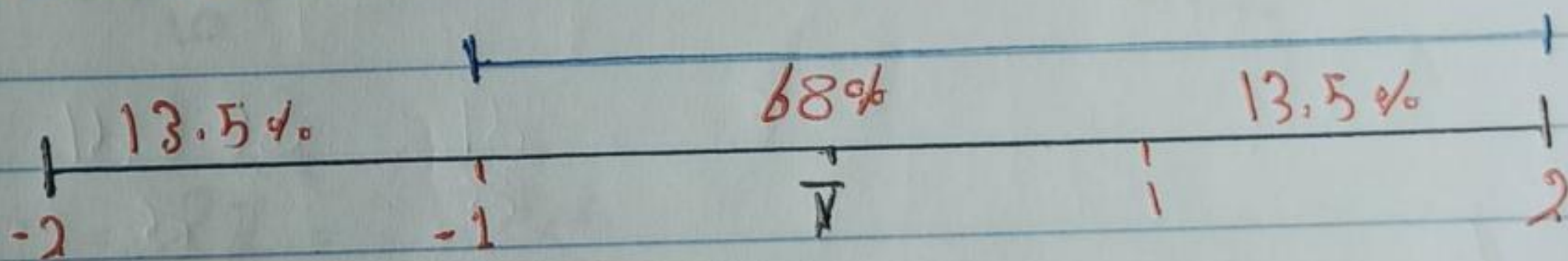
$$Z_{50} = \frac{50 - 60}{10} = \frac{-10}{10} = \underline{-1}$$

$$Z_{80} = \frac{80 - 60}{10} = \frac{20}{10} = \underline{2}$$



$$95 - 68 = 27$$

$$\frac{27}{2} = 13.5$$



So between -1 and 2 = 68% + 13.5% = 81.5%

(B) between 40 and 60

$$Z_{40} = \frac{40 - 60}{10} = \frac{-20}{10} = -2$$

$$Z_{60} = \frac{60 - 60}{10} = \frac{0}{10} = \underline{0}$$

So we have

↳ 47.5% of grades are between 40 and 60

$$\frac{95}{2} = \underline{47.5}$$

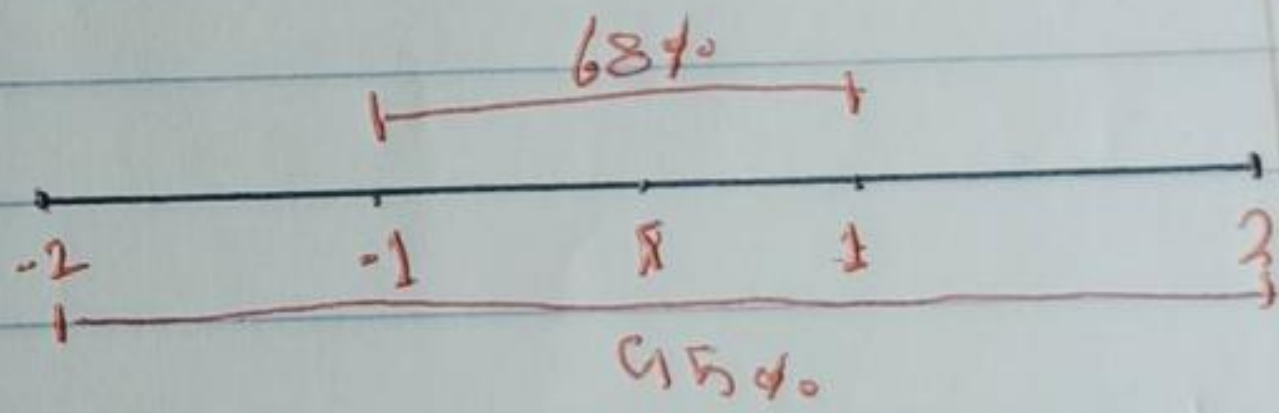
95%



(C) between 70 and 40

$$Z_{70} = \frac{70 - 60}{10} = 1$$

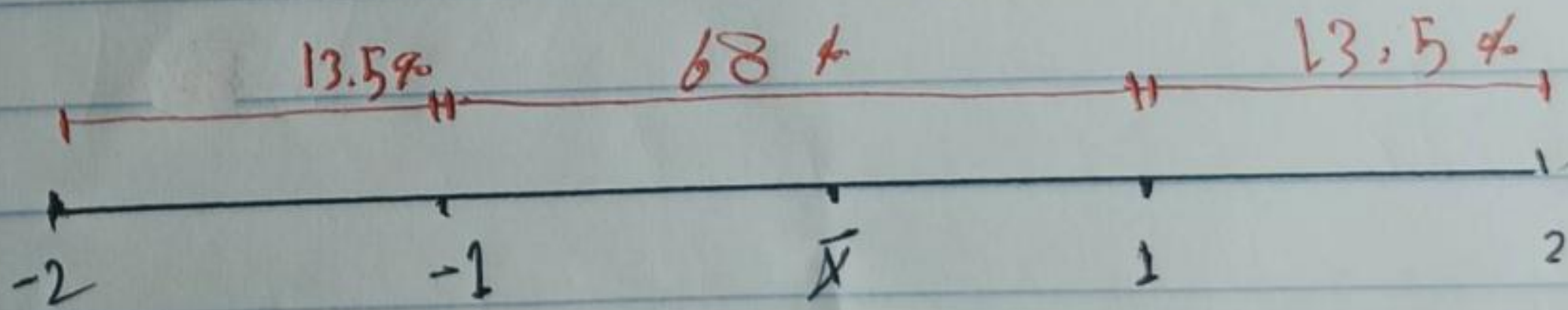
$$Z_{40} = \frac{40 - 60}{10} = -2$$



$$95 - 68 = \frac{27}{2} = 13.5\%$$

So between 1 and -2

$$= 68 + 13.5 = 81.5\%$$



Detecting outliers  $\Rightarrow$  Bell-shaped

30  
|| 80 ||  
10

$\hookrightarrow$  Extreme values are called also outliers

if  $Z > 3$  or  $Z < -3$

then this  $Z$ -score belong to outliers

Example  $\rightarrow \bar{x} = 80$

$s = 5$

$Z_{98}$ ?

$$Z_{98} = \frac{98 - 80}{5} = \frac{18}{5} = \underline{\underline{3.6 \text{ outliers}}}$$

$Z > 3$

### 3.4 || Exploratory Data Analysis $\leftrightarrow$ 5 numbers

□ 5-number summary

↳ The five following numbers are used to summarize data

- ① smallest value
- ② First Quartile  $Q_1$
- ③ Median
- ④ Third Quartile  $Q_3$
- ⑤ Largest value

□ Example  $\rightarrow$  Construct the five number summary of the following numbers

3310, 3355, 3450, 3480, 3490  
3520, 3540, 3550, 3650, 3925

□ Smallest value is  $\sim$  3310

□ First Quartile  $Q_1 = P_{25}$

↳ so we have to compute the position index  $i$

$$i = \frac{P}{100} \times n \quad \rightarrow \quad P_{25} = \frac{25}{100} \times 12 = 3$$

↳  $i$  is integer  $\rightarrow$  so the first quartile is the average of the third and fourth values

$$Q_1 = \frac{3450 + 3480}{2} = \underline{\underline{3465}}$$

### Median

↳ we have two ways to compute the median

طريقة

II we have 12 values so the median will be the average of the 6<sup>th</sup> and 7<sup>th</sup> values

$$\text{Median} = \frac{3490 + 3520}{2} = 3505$$

طريقة

$$\text{Median} = P_{50}$$

$$L = \frac{50}{100} \times 12 = \underline{6}$$

L is integer  $\Rightarrow$  so the median is the average of 6<sup>th</sup> and 7<sup>th</sup> values

$$\text{Median} = \frac{3490 + 3520}{2} = 3505$$

### Q<sub>3</sub>

$$\text{Q}_3 = P_{75}$$

$$L = \frac{75}{100} \times 12 = 9$$

L is the average of 9<sup>th</sup> and 10<sup>th</sup> values

$$\text{Q}_3 = \frac{3550 + 3650}{2} = 3600$$

### The largest value

$$\text{the largest value} = 3925$$



So the 5-number summary is

3310, 3465, 3505, 3600, 3925

Approximately 25% of the data are between any two adjacent "lag" #

ل 25% من البيانات تقع بين أي رقمين متجاورين

1. 2. 3.

" Box Plot "

Box Plot is a graphical summary of data that is based on Five Summary numbers

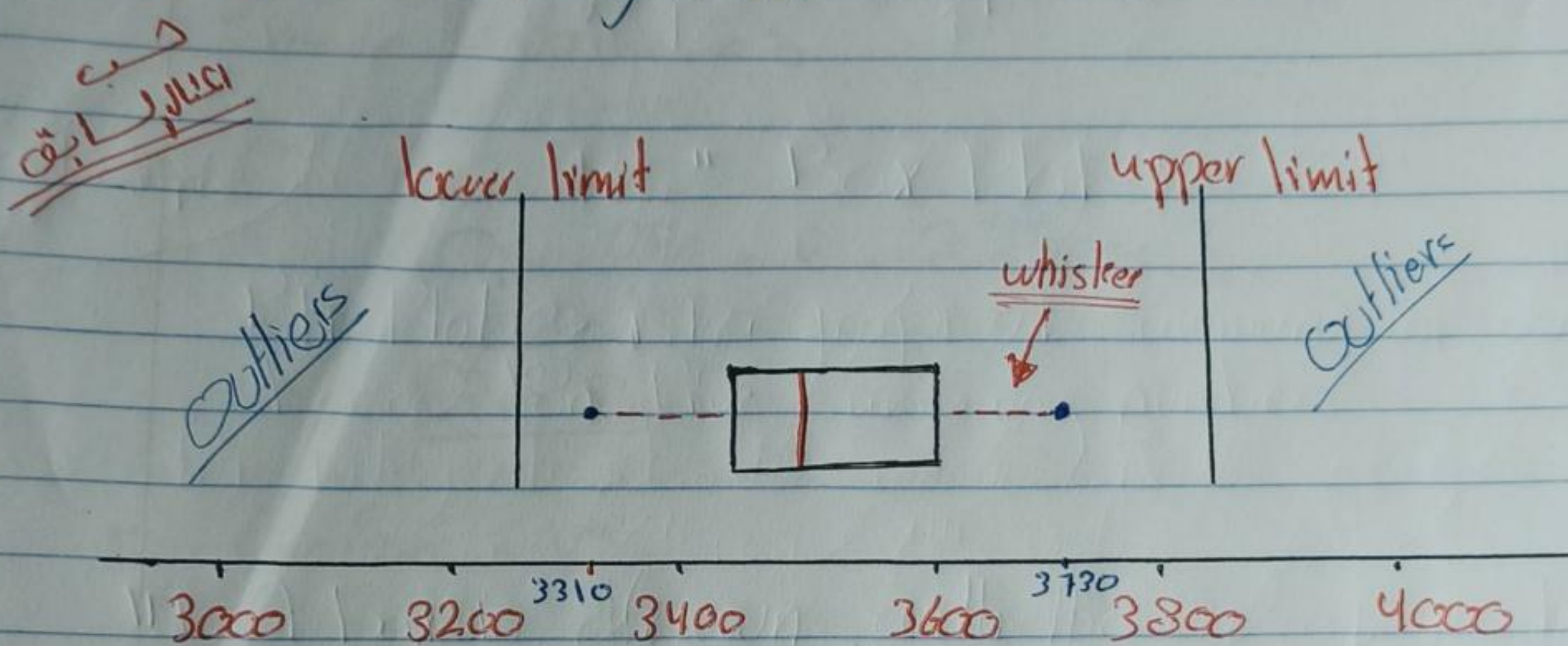
1] A box is drawn with the end of the box located at first and third quartiles  
 → The box contains 50% of data

2] A vertical line is drawn in the box at location of median

3] By using the interquartile range limits are located, the limits for the box plot are  $(1.5)(IQR)$  below  $Q_1$  and  $(1.5)(IQR)$  above  $Q_3$

- ④ dash lines are called whisker " كَبْلَة " .  
 The whiskers are drawn from the end of the box to the smallest and largest values inside the limits. Compute in step 3.

- ⑤ Finally the location of each outlier is shown with the symbol.



$$Q_1 = 3465$$

$$Q_3 = 3600$$

$$\text{median} = 3505$$

$$\text{IQR} = Q_3 - Q_1 = 3600 - 3465 = \underline{\underline{135}}$$

$$1.5 \text{ IQR} = (1.5)(135) = 202.5$$

$$Q_1 - 1.5 \text{ IQR} = 3465 - 202.5 = 3262.5 \text{ lower limit}$$

$$Q_3 + 1.5 \text{ IQR} = 3600 + 202.5 = 3802.5 \text{ upper limit}$$

Smallest # inside the limits is 3310 } → داتل كبر  
 Largest # inside the limits is 3730 }

3925 is an outlier

Example

5, 15, 18, 10, 8, 12, 16, 10, 6

Construct the Box plot

$\hookrightarrow$  5, 6, 8, 10, 10, 12, 15, 16, 18

Smallest value = 5

$Q_1 = P_{25}$

$$\hookrightarrow Q_1 = \frac{25}{100} \times 9 = 2.25 \approx \underline{\underline{3}}$$

So  $Q_1 = 8$

Median = 10

$Q_3 = P_{75}$

$$\hookrightarrow Q_3 = \frac{75}{100} \times 9 = 6.75 \approx \underline{\underline{7}}$$

$Q_3 = 15$

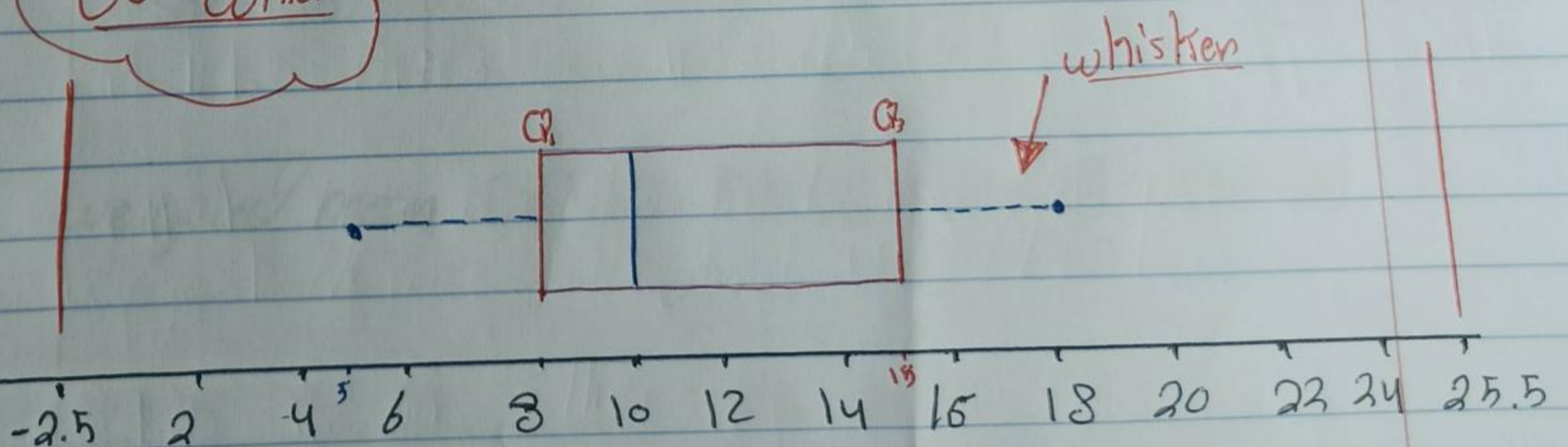
$$IQR = 15 - 8 = \underline{\underline{7}}$$

$$\sim \Rightarrow 1.5 \text{ IQR} = 1.5 \times 7 = \underline{\underline{10.5}}$$

$$\hookrightarrow \text{lower limit} = Q_1 - 10.5 = -2.5$$

$$\hookrightarrow \text{upper limit} = Q_3 + 10.5 = 25.5$$

No outlier



### 3.6 The weighted mean and working with grouped data

→ The mean of a sample

$$\bar{x} = \frac{\sum x_i}{n} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

weighted mean → عندما نريد ان نوزن معلول لقيم غير متكافئة

$$\bar{x} = \frac{\sum x_i w_i}{\sum w_i} \quad \text{مثال الطريقة هنا مع جدول كما ستعلمون}$$

$x_i$  = value of observation  $i$  → القيمة

$w_i$  = weight for observation  $i$  → وزن القيمة

Example → You have 3 Courses → لديك 3 مواد

| Courses  | Grade | # of hours |
|----------|-------|------------|
| Math 235 | 90    | 3          |
| Com 221  | 80    | 2          |
| stat 236 | 95    | 3          |

المتوسط المرجح

$$\text{weighted mean } (\bar{x}) = \frac{(90) \times (3) + (80) \times (2) + (95) \times (3)}{3 + 2 + 3}$$

$$\bar{x} = 89.375$$

Example →

Clothes

| Cost per meter<br>( $x_i$ ) | # of meter<br>( $w_i$ ) |
|-----------------------------|-------------------------|
| 3 \$                        | 20                      |
| 5 \$                        | 15                      |
| 2 \$                        | 70                      |

Compute the weighted mean (Cost per meter)

$$\bar{x} = \frac{(3 \times 20) + (5 \times 15) + (2 \times 70)}{20 + 15 + 70}$$

$$= \frac{60 + 75 + 140}{85} = \frac{275}{85}$$

$$= 3.23 \text{ $ per meter}$$

مع  
المتوسط  
الوزن

Grouped data      بيانات موزونة

| Classes | Frequency |
|---------|-----------|
| 10-14   | 4         |
| 15-19   | 8         |
| 20-24   | 5         |
| 25-29   | 2         |
| 30-34   | 1         |

sample mean for grouped data

$$\bar{x} = \frac{\sum P_i M_i}{n}$$

$M_i$  → The midpoint for class  $i$

$P_i$  → Frequency of class  $i$

$n$  → sample size

| Class | Midpoint               | $P_i$     | $M_i P_i$  |
|-------|------------------------|-----------|------------|
| 10-14 | $\frac{10+14}{2} = 12$ | 4         | 48         |
| 15-19 | $11 = 17$              | 8         | 136        |
| 20-24 | $11 = 22$              | 5         | 110        |
| 25-29 | $11 = 27$              | 2         | 54         |
| 30-34 | $11 = 32$              | 1         | 32         |
|       |                        | <u>20</u> | <u>380</u> |

$$\bar{x} = \frac{\sum P_i M_i}{n} = \frac{380}{20} = 19 \text{ days}$$

∴ Sample variance =  $\frac{\sum^2}{n-1}$   
 $\sum^2 = \frac{\sum P_i (c_{mp} - \bar{x})^2}{n-1}$

حسابات

| Classes | P <sub>i</sub> | M <sub>i</sub> | (M <sub>i</sub> - $\bar{x}$ ) | (M <sub>i</sub> - $\bar{x}$ ) <sup>2</sup> | P <sub>i</sub> (M <sub>i</sub> - $\bar{x}$ ) <sup>2</sup> |
|---------|----------------|----------------|-------------------------------|--------------------------------------------|-----------------------------------------------------------|
| 10-14   | 4              | 12             | 12-19 = -7                    | (-7) <sup>2</sup> = 49                     | 196                                                       |
| 15-19   | 8              | 17             | 17-19 = -2                    | (-2) <sup>2</sup> = 4                      | 32                                                        |
| 20-24   | 5              | 22             | 11 = 3                        | 11 = 9                                     | 45                                                        |
| 25-29   | 3              | 27             | 11 = 8                        | 64                                         | 128                                                       |
| 30-34   | 1              | 32             | 13 = 13                       | 169                                        | 169                                                       |
|         | 20             |                |                               |                                            | $\sum = 570$                                              |

19  
 $\bar{x}$   
 صيغة

$\sum^2 = \frac{570}{19} = 30$   
 $\frac{19}{20-1}$

∴ population mean For grouped data

$M = \frac{\sum P_i M_i}{n}$

معدل  
 الحسابات

∴ population variance For grouped data

$\sigma^2 = \frac{\sum P_i (c_{mp} - M)^2}{n}$

3.5 Measures of Association between two variables

↳ In this section we present descriptive measures of the relationship between two variables

Covariance

التباين

↳ For n-observations  $(x_1, y_1) (x_2, y_2)$

sample Covariance 
$$s_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

Example → we have two variables # of advertisement per month and sales for a company (monthly)

| $x$ | $y$ | $x - \bar{x}$ | $y - \bar{y}$ | $(x - \bar{x})(y - \bar{y})$ |
|-----|-----|---------------|---------------|------------------------------|
| 3   | 50  | -1            | 8             | -8                           |
| 1   | 20  | -3            | -22           | 66                           |
| 4   | 30  | 0             | -12           | 0                            |
| 5   | 60  | 1             | 18            | 18                           |
| 7   | 50  | 3             | 8             | 24                           |
|     |     |               |               | 100                          |

Compute the Covariances of the above sample

$$s_{xy} = \frac{\sum (x - \bar{x})(y - \bar{y})}{n-1}$$

Note that  
 $\sum (x - \bar{x}) = 0$   
 $\sum (y - \bar{y}) = 0$



Example 1

$$\bar{x} = \frac{3+1+4+5+7}{5} = \frac{20}{5} = 4$$

$$\bar{y} = \frac{50+20+30+60+50}{5} = \frac{210}{5} = 42$$

$$s_{xy} = \frac{100}{5-1} = \frac{100}{4} = 25 \quad \text{two relationship}$$

□ Sample Covariance  $s_{xy} = \frac{\sum (x-\bar{x})(y-\bar{y})}{n-1}$

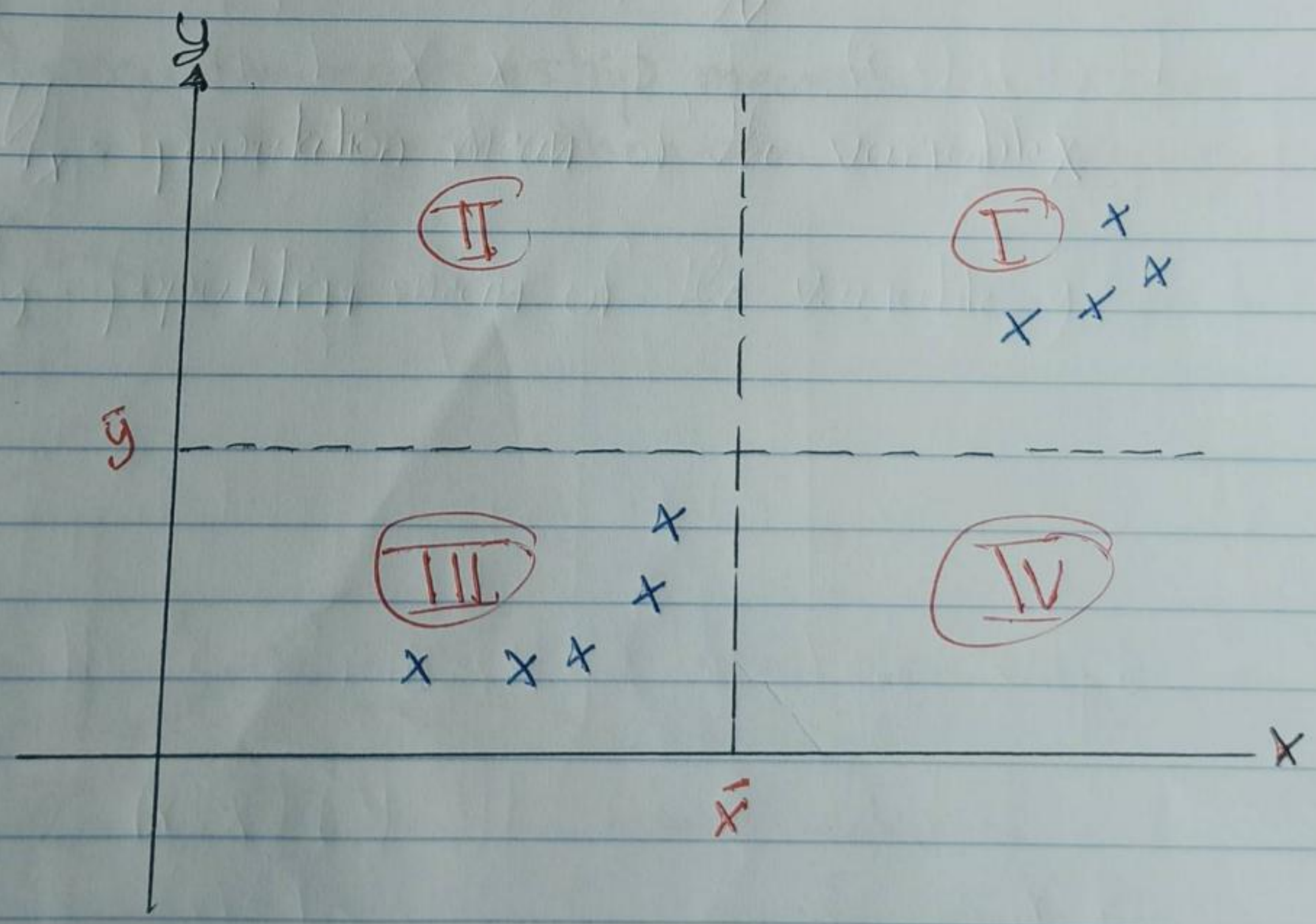
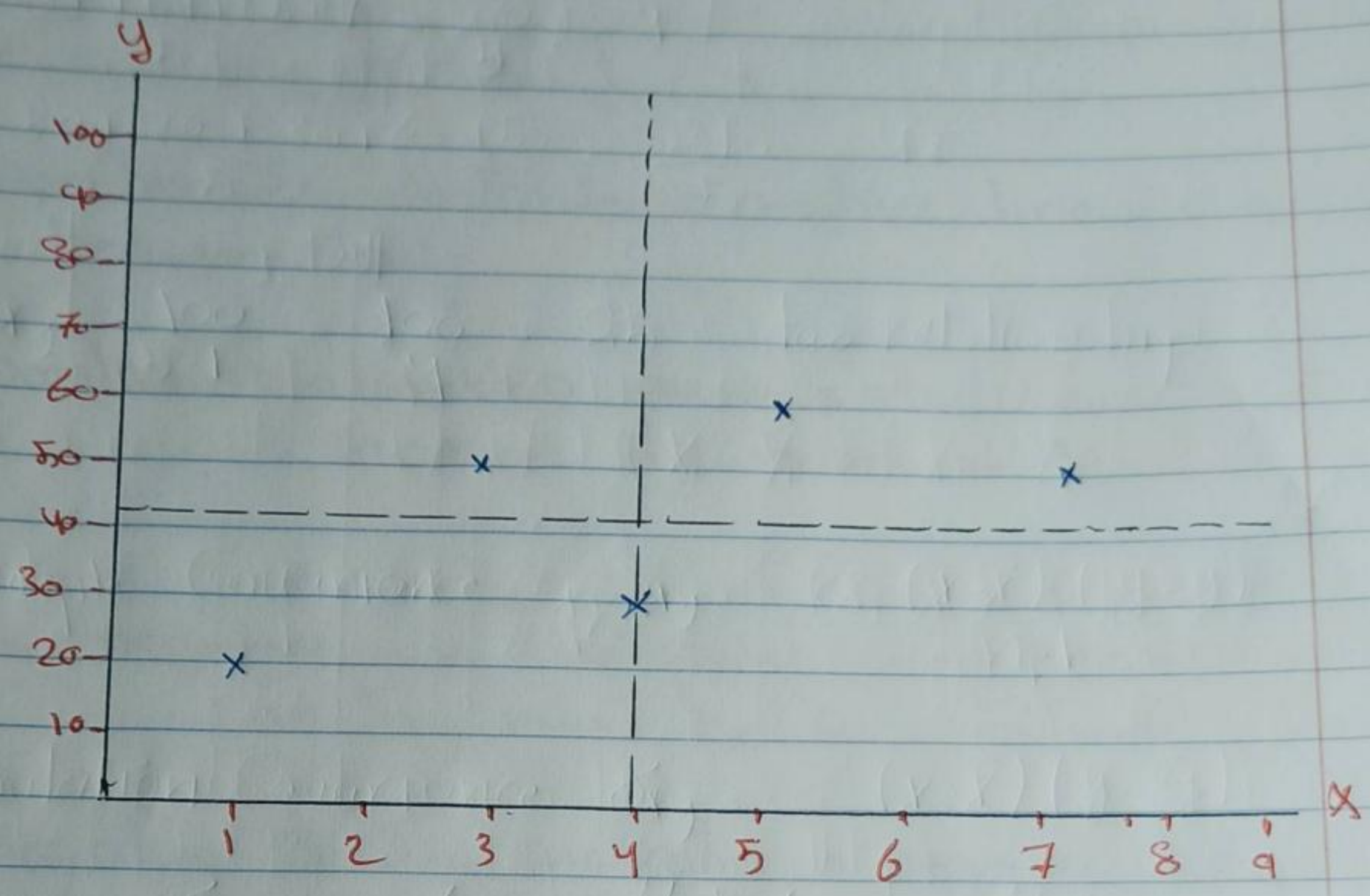
□ population Covariance  $\sigma_{xy} = \frac{\sum (x_i - \mu_x)(y_i - \mu_y)}{N}$

$$\sigma_{xy} = \frac{\sum (x_i - \mu_x)(y_i - \mu_y)}{N}$$

$\mu_x$  = population mean of the variable x

$\mu_y$  = population mean of the variable y

The Covariance is a measurer of linear Association between two variables



18/11

if  $\sum xy$  positive  $\rightarrow$  The point with greatest influence on  $\sum xy$  must be in quadrants I and III

مثال

A positive  $\sum xy$  indicates a positive linear association between  $x$  and  $y$

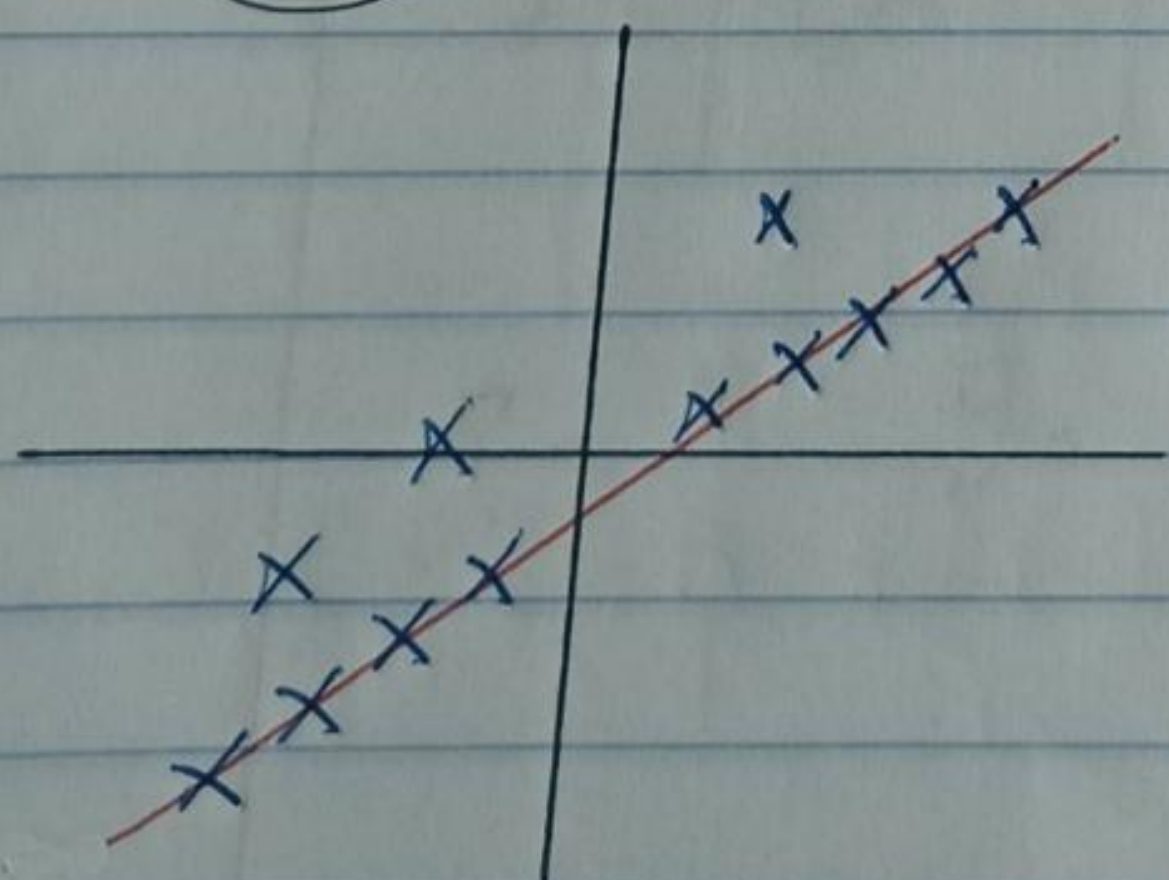
A positive relationship means  $\&$  As values of  $x$  increase the values of  $y$  also increase.

if  $\sum xy$  negative  $\rightarrow$  The point with greatest influence on  $\sum xy$  must be in quadrants II and IV

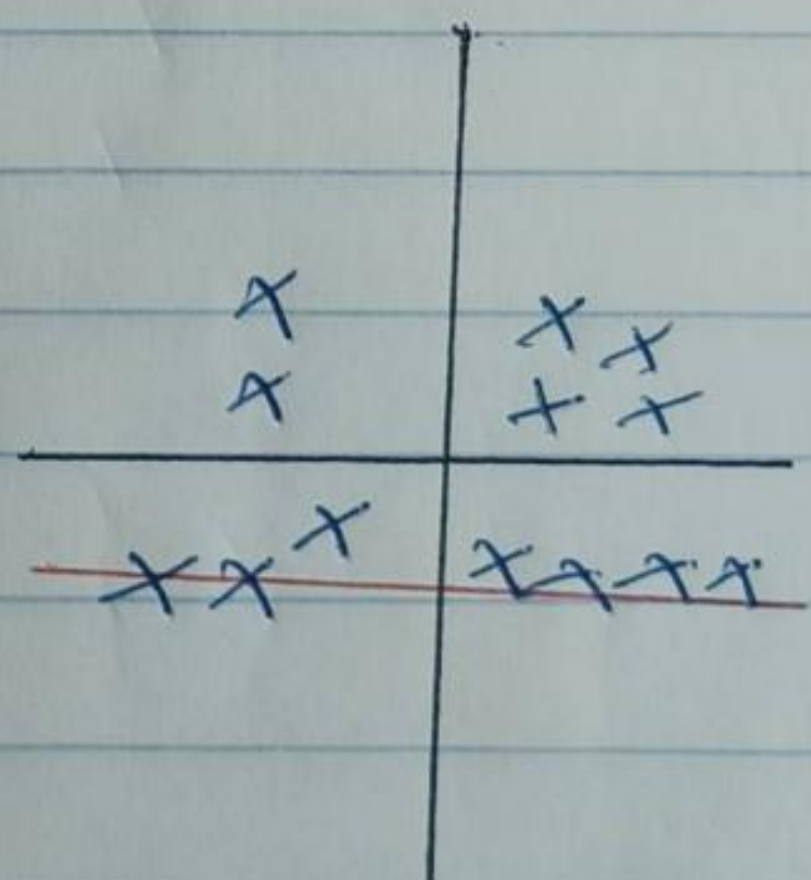
A negative  $\sum xy$  indicates a negative linear association between  $x$  and  $y$

A negative relationship means  $\&$  As values of  $x$  increase the value of  $y$  decreases

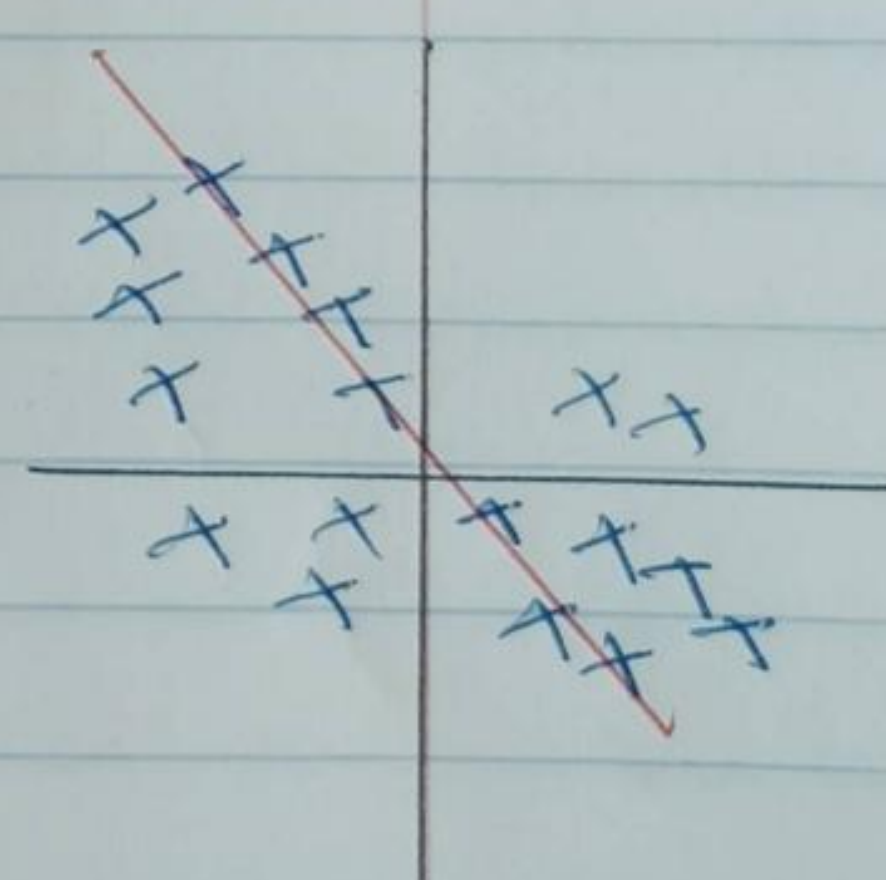
مثال



$\sum xy = +$



$\sum xy \approx 0$



$\sum xy = -$

One problem using Covariance as a measure of the strength of the linear relationship is that the value of the Covariance depends on the units of measurement of  $x$  and  $y$ .

أحد مشاكل استخدام التغاير كقياس لمتى قوة العلاقة بين متغيرين، إذ أنها تعتمد على الوحدة

A measure of the relationship between two variables that is not affected by the unit of measurement for  $x$  and  $y$  is the

مقياس لا يتأثر بالوحدة ← Correlation Coefficient  
 مقياس العلاقة كما سبق

$$r_{xy} = \frac{\sum xy}{s_x s_y}$$

| $x$ | $y$ | $x - \bar{x}$ | $y - \bar{y}$ | $(x - \bar{x})^2$ | $(y - \bar{y})^2$ |
|-----|-----|---------------|---------------|-------------------|-------------------|
| 3   | 50  | -1            | 8             | 1                 | 64                |
| 1   | 20  | -3            | -22           | 9                 | 484               |
| 4   | 30  | 0             | -12           | 0                 | 144               |
| 5   | 60  | 1             | 18            | 1                 | 324               |
| 7   | 50  | 3             | 8             | 9                 | 64                |
|     |     |               |               | 20                | 1080              |

$\bar{x} = 4$      $\bar{y} = 42$

$$s_x = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} = \sqrt{\frac{20}{4}} = \sqrt{5} \approx 2.23$$

$$s_y = \sqrt{\frac{\sum (y_i - \bar{y})^2}{n-1}} = \sqrt{\frac{1080}{4}} = \sqrt{270} = 16.43$$

$$\sum xy = \sum \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1} = 25$$

$$\text{So } r_{xy} = \frac{25}{(2.23)(16.43)} = 0.682$$

□ Correlation Coefficient of population data

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

□ Correlation Coefficient of sample data

$$r_{xy} = \frac{\sum xy}{\sum x \sum y}$$

↳ The Correlation Coefficient ranges from -1 to 1

□ values close to -1  $\rightsquigarrow$  strong negative relationship

□ values close to 1  $\rightsquigarrow$  strong positive relationship

□ values close to 0  $\rightsquigarrow$  weaker relationship

بزرگ

د کيفيتہ التطبيق على الآلة كما به

Mode, 3, =, = ← بجعل تنظيم ← shift

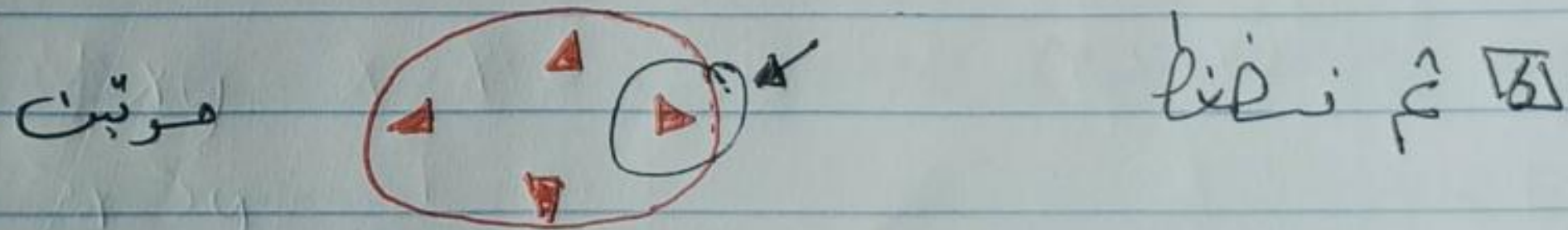
Mode بزرگ

"REG" 3 بزرگ

Linear من "Lin" 1 بزرگ

نم ندفل البيانات رسم لجدول هذه الطريقة

| x | y  | نفسه  |        |
|---|----|-------|--------|
| 3 | 50 | 3, 50 | M+ n=1 |
| 1 | 20 | 1, 20 | M+ n=2 |
| 4 | 30 | 4, 30 | M+ n=3 |
| 5 | 60 | 5, 60 | M+ n=4 |
| 7 | 50 | 7, 50 | M+ n=5 |



| A | B | r |
|---|---|---|
| 1 | 2 | 3 |

لم قائله سائله

Correlation Coefficient "r" 3 بزرگ

r<sub>xy</sub>

0.680 بزرگ