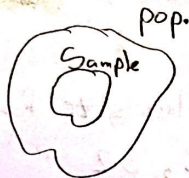


* Ch 3: Descriptive Statistics: Numerical Measures.



سواء كانت البيانات مأخوذة من عينة أو مجتمع جازء القراءات الرقمية التي نجد لها هذه البيانات تكون ضمن الإحصاء الوصفي (Descriptive statistic).

- Measures are computed for data from a sample: sample statistics.
- Measures are computed for data from a population: parameters.
- A sample statistic is the point estimator of the corresponding population parameter. } inferential statistic هنا تكون ضمن

أي قراءاتها بلغة population تسمى parameter
 أي قراءتها بلغة sample تسمى statistic

عند صعوبة دراسة المجتمع الكامل «pop» نلجأ إلى عينة (sample) وتعتبر قراءات العينة نقاط تقريبية لقراءات المجتمع كإجمالي.

- Sec 3.1: Measures of location:-

Measures of central tendency
 مقاييس الموضع أو مقاييس النزعة المركزية

مقاييس تستخدم لقياس موضع تركيز أو تجمع البيانات.

1 Mean: «Arithmetic average» الوسط الحسابي

- If the data are for a sample, the mean is denoted by \bar{X} .
- If the data are for a population, the mean is denoted by M .

→ the sample mean $\bar{X} = \frac{\sum X_i}{n}$; \sum : the summation sign التجميع ; X_i : the data values, n : the sample size.

→ the population mean $M = \frac{\sum X_i}{N}$; N : the population size.

المتوسط الحسابي = مجموع القيم / العدد الكلي

-Ex: The following data represents the class size for a sample of 5 classes. 46, 54, 42, 46, 32

Find the sample mean.

$$\rightarrow \bar{X} = \frac{46 + 54 + 42 + 46 + 32}{5} = 44$$

-Note: \bar{X} is a point estimator for M .

2) Median: الوسيط
The value in the middle.

To find the median:

- a- order the data (from smallest to largest) ^{أو الصغرى إلى الكبرى}.
- b- If the size (n) is odd, the median is the middle value.
- If the size (n) is even, the median is the average of of two middle values.

-Ex: Given the sample: 46, 54, 42, 46, 32. Find the median.

$n = 5$ odd

ترتيب البيانات \rightarrow 32, 42, 46, 46, 54

the median is 46

بما كانتنا بعد ترتيب البيانات، إيجاد رتبة الوسيط بالطريقة التالية إذا n فردية $\frac{n+1}{2}$ ونبحث عنه القيمة

- Ex: Given the sample: 13, 10, 15, 10, 10, 9, 12, 14,
Find the median.

→ $n = 8$ (even)

ترتيب البيانات 9, 10, 10, 10, 12, 13, 14, 15

the median = $\frac{10 + 12}{2} = 11$

the rank of the median = 4, 5

* Note that: 50% of data value are less than or equal the median, and 50% are greater than or equal the median.

بما انك لا تجد
رتبة الوسيط اذا
كان العدد زوجي:-
 $\frac{n}{2}, \frac{n}{2} + 1$

* Extreme values: outliers القيم المتطرفة
أو الشاذة تكون بعيدة عن البيانات.

- Ex: ① Given the sample: 4, 2, 1, 5, 4

→ 1, 2, 4, 4, 5

the median is 4

② Given the sample: 4, 2, 1, 95, 4

→ 1, 2, 4, 4, 95

the median is 4

نلاحظ أن البيانات في المثال ① لا يوجد فيها قيم متطرفة وكان

Median = 4 في المثال ② كانت 95 عبارة عن قيمة متطرفة

و بقي Median = 4 ← الوسيط لا يتأثر بالقيم المتطرفة.

- Note that the median isn't affected by extreme values, but the mean ^{على عكس الوسط الحسابي} is affected.

3) Mode: أكثر البيانات تكررًا (التكرار).
the highest frequency data (that is, the data occurs with greatest frequency).

- Ex: (A) Given the sample: 4, 2, 2, 1, 3, 10.

→ The mode is 2. تكرر مرتين

(B) Given the sample: 4, 2, 2, 1, 4, 10

→ the mode is 2, 4 تكرر مرتين

(C) Give the sample: 4, 2, 3, 1, 10.

→ there is no mode جميع البيانات لها نفس التكرار
بالتالي لا يوجد صواب.

- Note that:

We could have no modes, 1 mod, \geq modes, 3 modes, ...

- If the sample has 1 mode: unimodal data.
- If the sample has \geq modes: bimodal data.
- If the sample has more than \geq modes: multimodal.

- Ex: The following data represents the blood type for 10 students. A, O, O, O, AB, O, B, B, A, O

→ The mode is O. لهذا في هذا المثال لا نستطيع إيجاد mean + median.

- Note that: If the data is Qualitative, then we can just find the mode.
 في البيانات النوعية، يمكن إيجاد المود.

④ Percentiles: المئينات توضع لتأسيمة تزايد البيانات من أقل قيمة لأعلى قيمة.

the pth percentile is a value from the ordered such that at least p% are less than or equal to this value and at least (100-p)% are greater than or equal to this value.

→ To find the pth percentile:-

1) order the data from smallest to largest.

2) Compute $i = \left(\frac{p}{100}\right)n$; p: the percent. ; n: the size.

3) If i isn't integer, round up: العدد الذي يليه «عند» هو موقع percentile.

If i is integer, the pth percentile = $\frac{X_i + X_{i+1}}{2}$ القيمة التي موقعها i و القيمة التي موقعها i+1 نأخذ معدلهم.

- Ex: Given the sample:

- 3450, 3550, 3650, 3480, 3355, 3310, 3490,
- 3730, 3540, 3925, 3520, 3480.

a) Find the 85th percentile (P₈₅)

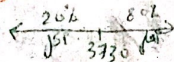
أولاً نرتب البيانات

- 3310, 3355, 3450, 3480, 3480, 3490, 3520,
- 3540, 3550, 3650, 3730, 3925

$i = \frac{p}{100} \times n = \frac{85}{100} \times 12 = 10.2$ نجد i
 $\rightarrow 11$

∴ the 85th percentile is (3730) .

موقعها رقم 11 ←



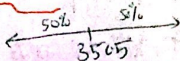
توضيح الجواب: 85% من البيانات أقل من 3730، و 20% أكثر منها.

b) Find the 50th percentile (P_{50})

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

← معدل القيمة السادسة والسابعة

∴ the 50th percentile = $\frac{3490 + 3520}{2} = 3505$



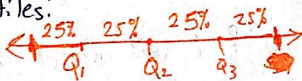
- Note that: the 50th percentile is the median.

5) Quartiles:

الربيعات

لو جزأنا البيانات بعد ترتيبها إلى أربعة أجزاء كل جزء يأخذ 25% تقريباً.

- The division points are the quartiles:



$Q_1 =$ first quartile = P_{25} (« 25th percentile »)

$Q_2 =$ Second quartile = P_{50} (« 50th percentile ») = median

$Q_3 =$ third quartile = P_{75} (« 75th percentile »)

- Ex: Give the sample: 1, 4, 4, 3, 7, 3

a) Find the first quartile. → ترتيب البيانات 1, 3, 3, 4, 4, 7

$$Q_1 = P_{25}$$

$$\rightarrow i = \frac{25}{100} \times 6 = 1.5 \rightarrow \underline{2}$$

∴ $Q_1 = \underline{3}$ القيمة التي موقعها 2

b) Find the second quartile.

$$Q_2 = P_{50}$$

$$i = \frac{50}{100} \times 6 = 3$$

$$\therefore Q_2 = \frac{3 + 4}{2} = 3.5$$

معدل القيمة الثالثة مع القيمة الرابعة

c) Find the third quartile.

$$Q_3 = P_{75}$$

$$i = \frac{75}{100} \times 6 = 4.5 \rightarrow 5 \quad \text{القصة التي موكدها 5}$$

$$\therefore Q_3 = 4$$

1, $\boxed{3}$, 3, $\boxed{4}$, 4, $\boxed{4}$, 7
 Q_1 Q_2 Q_3

- Sec 3.2: Measures of variability: مقياس التشتت
«Measures of dispersion»
تقيس هذه القراءات مقدار تشتت وتباعد القيم بحيث أنه مقياس النزعة المركزية لوحدها لا تكفي ولا تعطي صورة كاملة عن البيانات لذا يجب قياس درجة تشتت القيم وتبعرها.

① Range

→ Range = largest value - smallest value

للمدى = أكبر قيمة - أصغر قيمة

- Ex: Give the sample: 52, 33, 75, 80, 91, 45.
Find the range.

$$\rightarrow \text{Range} = 91 - 33 = 58$$

- Note: The range is affected by outliers. شكل كبير جداً

② Interquartile range (IQR): نصف المدى الربيعي
the range for the middle 50% of the data.
«the difference between the third quartile Q_3 and the first quartile»
الجال الذي تنتشر فيه 50% من البيانات

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

- Ex: Given the sample 20, 40, 10, 25, 10.

Find the interquartile range.

$$\rightarrow Q_3 = P_{75}$$

$$i = \frac{75}{100} \times 5 = 3.75 \rightarrow \textcircled{4}$$

القيمة التي يوافقها 4
بعد ترتيب البيانات
تصاعدي

10, 10, 20, 25, 40

$$\therefore Q_3 = 25$$

$$Q_1 = P_{25}$$

$$i = \frac{25}{100} \times 5 = 1.25 \rightarrow 2$$

القيمة التي يوافقها 2
بعد ترتيب البيانات
تصاعدي

10, 10, 20, 25, 40

$$\therefore Q_1 = 10$$

$$\text{Now } IQR = Q_3 - Q_1 = 25 - 10 = \underline{\underline{15}}$$

- Note:

IQR isn't affected by outliers.

③ Variance: التباين
متوسط مربعات انحرافان القيم عن وسطها
الحسابي

the average of the squared deviations about the mean.
(measure whether the data cluster around the mean).

→ Population variance (σ^2) = $\frac{\sum (X_i - M)^2}{N}$; X_i : data values.
 M : pop mean.
 N : pop size

→ Sample variance (S^2) = $\frac{\sum (X_i - \bar{X})^2}{n-1}$; \bar{X} : Sample mean.
 n : sample size.

→ S^2 is a point estimator of σ^2 .

$$S^2 = \frac{\sum X_i^2 - n\bar{X}^2}{n-1}$$

- Ex: Given the sample: 46, 54, 42, 46, 32.

$$\rightarrow \bar{X} = \frac{46 + 54 + 42 + 46 + 32}{5} = 44.$$

X	$X - \bar{X}$	$(X - \bar{X})^2$
46	2	4
54	10	100
42	-2	4
46	2	4
32	-12	144

$$\text{Now } S^2 = \frac{\sum (X - \bar{X})^2}{n-1} = \frac{256}{4} = \boxed{64}$$

- Note:

$\sum (X - \bar{X}) = 0$ « the sum of the deviations about the mean = 0 always ».

4) Standard deviation: الانحراف المعياري
الجذر التربيعي لمجموع مربعات انحرافات القيم عن وسطها.
the positive square root of the variance.

→ Sample standard deviation (s) = $\sqrt{s^2}$

→ Population standard deviation (σ) = $\sqrt{\sigma^2}$

→ the sample standard deviation is a point estimator of the population standard deviation.

- Ex: Find the standard deviation for the last example.

$$\rightarrow s = \sqrt{64} = 8$$

- Note:

The variance and standard deviation are affected by outliers.

5) Coefficient of variation. معامل الاختلاف.
مقاييس التشتت تشير إلى مدى تقارب أو تباعد البيانات عن الوسط الحسابي. أما بالنسبة لمعامل الاختلاف فهو يستخدم المقارنة بين مجموعتين مختلفتين.

→ the coefficient of variation is a relative measure of variability, it measures the standard deviation relative to the mean.

$$\rightarrow \text{Coefficient of variation (CV)} = \frac{\text{Standard deviation}}{\text{mean}} \times 100\%$$

- Ex: The mean and standard deviation for the grades of 2 sections in Stat 2311 were:

Section 1: mean = 64.87, standard deviation = 14.8

Section 2: mean = 70.45, standard deviation = 12.71

Find the coefficient of variation, then determine which section has more variable grades.

$$\rightarrow \text{Section 1: } CV = \frac{14.8}{64.87} \times 100\% = 22.81\%$$

$$\text{Section 2: } CV = \frac{12.71}{70.45} \times 100\% = 18.04\%$$

• Section 1 has more variability in grades than section 2.

CV أقل كانت البيانات متجانسة رتبة أعلى بمقارنة بالمجموع الأخرى.

* Using calculator:-

1) to find the mean:

بعض النظر كانت sample
و population

Mode 2

x_1 M+

x_2 M+

\vdots
 x_n M+

shift 2 1 =

2) to find the sample standard deviation.

shift 2 3 =

3) to find the population standard deviation.

shift 2 2 =

4) to find the variance.

Just square S or σ .

- Ex: Given the data 25, 30, 47, 80, 56, 62, 74, 80

a) Find the mean. (\bar{x} , M).

Mode 2

25 $M+$ 30 $M+$ 47 $M+$ 80 $M+$ 56 $M+$
62 $M+$ 74 $M+$ 80 $M+$

$$\text{shift } 2 \quad 1 = 56.75$$

b) Find the sample standard deviation. (S)

$$\text{shift } 2 \quad 3 = 21.47$$

c) Find the pop. standard deviation. (σ)

$$\text{shift } 2 \quad 2 = 20.08$$

d) Find the sample variance. (S^2)

$$S^2 = 21.47^2 = 460.96.$$

e) Find the pop. variance (σ^2)

$$\sigma^2 = 20.08^2 = 403.21$$

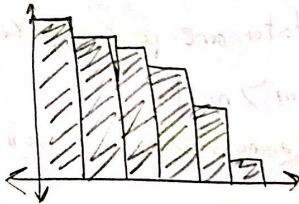
- Sec 3.3: Measures of distribution shapes Relative location, and detecting outliers:-

هناك قراءات تحدد لنا شكل التوزيع للبيانات.

* Skewness:

A histogram is a graphical display showing the shape of a distribution. We have 3 main distribution shape:-

1) positively skewed



2) negatively skewed.



3) Normal (Bell-shaped)

(Symmetric)



→ Skewness: is a measure of distribution shape.

‡ The formula for the skewness of sample data =

$$\frac{n}{(n-1)(n-2)} \sum \left(\frac{X_i - \bar{X}}{s} \right)^3 \quad \left\{ \begin{array}{l} \text{القانون} \\ \text{غير مطلق} \end{array} \right.$$

• For a symmetric distribution (normal or bell shaped):

mean = median = mode (at the center).

$$\text{skewness} = 0$$

• When the data are ^(rightly) positively skewed:

$$\text{mean} > \text{median}$$

$$\text{skewness} > 0 \quad (\text{+ve})$$

• when the data are ^(leftly) negatively skewed:

$$\text{mean} < \text{median}$$

$$\text{skewness} < 0 \quad (\text{-ve}).$$

- Note:

When data are highly skewed, the median is better than the mean.

* Z-Scores:

بالإضافة إلى مقاييس التشتت ومقاييس الموضع وسكل التوزيع σ هناك قراءات توضح موقع البيانات بالنسبة لبعضها البعض.

Measures of relative location help us determine how far a particular value is from the mean.

Suppose a sample of n observations $\rightarrow X_1, X_2, \dots, X_n$ and assume the sample mean is \bar{X} and standard deviation S

\rightarrow Z-score:

$$Z_i = \frac{X_i - \bar{X}}{S} \quad (\text{standardized value})$$

(the number of standard deviations X_i from the mean \bar{X}).

مقدار بعد القيمة عن الوسط الحسابي وهذه المقدار يكون كسبة من الانحراف المعياري.

- Ex: The following data represent the number of students in 5 classes. 46, 54, 42, 46, 32

a) Find the Z-score of the class size 46. $\bar{X} = 44.7$

$$Z = \frac{46 - 44}{8} = 0.25$$

$S = 8$
from the calculator

this Z indicates that $X_1 = 46$ is 0.25 standard deviations greater than the mean

أي أن $X = 46$ تبعد عن الوسط الحسابي بمقدار 0.25 من الانحراف المعياري.

-Notes:-

① if $x_i > \bar{x} \rightarrow z_i > 0$

② if $x_i < \bar{x} \rightarrow z_i < 0$

③ if $x_i = \bar{x} \rightarrow z_i = 0$

• If z_i and \bar{x} , s are given then we can find x_i by:-

$$z_i = \frac{x_i - \bar{x}}{s} \rightarrow \boxed{x_i = z_i s + \bar{x}}$$

-Ex: If the mean = 44.81 and standard deviation = 5.41, Find the data value whose Z-score is -2.01.

$$x_i = z_i(s) + \bar{x}$$

$$\rightarrow x_i = (-2.01)(5.41) + 44.81 = 33.94$$

* Chebyshev's theorem: skip.

* Empirical Rule: « based on the normal distribution »

the empirical rule can be used to determine the percentage of data values that must be within a specified number of standard deviations of the mean.

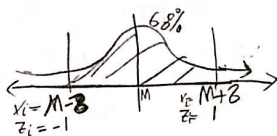
• Bell-shaped ^{normal} توزيع

→ Empirical rule:

For data having a bell-shaped distribution:-

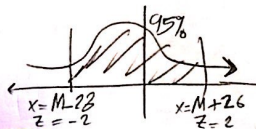
① Approximately 68% of the data values will be within one standard deviation of the mean.

• $x_i = M - \sigma$, $x_i = M + \sigma$ تقع 68% من البيانات

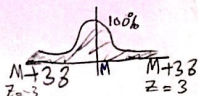


② Approximately 95% of the data values will be within 2 standard deviations of the mean.

$x_i = M - 2\sigma$, $x_i = M + 2\sigma$ تقع 95% من البيانات



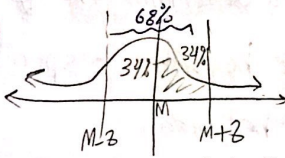
③ Almost all of the data values will be within 3 standard deviations of the mean.



• أن 99.7% من البيانات تقع بين $M - 3\sigma$ و $M + 3\sigma$
 وفي قيمة تقع خارجها تعتبر outliers

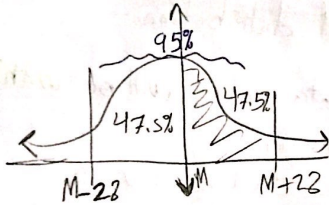
-Notes: since the bell shaped is symmetric, we have.

①

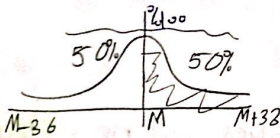


أي بين M و $M+2$ يقع 34% من البيانات، و ما نسبته 34% من البيانات يقع بين $M-2$ و M .

②

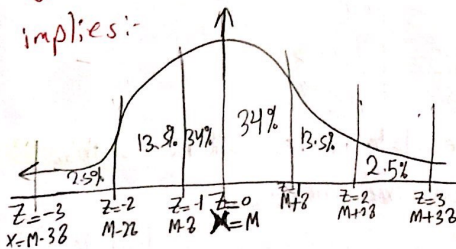


③



~~scribble~~

. this implies:



- Ex: Assume the ages of employees in BZU have a bell-shaped distribution. with mean = 44, and standard deviation = 8.

a) What is the percentage of employees with ages between 28 and 68 years.

① نرسيه Bell-shaped مستطيلة Empirical rule.

$$M = 44, \quad \sigma = 8$$

$$\rightarrow Z = -1, \quad X = M - \sigma = 36$$

$$Z = 1, \quad X = M + \sigma = 52$$

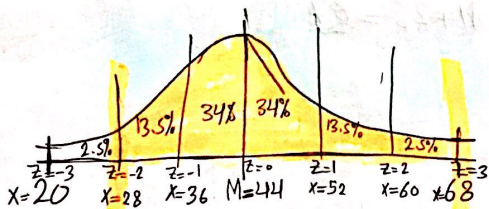
$$\rightarrow Z = -2, \quad X = 28 \quad \left. \begin{array}{l} M - 2\sigma \\ M - 2\sigma \end{array} \right\}$$

$$Z = 2, \quad X = 60 \quad \left. \begin{array}{l} M + 2\sigma \\ M + 2\sigma \end{array} \right\}$$

$$\rightarrow Z = -3, \quad X = 20 \quad \left. \begin{array}{l} M - 3\sigma \\ M - 3\sigma \end{array} \right\}$$

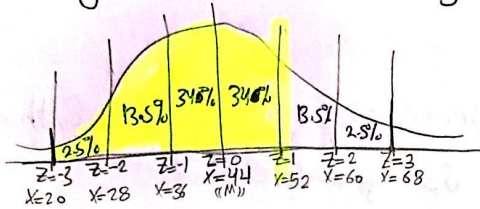
$$Z = 3, \quad X = 68 \quad \left. \begin{array}{l} M + 3\sigma \\ M + 3\sigma \end{array} \right\}$$

$$\rightarrow Z = 0, \quad X = M = 44 \quad (\text{at the center})$$



$$\therefore \text{percentage} = (13.5 + 34 + 34 + 13.5 + 2.5)\% = 97.5\%$$

b) What is the percentage of ~~students~~ employees with ages less than 52 years?



∴ the percentage is 84%

-Ex: If the weights of students in BZU have a bell shaped distribution with $M=70$ Kg and $\sigma=8.5$.

What is the percentage of students with weights between 61.5 Kg and 87 Kg? ~~M=70~~ $M=70$, $\sigma=8.5$.

$$\rightarrow Z = -1, \quad X = M - \sigma = 70 - 8.5 = \underline{61.5}$$

$$Z = 1, \quad X = M + \sigma = 70 + 8.5 = \underline{78.5}$$

$$\rightarrow Z = -2, \quad X = M - 2\sigma = \underline{53}$$

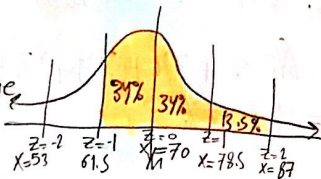
$$Z = 2, \quad X = M + 2\sigma = \underline{87}$$

$$\rightarrow Z = -3$$

$$Z = 3$$

$$\rightarrow Z = 0$$

∴ the percentage = 81.5%



* Detecting outliers:

Sometimes a data set will have one or more observations with unusually large or small values. These are extreme values or outliers.

→ If the data has a bell-shaped distribution, we can use the empirical rule to detect outliers:-

• any value is greater than $M+3\sigma$ or less than $M-3\sigma$ is an outlier.

→ that is $(X > M+3\sigma, X < M-3\sigma)$ outliers.
 $Z > 3, Z < -3$

- Ex: Given a sample:

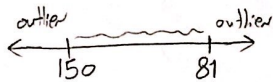
310, 500, 470, 1920, 450, 380, and

assume this sample is taken from a bell-shaped distribution with mean $M = 480$ and $\sigma = 110$. Use the empirical rule to find the outliers, if any.

~~$M = 480$~~

$$X_{\alpha} = M - 3\sigma = 480 - 3(110) = 150$$

$$X^{\alpha} = M + 3\sigma = 480 + 3(110) = 810$$



∴ $X = 1920$ is an outlier.

- Sec 3.4 : Exploratory data analysis:

• Exploratory data analysis:-

we use simple arithmetic to draw pictures to summarize data.

→ by considering 5 number summaries and box plots.

• the ^{following} five numbers summary are used to summarize the data:-

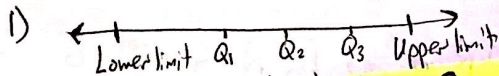
- 1) smallest value
- 2) First quartile (Q_1)
- 3) Median (Q_2)
- 4) Third quartile (Q_3)
- 5) largest value.

then we place the data in ascending order.
(من أصغر إلى أكبر)
(From smallest to largest)

• Box plot:-

is a graphical summary of data that is based on a five-number summary

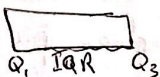
- The steps to construct the box plot:-



where lower limit $LL = Q_1 - 1.5IQR$

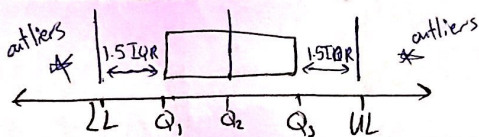
Upper limit $UL = Q_3 + 1.5 IQR$; $IQR = Q_3 - Q_1$

2) a box is drawn with the ends located at Q_1 and Q_3 .



3) A vertical line is drawn at the location of the median (Q_2)

4) data outside the limits are considered outliers. (*)



- Ex: Construct a box plot for:

12, 10, 18, 13, 25, 18

$$\rightarrow Q_1 = P_{25}$$

10, 12, 13, 18, 18, 25

$$i = \frac{25}{100} \times 6 = 1.5 \rightarrow Q_1 = 12$$

$$\rightarrow Q_2 = \frac{13 + 18}{2} = 15.5 \text{ median of } P_{50}$$

$$\rightarrow Q_3 = P_{75}$$

$$i = \frac{75}{100} \times 6 = 4.5 \rightarrow Q_3 = 18$$

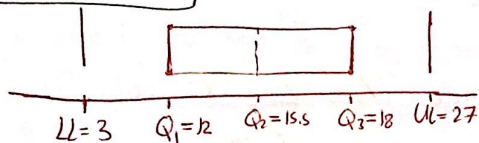
$$\rightarrow LL = Q_1 - 1.5 IQR ; IQR = Q_3 - Q_1 = 6$$

$$\therefore LL = 12 - 1.5(6) = 3$$

$$\rightarrow UL = Q_3 + 1.5 IQR$$

$$= 18 + 1.5(6)$$

$$\therefore UL = 27$$



there is no outlier.

- Ex: Consider the sample: 110, 150, 210, 180, 210, 70, 400. Use the box plot to find outliers.

\rightarrow 70, 110, 150, 180, 210, 210, 400

$$Q_1 = P_{25} \rightarrow i = \frac{25}{100} \times 7 = 1.75 \rightarrow 2 \quad \therefore Q_1 = 110$$

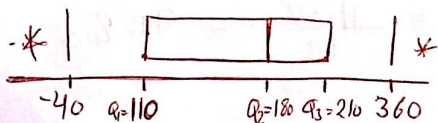
$$Q_2 = P_{50} \rightarrow i = \frac{50}{100} \times 7 = 3.5 \rightarrow 4 \quad \therefore Q_2 = 180$$

$$Q_3 = P_{75} \rightarrow i = \frac{75}{100} \times 7 = 5.25 \rightarrow 6 \quad \therefore Q_3 = 210$$

$$IQR = Q_3 - Q_1 = 210 - 110 = 100$$

$$LL = Q_1 - 1.5 IQR = 110 - 1.5(100) = -40$$

$$UL = Q_3 + 1.5 IQR = 210 + 1.5(100) = 360$$



so 400 is an outlier because $400 > UL$.

- Sec 3.6: The weighted mean and working with grouped data:-

① Weighted mean:

→ The weighted mean $\bar{x} = \frac{\sum w_i x_i}{\sum w_i}$

where x_i = value of observation i
 w_i = weight for observation i .

if we have population, the weighted mean is M .

- Ex: The grades of 4 courses were:-

	"w"	"x _i "
English	2 hours	95
Arabic	3 hours	80
Calculus 2	3 hours	92
Stat 2311	4 hours	99

Find the mean of grades.

$$\begin{aligned}\rightarrow \bar{x} &= \frac{\sum x_i w_i}{\sum w_i} = \frac{95 \times 2 + 80 \times 3 + 92 \times 3 + 99 \times 4}{2 + 3 + 3 + 4} \\ &= \frac{1102}{12} = 91.83\end{aligned}$$

*2) Grouped data:-

data are available only in a grouped or frequency distribution. We show how the weighted mean formula can be used to obtain an approximation of the mean, variance and standard deviation for grouped data.

→ Sample mean for grouped data:-

$$\bar{X} = \frac{\sum f_i M_i}{n}$$

M_i = the midpoint for class i .

f_i = the frequency for class i .

n = the sample size.

→ Sample variance for grouped data:-

$$S^2 = \frac{\sum f_i (M_i - \bar{X})^2}{n-1}$$

→ population variance for grouped data:-

$$\sigma^2 = \frac{\sum f_i (M_i - \mu)^2}{N}$$

- Ex: Consider the following data:-

class	frequency
10-14	4
15-19	8
20-24	5
25-29	2
30-34	1

1) Find the **sample mean**.

lower limit + upper limit = $\frac{\text{نقطة الوسطى}}{2}$ midpoints

class	f	mid point M	Mf
10-14	4	$\frac{10+14}{2} = 12$	$4 \times 12 = 48$
15-19	8	17	$8 \times 17 = 136$
20-24	5	22	$5 \times 22 = 110$
25-29	2	27	$2 \times 27 = 54$
30-34	1	32	$1 \times 32 = 32$

$$\therefore \bar{x} = \frac{\sum Mf}{\sum f} = \frac{48 + 136 + 110 + 54 + 32}{4 + 8 + 5 + 2 + 1} = \frac{380}{20} = 19$$

2) Find the **sample variance**. $\bar{X} = 19$ (from part 1)

class	f	midpoint M	مقطع - كل الوسطية $M - \bar{X}$	مربع القطوع $(M - \bar{X})^2$	$f(M - \bar{X})^2$ → دمجها في الجدول الثاني
10-14	4	12	12-19 = -7	49	196
15-19	8	17	17-19 = -2	4	32
20-24	5	22	22-19 = 3	9	45
25-29	2	27	27-19 = 8	64	128
30-34	1	32	32-19 = 13	169	169

∴ $S^2 = \frac{\sum f(M - \bar{X})^2}{\sum f - 1} = \frac{196 + 32 + 45 + 128 + 169}{(4 + 8 + 5 + 2 + 1) - 1}$

نضع الجدول السابق في مجموع الجدول الثاني

→ $S^2 = \frac{570}{19} = \underline{30}$

3) Find the **sample standard deviation**. $\bar{X} = 19$

→ $S = \sqrt{30} = \underline{5.48}$

4) If these data are from **population**, find the **population variance**.

→ $\sigma^2 = \frac{\sum f(M - \bar{M})^2}{\sum f} = \frac{570}{20} = \underline{28.5}$

نجدها كما أسبقا، ولتة نفس $\sum f$ في الوسط

5) find the **population standard deviation**.

→ $\sigma = \sqrt{28.5} = \underline{5.34}$

* Using SD mode. «calculator fx-82MS».

1) To find the weighted mean.

X_i	w_i
x_1	w_1
x_2	w_2
\vdots	\vdots
x_n	w_n

① Mode 2

② x_1 M+ x_2 M+ ... x_n M+
 يدخل القيمتين و بين كل قيمة
 و اخرى تضغط M+

③ $\nabla \nabla w_1 = \nabla \nabla w_2 = \nabla \nabla w_3 = \dots$
 $\nabla \nabla w_n =$

④ shift 2 1 =

2) To find the mean of grouped data.

classes f

$$M_i = \text{midpoint} = \frac{\text{lower limit} + \text{upper limit}}{2}$$

① We find the midpoints (M_1, M_2, \dots, M_n)

② mode 2

③ M_1 M+ M_2 M+ ... M_n M+
 يدخل midpoint و بين كل قيمة
 و اخرى تضغط M+

④ $\nabla \cdot \nabla f_1 = \nabla \nabla f_2 = \dots = \nabla \nabla f_n =$

⑤ shift 2 1 =

3) To find the ^{sample} standard deviation of grouped data.

shift 2 3 =

4) To find the ~~sample~~ population standard deviation of grouped data.

shift 2 2 =

5) To find the sample variance of grouped data.
we just square s.

6) To find the population variance of grouped data.
we just square s.

- Ex: Solve the previous example using calculator.

class	f	Midpoints
10-14	4	12
15-19	8	17
20-24	5	22
25-29	2	27
30-34	1	32

① Find the sample mean.

① $\text{mode} = 2$

② $12 \text{ [M+]} 17 \text{ [M+]} 22 \text{ [M+]} 27 \text{ [M+]} 32 \text{ [M+]}$

③ $\nabla 4 = \nabla \nabla 8 = \nabla \nabla 5 = \nabla \nabla 2 = \nabla \nabla 1 =$

④ $\text{shift} 2 1 = \underline{19}$

② Find the sample standard deviation.

$\text{shift} 2 3 = 5.48$

③ Find the population standard deviation.

$\text{shift} 2 2 = 5.34$

④ Find the sample variance, population variance.

$S^2 = 30$ "square part (2)", $\sigma^2 = 28.5$ "square part (3)!"

- Sec 3.5: Measures of Association Between two variables

العلاقة بين متغيرين

أخذنا بيانات لتلخيص بيانات بمتغير واحد، وهناك نوعين من العلاقات التي تربط بين متغيرين.

① Covariance:

التباين بين متغيرين
أو التباين بين متغيرين

For a sample of size n with the observations:-

$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

→ the sample covariance $S_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$;

\bar{x} : sample mean for variable X .

\bar{y} : sample mean for variable Y .

→ the population covariance $\sigma_{xy} = \frac{\sum (x_i - M_x)(y_i - M_y)}{N}$;

M_x : population mean for X .

M_y : population mean for Y .

N : population size

- Ex: The following table summarizes the number of absences (x) and the grade (y) of a sample of students.

x_i	6	3	1	4	5
y_i	65	81	94	76	69

Find the sample covariance.

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

$$\rightarrow \textcircled{1} \quad \bar{x} = \frac{\sum x_i}{n} = \frac{6+3+1+4+5}{5} = 3.8$$

$$\textcircled{2} \quad \bar{y} = \frac{\sum y_i}{n} = \frac{65+81+94+76+69}{5} = 77$$

$$\textcircled{3}$$

x_i	y_i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$
6	65	$6 - 3.8 = 2.2$	$65 - 77 = -12$	$(2.2)(-12) = -26.4$
3	81	-0.8	4	-3.2
1	94	-2.8	17	-47.6
4	76	0.2	-1	-0.2
5	69	1.2	-8	-9.6

$$\textcircled{4} \quad S_{xy} = \frac{-26.4 + -3.2 + -47.6 + 0.2 + -9.6}{5 - 1} = -21.75$$

ملاحظة: لو كانت population قسم على N

• Interpretation of the covariance:-

The covariance is a measure of the linear association between 2 variables. يقوس نوع العلاقة الخطية بين متغيرين

- ① if $S_{xy} > 0$ («positive»), indicates a positive linear relation between X and Y. («if X increases y increases»)
- ② if $S_{xy} < 0$ («-ve»), indicates a negative linear relation between X and Y. («if X increases, y decreases»)
- ③ if $S_{xy} = 0$, indicates no linear relation between X and Y.

② Correlation Coefficient:

معامل الارتباط

measures the type and strength of the relation between X and Y.

→ the sample correlation coefficient $r_{xy} = \frac{S_{xy}}{S_x S_y}$;

S_{xy} : sample covariance

S_x : sample standard deviation of X.

S_y : sample standard deviation of Y.

→ the population correlation coefficient $\rho_{xy} = \frac{\Sigma xy}{\Sigma x \Sigma y}$;

Σxy : population covariance.

Σx : population standard deviation of X.

Σy : population standard deviation of Y.

→ The sample correlation coefficient r_{xy} is the point estimator of the population correlation coefficient. (r_{xy} is a point estimator of ρ_{xy})

• Interpretation of the correlation coefficient:-

① $-1 \leq r_{xy} \leq 1$

② if $r_{xy} = 1$, then there is a perfect positive linear relationship between X and Y.

③ if $r_{xy} \cong 1$, a strong positive linear relationship.

④ if $r_{xy} = -1$, a perfect negative linear relationship.

⑤ if $r_{xy} \cong -1$, a strong negative linear relationship.

⑥ if $r_{xy} = 0$, there is no linear relationship.

⑦ if $r_{xy} \cong 0$, there is a weak linear relationship.

~~Ex: Consider the following sample data.~~

-Ex: Consider the following sample data.

X_i	Y_i
5	10
10	30
15	50

Find the sample correlation coefficient.

$$\textcircled{1} \bar{x} = \frac{\sum x_i}{n} = \frac{5+10+15}{3} = 10$$

$$\bar{y} = \frac{\sum y_i}{n} = \frac{10+30+50}{3} = 30$$

$$\textcircled{2} S_x, S_y$$

x_i	y_i	$x_i - \bar{x}$	$(x_i - \bar{x})^2$	$y_i - \bar{y}$	$(y_i - \bar{y})^2$
5	10	-5	25	-20	400
10	30	0	0	0	0
15	50	5	25	20	400

$$\therefore S_x^2 = \frac{\sum (x_i - \bar{x})^2}{n-1} = \frac{25+0+25}{3-1} = \frac{50}{2} = 25$$

$$\rightarrow S_x = \sqrt{25} = 5$$

$$\therefore S_y^2 = \frac{\sum (y_i - \bar{y})^2}{n-1} = \frac{400+0+400}{3-1} = \frac{800}{2} = 400$$

$$\rightarrow S_y = \sqrt{400} = 20$$

$$\textcircled{3} S_{xy}$$

$(x_i - \bar{x})(y_i - \bar{y})$
$(-5)(-20) = 100$
0
100

$$\therefore S_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1} = \frac{200}{2} = 100$$

$$\textcircled{4} r_{xy} = \frac{S_{xy}}{S_x S_y} = \frac{100}{(5)(20)} = 1 \quad \therefore \text{there is a perfect linear relationship.}$$

* Using calculator (SD mode) to find r_{xy} , S_{xy} .

(x_1, y_1) (x_2, y_2) (x_n, y_n) .

~~1) $\frac{1}{n} \sum x_i$~~

① mode 3 1

② x_1 2 y_1 $M+$

x_2 2 y_2 $M+$

⋮

x_n 2 y_n $M+$

③ To find the sample mean of $X \rightarrow \bar{X}$

shift 2 1 =

لو كان pop نفس الخطوات M_x \bar{X}

④ To find the sample mean of $Y \rightarrow \bar{Y}$

shift 2 ~~1~~ 1 =

لو كان pop نفس الخطوات M_y

⑤ To find the sample standard deviation of $X \rightarrow S_x$

shift 2 3 =

لو كان pop نفس الخطوات S_x $\frac{1}{n} \sum x_i^2$ $\frac{1}{n} \sum x_i$

⑥ To find the sample standard deviation of $Y \rightarrow S_y$

shift 2 ~~1~~ 3 =

لو كان pop نفس الخطوات S_y $\frac{1}{n} \sum y_i^2$ $\frac{1}{n} \sum y_i$

⑦ To find the ~~est~~ sample correlation coefficient $\rightarrow r_{xy}$.

shift 2 ▷ ▷ 3 = (اللوكانه P.P. (Sxy) نفس الخطوات))

⑧ To find the sample covariance $\rightarrow S_{xy}$.

$$S_{xy} = r_{xy} S_{xx} S_y.$$

(اللوكانه P.P. (Sxy))

$$\rightarrow S_{xy} = S_{xx} S_y$$