

Statistics

Chapter 1

The Nature of Probability and Statistics

Statistics is the science of conducting studies to collect, organize, summarize, analyze, and draw conclusions from data.

Main Areas of Statistics

Descriptive statistics consists of the collection, organization, summarization, and presentation of data.

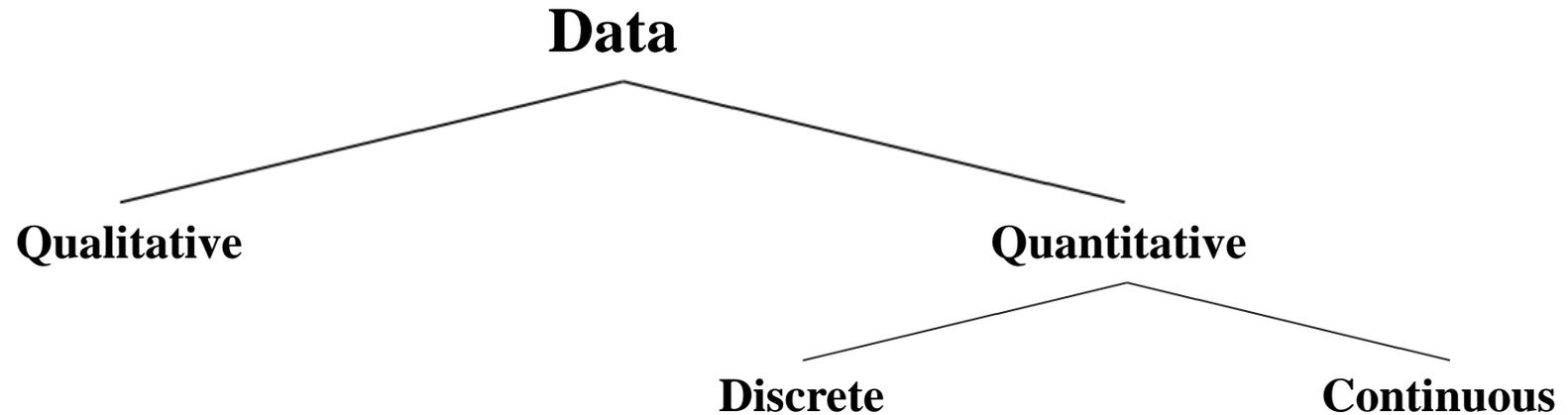
Inferential statistics consists of *generalizing* from samples to populations, performing estimation and hypothesis tests, determining relationships among variables, and making predictions.

Inferential statistics uses **probability** (the likelihood of an outcome occurring) to make conclusions and predictions.

A **population** consists of all subjects that are being studied.

A **sample** is a group of subjects selected from a population.

Data are the values (measurements or observations) that the variables can assume.



A **data set** is a collection of **data values**.

A **variable** is a characteristic or attribute that can assume different values.

A **random variable** is a variable whose values are determined by chance.

Types of Variables

Qualitative variables are variables that can be placed into categories, according to some characteristic or attribute.

Quantitative variables are numeric in nature and can be ordered or ranked. Quantitative variables can be either **discrete** or **continuous**.

A **continuous variable** is a quantitative variable that can assume ANY numerical value between any two specific values.

- Obtained by measuring.
- May include fractions and decimals.

A **discrete variable** is a quantitative variable that has either a finite number of possible values or a countable number of possible values.

- Countable means that the values result from counting, such as 0, 1, 2, 3...

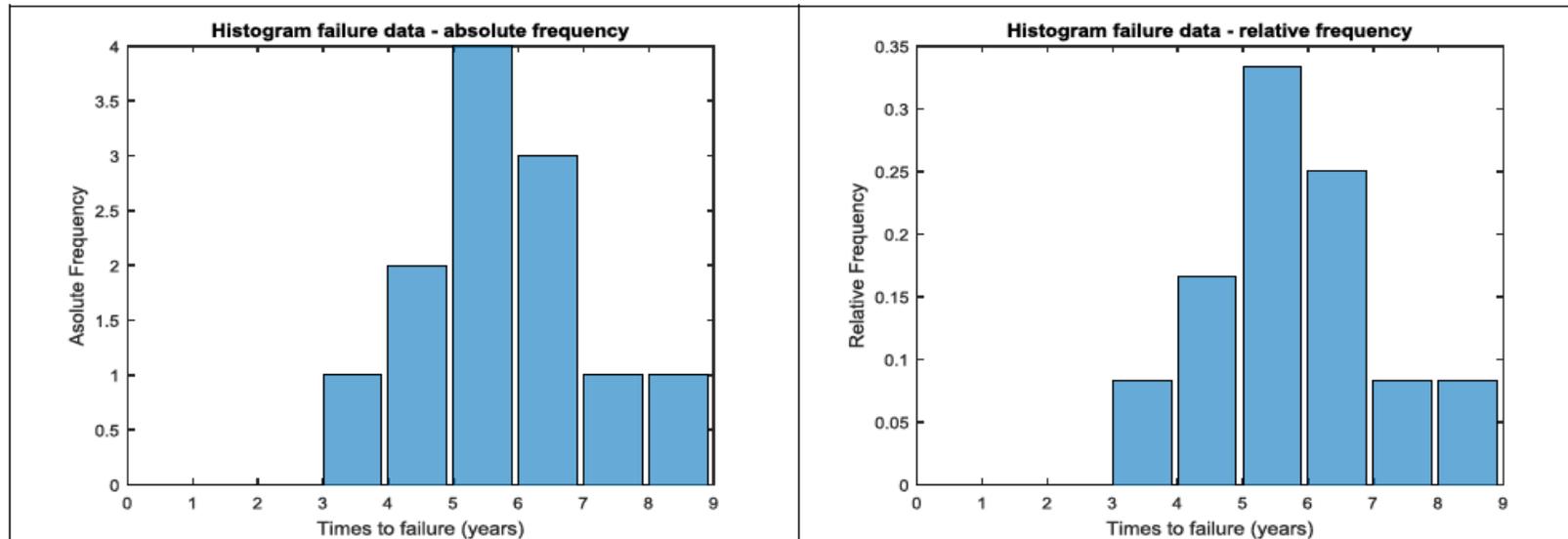
Histogram

Histogram: set of columns describing the distribution of specific data. Each column represents a number of items or frequency.

This graphic gives an idea of sample shape.

Absolute frequency: It indicates the number of data points.

Relative frequency: It indicates the percentage of the total number of data points.



The characteristics of any data can be summarized by knowing some statistical properties:

1. The data center is measured using the median or mean.
2. Data spread is measured using variance and standard deviation.
3. Knowing that most of the data is located using the quartiles.
4. Find outlie

Median

When arranging the data in ascending or descending order, the median is the statement in which 50% of the data fall before and 50% after it, and it is calculated:

$$X(1), X(2), \dots, X(n-1), X(n)$$

$$\tilde{X} = \begin{cases} X\left(\frac{n+1}{2}\right) & \text{if } n \text{ is odd} \\ \frac{X\left(\frac{n}{2}\right) + X\left(\frac{n}{2} + 1\right)}{2} & \text{if } n \text{ is even} \end{cases}$$

Example 1: Find the median for the data listed below 60 72 40 80 63.

Sol:

We arrange the data in ascending order

40 60 63 72 80

Since the number of data is odd, the median value (third) represents the median

$$\tilde{X} = X_{(5 + 1)/2} = X_3 = 63$$

Example 2: Find the median for the data listed below 72 60 72 40 80 63.

Sol:

We arrange the data in ascending order

40 60 63 72 72 80

Since the number of data is even, the median value is the average of the two readings in the middle of the data

$$\tilde{X} = (63 + 72) / 2 = 67.5$$

Statistics

Chapter 1

The Nature of Probability and Statistics

Data Collection and Sampling Techniques

Data is often collected via surveys.

1. Telephone Surveys
 2. Mailed Questionnaire
 3. Personal Interview
 4. Internet Survey
- What are advantages and disadvantages of data collection through surveys?

	Advantages	Disadvantages
Telephone Surveys	<ol style="list-style-type: none"> 1. Less costly than personal interview. 2. People may be more candid. 	<ol style="list-style-type: none"> 1. Not everyone has a phone. 2. Cell phones typically not included. 3. Tone of interviewer's voice may affect response.
Mailed Questionnaire	<ol style="list-style-type: none"> 1. Can cover wider geographic area than phone survey or personal interview 2. Respondents can remain anonymous 3. Less expensive 	<ol style="list-style-type: none"> 1. Low number of responses 2. Inappropriate answers to questions 3. Low reading abilities or not understanding questions may create useless responses.

	Advantages	Disadvantages
Personal Interview	In-depth responses to questions	<ol style="list-style-type: none">1. Interviewers must be trained in asking questions and recording responses (which is costly).2. Interviewer may be biased in the selection of participants.
Internet survey	As with telephone and mail surveys, <ol style="list-style-type: none">1. Inexpensive, often free2. Candor3. Large geographic coverage4. Anonymity	<ol style="list-style-type: none">1. Will miss demographics without computer access2. May have inappropriate answers if questions are misunderstood

Sampling Techniques

- Researchers use samples to collect data and information about a particular variable from a population.
- Samples save time, money, and may actually allow a researcher to collect better information.
- Samples need to be representative of the population or they are meaningless in drawing conclusions about the population.
- Sampling must be done in a way that the samples are **unbiased**—that each subject in the population has an equal chance of being in the sample.

Technique	Description
Random sampling	Uses chance methods or random numbers to select the sample. Everyone or everything from the population has the same chance of being selected for the sample and it is the best way of obtaining a representative sample.
Systematic sampling	Numbers each subject of the population and then selects every kth subject.
Convenience sampling	Selects subjects that are convenient for the researcher. These samples are typically of not statistical value.
Stratified sampling	Divides the population into groups (called strata) according to some characteristic that is important to the study, then randomly samples subject from each group.
Cluster Sampling	Divides the population into groups called clusters by some means such as geographic area or schools in a school district, etc. Then randomly select some of the clusters and use ALL members of the selected clusters.

Observational and Experimental Studies

Observational study - the researcher merely observes what is happening or what has happened in the past and tries to draw conclusions based on these observations.

Experimental study - the researcher manipulates one of the variables and tries to determine how the manipulation influences other variables.

Statistics

The summary of lecture 1 and lecture 2

The importance and concept of statistics

اهمية ومفهوم علم الاحصاء

يعد علم الاحصاء اليوم من اهم العلوم التي تتوقف عليها التنمية السياسية و الاقتصادية و الثقافية ... الخ و للاحصاء حصة اساسية من عمل الدول و المؤسسات و المنظمات السياسية و الاقتصادية و الاجتماعية ، و كثيرا ما يرتهن مصير مشاريع او قرارات كبرى بالنتائج التي يقدمها الاحصاء في مجال معين. و بصورة عامة، فأن افتقاد الجهد الاحصائي، في مجال من المجالات، يمنع من التأكد و تحصيل الضمان في استجابة اي مشروع للواقع، كما يحول دون تحديد مدى نجاحه او اخفاقه، و يجعل في الاقدام عليه شيئا من المخاطر.

ويمكننا فهم الاحصاء بصورة مختصرة كالآتي:

اولا: الاحصاء قادر على توصيف الظواهر العلمية توصيفا رقميا كميا دقيقا و اكثر وضوحا و قربا من الواقع.

ثانيا: يستطيع الاحصاء ان يفسر الظواهر، و ان يحدد مدى تأثير العوامل المفترضة، كما يمكنه التنبؤ بالمستقبل بالمعنى العلمي للكلمة. مثلا يمكن معرفة حالة الطقس في يوم ما بالاعتماد على قاعدة بيانات تظهر حالة الطقس لنفس اليوم بالسنين السابقة.

ثالثا: يعتمد الاحصاء المعادلات الرياضية و حساب الاحتمالات، و يركز على اسس علمية رياضية مبرهن عليها.

رابعا: الاحصاء هو منهج و عقل و تفكير و الية تأمل و نمط قراءة. فهو غير محدد بمادة علمية سوى تلك التي تحتوي نظام العينات.

Definition of statistics

تعريف علم الاحصاء

نعني بالاحصاء العلم الذي يكون موضوعه جمع وتنسيق وتصنيف وتعليل الوقائع او المعطيات الرقمية او العددية في كل صنف، بحيث يمكن الحصول على نسب عدديّة مستقلة استقلالاً ملموساً عن المصادفة و استثناءاتها، و دالة على وجود العلل المنتظمة التي ينعدم حصولها بوجود العلل الفجائية، فهو يتخذ طريقة للتحليل في العلوم الدقيقة و العلوم الاجتماعية و في المشروعات الاقتصادية على اختلافها. وهو يعني بالتنبؤ باحتمالات حدوث امر بعينه او حاله بعينها. و تطور علم الاحصاء تطورا كبيرا بعد استحداث الحاسبة الالكترونية التي تتعامل مع كميات من الارقام الضخمة تعاملًا سريعًا. ان المعنى الحديث لعلم الاحصاء انه العلم الذي يبحث في جمع البيانات، و تنظيمها و عرضها و تحليلها و استقراء النتائج و اتخاذ القرار بناء عليها.

المجتمع الاحصائي (population): هو جميع المفردات و المعلومات التي تتمتع بصيغة ما أو خاصية ما مشتركة ضمن اطار عام واحد ويكون المجتمع الاحصائي مقيدا بزمان و مكان او لا يكون و قد يكون محددًا مثل عدد طلاب المرحلة الاولى في كلية الهندسة او يكون غير محدد مثل عدد الاسماك في المحيط الهادي.

العينة (Sample): و هي جزء من المجتمع الاحصائي محل الدراسة يتم اختياره اما عشوائيا او بصورة علمية خاصة.

المتغير (Variable): خاصية او صفة للأفراد او الاشكال تختلف من عنصر الى الاخر. وهو بنوعين:

- 1- متغير نوعي (وصفي): ويشمل اي صفة او ظاهرة تتغير نوعيا و تسجل بوصف لفظي مثل ذكر ، انثى، حار ، بارد.
- 2- متغير كمي (رقمي): ويشمل صفة او ظاهرة تتغير كميا وتسجل بارقام عددية مثل معدل درجة الحرارة ، منسوب سقوط الامطار. و المتغير الرقمي يقسم الى قسمين:
 - الف- متغير كمي متصل (continuous variable): و يكون فيه استمرارية او بمعنى اخر تقبل القيم الكسرية مثل درجة الحرارة والوزن فنستطيع ان نقول 23.5 درجة او 75.3 كغم.
 - ب- متغير كمي منفصل (discrete variable): و هي تكون متقطعة او لا تقبل قيم كسرية مثل عدد الطلاب في كلية الهندسة.

The Data

البيانات: هي مجموعة من الحقائق و المعلومات التي تتعلق بظاهرة ما و تشكل مادة خام للاحصاء.

مصادر البيانات الاحصائية:

- 1- مصادر تاريخية: من السجلات و الوثائق الاصلية و الثانوية و الارشيفات الرسمية و الخاصة سواء كانت مبوبة في قواعد بيانات او عشوائية.
- 2- مصادر ميدانية: و هي تجمع مباشرة بالمشاهدات الميدانية من وحدات المجتمع الاحصائي.

اساليب جمع البيانات:

- 1- اسلوب الحصر الشامل: جمع المعلومات من جميع مفردات المجتمع الاحصائي بدون استثناء. (حساب نسبة النجاح للمرحلة الاولى).
- 2- اسلوب العينة: يتم جمع المعلومات من جزء مختار من المجتمع الاحصائي عشوائيا او غير عشوائي لتمثيل جميع عناصر المجتمع (اي الحكم على كل المجتمع من خلال صفات و خصائص عينة ما مأخوذة منه).

اقسام الاحصاء: يقسم الاحصاء الى قسمين هما:

- 1- الاحصاء الوصفي: يشمل جمع و تبويب و تصنيف البيانات الاحصائية، و جعل البيانات المتعددة بشكل مفهوم و ذات مدلول يسهل التعامل معها.
- 2- الاحصاء الاستقرائي (الاستدلالي): يشمل استقراء النتائج و اتخاذ القرارات، و يتعامل مع طرق تحليل و تفسير و استخلاص الاستنتاجات لاتخاذ افضل قرار ممكن عندما تكون المعلومات غير وافية.

Statistics

Lecture 3

Data display methods

طرق عرض البيانات

ان ترتيب وتنظيم البيانات وعرضها بشكل يسهل تمييزها والتعرف على خصائصها يعتبر الحجر الاساس والخطوة الاولى في التحليل الاحصائي وفي بعض الاحيان يكون الغرض من استخدام طرق عرض البيانات هو جذب الانتباه ليس الى الارقام بحد ذاتها وانما الى ميول واتجاه تلك الارقام بالزيادة والنقصان وبالتالي معرفة ماهية او طبيعة الظاهرة المدروسة دون الخوض في تفاصيلها. طرق عرض البيانات هي:

a- Tables Method (طريقة الجداول):

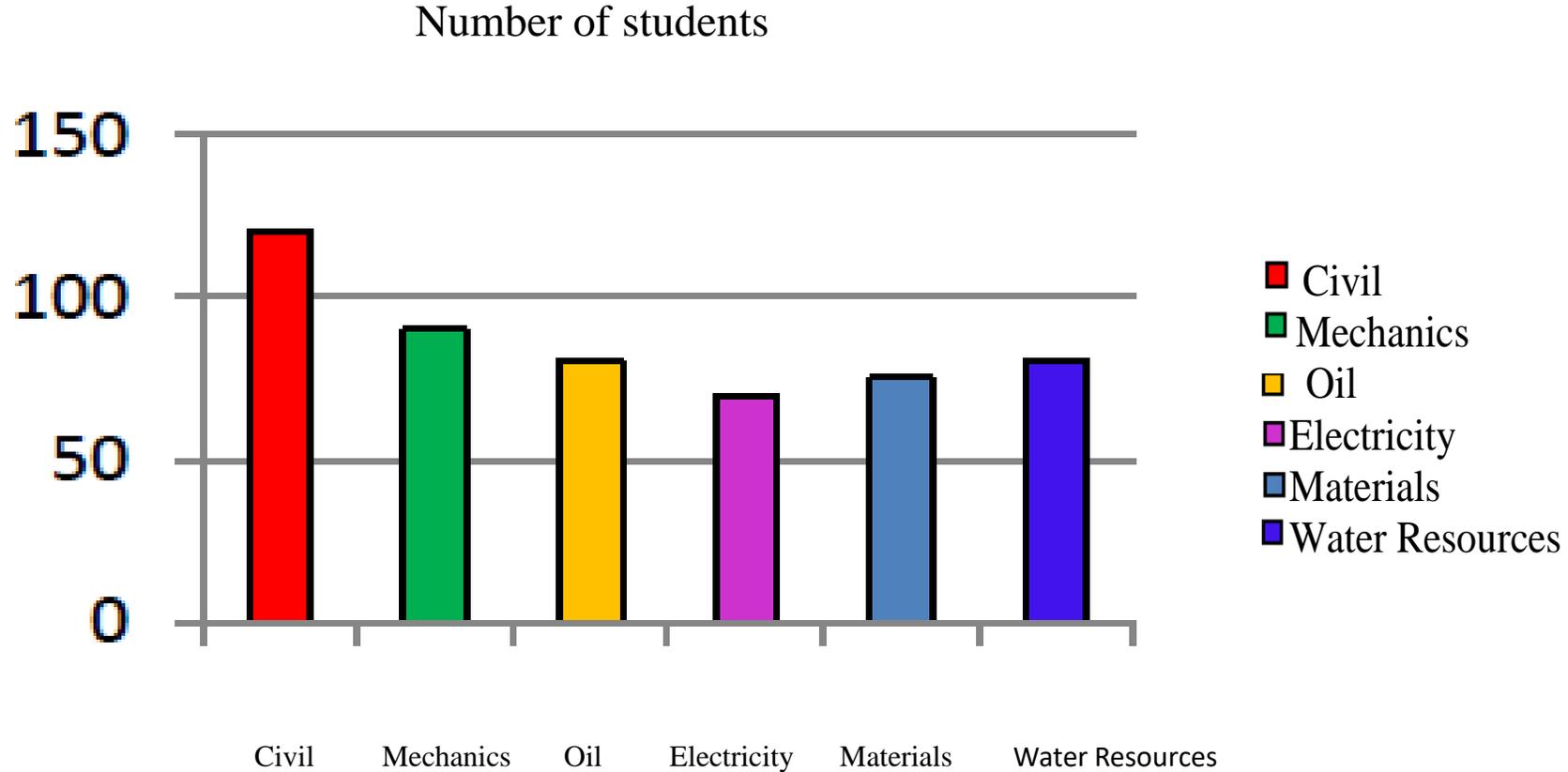
وتشمل وضع البيانات في جداول وتستهمل في عرض تغيير ظاهرة ما مع الزمن او مع مسميات اخرى معينة.

Ex: Represent the distribution of the students of the engineering college department that admitted for the year 2013-2014 among the scientific departments:

Scientific department	Civil	Mechanics	Oil	Electricity	Materials	Water Resources
Number of students	120	90	80	70	75	80

b- Rectangles or Columns Method (Bar charts) (طريقة المستطيلات او الاعمدة):

وتشمل وضع المسميات على المحور الافقي او المحور العمودي ورسم مستطيل على كل مسمى بحيث يكون ارتفاع كل مستطيل ممثلا للقيمة المقابلة لذلك المسمى وذلك باستعمال مقياس رسم المناسب. ان هذه الطريقة هي من اكثر الطرق شيوعا ويمكن عرض عدة ظواهر مع الزمن.



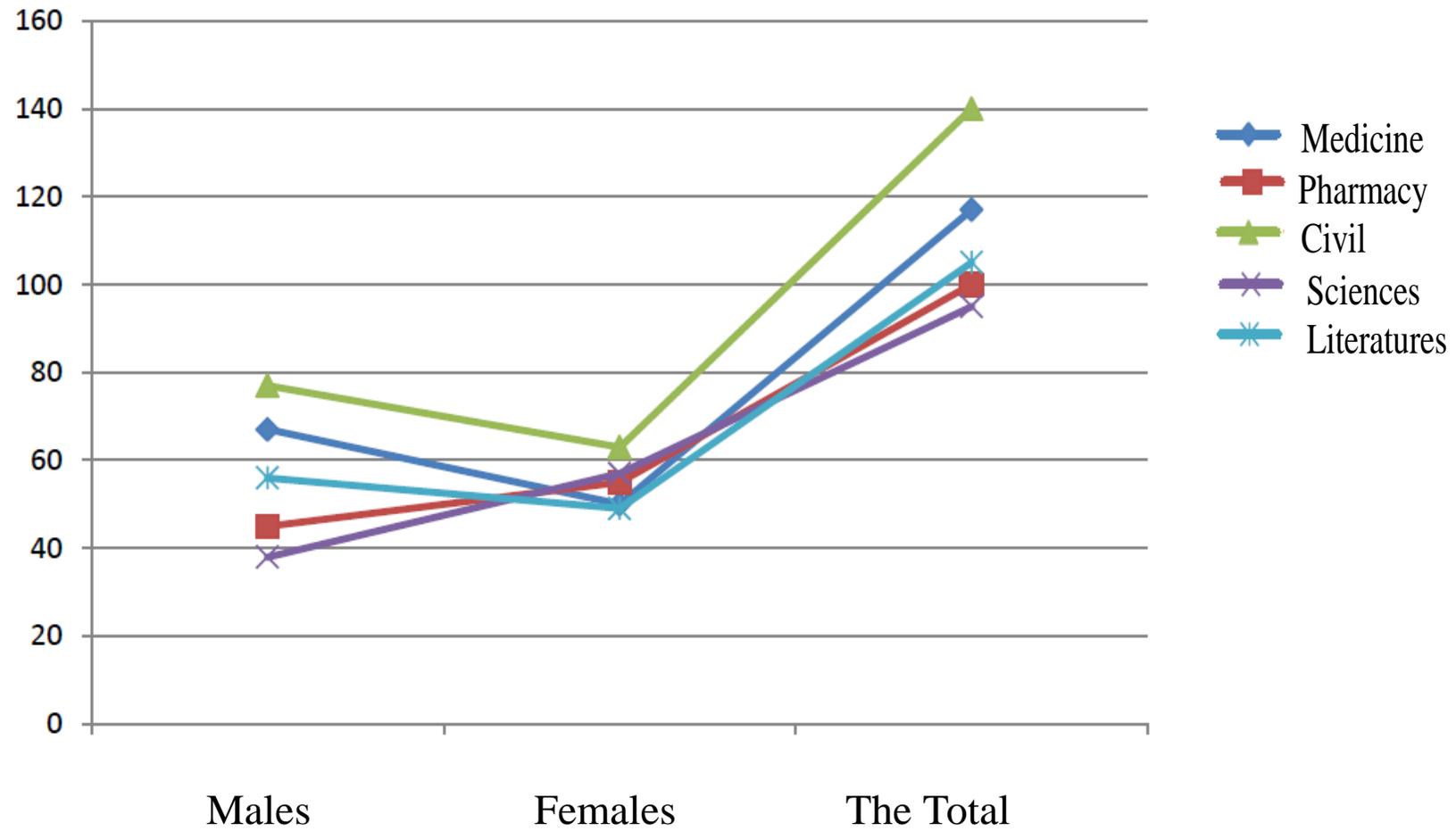
c- Included line method (Broken) (طريقة الخط المدرج او المتكسر):

وتشمل عرض البيانات الناتجة عن تغيير ظاهرة او عدة ظواهر مع الزمن مثل تغيير عدد طلبة المدارس الثانوية مع السنوات، تغيير درجة الحرارة للمريض خلال ساعات، مجموع ارصدة البنوك الخاصة مع السنوات... الخ. وعادة ما تستخدم من قبل خبراء البورصة في متابعة صعود وهبوط الاسهم بشكل يومي او كل ساعة.

Ex: The following table represents the distribution of graduates from Baghdad University. Represent the following data in the table in the included line method.

Specialization	medicine	Pharmacy	civil	Sciences	Literatures
Males	67	45	77	38	56
Females	50	55	63	57	49
The Total	117	100	140	95	105

Sol:



d- Curved line method (طريقة الخط المنحني):

تمثل طريقة الخط المتدرج و نحصل عليها بتمديد الخط المتكسر ليصبح منحنى، و تستعمل هذه الطريقة عندما تتغير الظاهرة على فترات قصيرة و كثيرة.

e- Pictorial method: (الطريقة التصويرية)

وتستعمل كه بديل لطريقة المستطيلات وتستعمل في الاعلانات والتقارير الحكومية وكتب الاطفال.

f- Divided circuit method (bie-charts) (طريقة الدائرة المقسمة) :

وتكون بتمثيل البيانات بدائرة كاملة مقسمة الى اجزاء بحسب الظواهر، واكثر ما تستخدم في المجالات الاقتصادية و تستخدم عندما يكون الهدف ابراز اجزاء الظاهرة المدروسة.

Ex: The following table represents the qualifications of the working team of a contracting company. Represent the following using the bie-charts method.

Qualification	Engineer	Supervisor	Expert worker	Worker
The number	3	6	5	20

Sol:

Total summation for team = 3+6+5+20 = 34

Central angle for the engineers qualification = $(3/34)*360^\circ = 31.764^\circ$

Central angle for the supervisor qualification = $(6/34)*360^\circ = 63.530^\circ$

Central angle for the expert worker qualification = $(5/34)*360^\circ = 52.941^\circ$

Central angle for the worker qualification = $(20/34)*360^\circ = 211.765^\circ$

نحسب الزوايا المركزية لكل قيمة من لبيانات:

الان يتم رسم دائرة تقسيم البيانات.



Frequency distributions

التوزيعات التكرارية

وتشمل تنظيم البيانات الكثيرة بحيث لا تفقد هذه البيانات من أهميتها الا الشيء القليل او لا تفقد ابدا. ويقسم التوزيع التكراري الى نوعين:

a- Frequency distributions (توزيع تكراري بسيط):

Ex: The following data represent the marks of 22 students in one of the exams. Represent the following data in a simple frequency table.

(19, 20, 19, 17, 15, 20, 18, 14, 12, 15, 18, 15, 17, 19, 11, 20, 14, 15, 13, 20, 15, 18)

Sol:

في التوزيع التكراري البسيط نحتاج الى:

1- قيم الفئات او القياسات (x_i) . (Class values).

2- التكرارات المقابلة لهذه الفئات (f_i) . (Frequencies).

الدرجة (X_i)	11	12	13	14	15	16	17	18	19	20	
التكرار (f_i)	1	1	1	2	5	0	2	3	3	4	$\sum f_i = 22$

ملاحظة: ان مجموع التكرارات يجب ان تساوي عدد البيانات الاصلية.

b- Distribution of classified frequency (توزيع تكراري ذو فئات) :

يستخدم هذا النوع من التوزيعات التكرارية في تنظيم البيانات ذات العدد الكبير جدا و ذات المدى (الفرق بين أعلى وأقل قيمة بالبيانات) كبير.

Ex: The following data represent the results of the compression test for fifty concrete cubes in units (MPa). Arrange the following data in the frequency table with the classes.

(29 30 32 39 38 38 44 28 33 31 35 37 42 49 25 34 30 31 35 37 40 26 32 33 33 39 44 45 26 31 34 31 36 36 41 27
30 31 32 37 40 39 25 34 31 30 38 35 43 48)

Sol:

1- نفرض عدد فئات (Number of classes) مناسب (عدد الفئات المفروض في اي جدول تكراري ذو فئات يجب ان لا تقل عن 5 فئات و لا يزيد عن 15 فئة).

في هذا المثال سنفرض عدد الفئات = 5 .

2- يستخرج مدى البيانات (rang of the data) = أكبر قيمة بالبيانات – أقل قيمة بالبيانات = 49 – 25 = 24

(rang) المدى

3- نحسب طول الفئة (Class length) الواحدة = $\frac{\text{المدى}}{\text{عدد الفئات (Number of classes)}}$

$$= \frac{24}{5} = 4.8 = 5 \text{ (يقرب الى اعلى عدد صحيح قريب)}$$

4- نفرض الحد الادنى للفئة الاولى (The minimum limit for the first class) = اصغر قيمة بالبيانات = 25

5- نحسب الحد الادنى الفعلي للفئة الاولى (The actual minimum for the first class) = الحد الادنى للفئة الاولى - 0.5 = 24.5

6- نحسب الحد الاعلى الفعلي للفئة الاولى (The actual upper limit for the first class) = الحد الادنى الفعلي للفئة الاولى + طول الفئة = 24.5 + 5 = 29.5

7- نحسب الحد الاعلى للفئة الاولى (The upper limit for the first class) = الحد الاعلى الفعلي للفئة الاولى - 0.5 = 29.5 - 0.5 = 29

8- يتم تحديد الحدود الدنيا و العليا لباقي الفئات بأضافة طول الفئة لكل حد.

9- تعيين مراكز الفئات (X_i) (center of classes) = $\frac{\text{الحد الادنى للفئة} + \text{الحد الاعلى للفئة}}{2} = \frac{29+25}{2} = 27$

10- يتم حساب مراكز الفئات الباقية = مركز الفئة السابقة + طول الفئة

11- جمع التكرارات (Collect of Frequencies) و الناتج يجب ان يساوي عدد البيانات الاصلي.

Class limits حدود الفئات	Actual limits for classes الحدود الفعلية للفئات	Centers of classes مراكز الفئات	Frequencies التكرارات
25 - 29	24.5 – 29.5	27	7
30 - 34	29.5 – 34.5	32	19
35 - 39	34.5 – 39.5	37	14
40 - 44	39.5 – 44.5	42	7
45 - 49	44.5 – 49.5	47	3
			$\sum f_i = 50$

Other types of frequency distributions

انواع اخرى من التوزيعات التكرارية

a- Congregate frequency distribution (Cumulative) ((التوزيع التكراري المتجمع(التراكمي))):

التوزيع التكراري المتجمع (التراكمي): وهو بنوعين التوزيع التكراري المتجمع الصاعد و التوزيع التكراري المتجمع النازل.
يتم حساب التوزيع التكراري المتجمع الصاعد (distribution of up word combined frequency) = تكرار الفئة + تكرار الفئات السابقة.

Same as the previous example, calculate the up word combined frequency.

Actual limits for classes	Frequencies (f_i)	The up word combined frequency
Less than 24.5	0	فئة غير موجودة 0
24.5 – 29.5	7	7
29.5 – 34.5	19	26
34.5 – 39.5	14	40
39.5 – 44.5	7	47
44.5 – 49.5	3	The total number of data (n) = 50

اما لحساب التوزيع التكراري المتجمع النازل (distribution of down word combined frequency) = تكرار الفئة – تكرار الفئات السابقة.

For the seam example:

Actual limits for classes	Frequencies (f _i)	The down word combined frequency
24.5 – 29.5	7	50
29.5 – 34.5	19	43
34.5 – 39.5	14	24
39.5 – 44.5	7	10
44.5 – 49.5	3	3
More than 49.5	0	لأنها فئة غير موجودة 0

b- Distribution of Relative Frequency (p) : (التوزيع التكراري النسبي)

$$P = \frac{f}{n}$$

f: تكرار الفئة

n: مجموع التكرارات

ملاحظة مجموع التكرارات النسبية يجب ان يساوي 1.

c- Distribution of percentile Frequency (التوزيع التكراري المئوي):

Distribution of percentile Frequency = $P * 100\%$ التوزيع التكرار المئوي

For the seam example:

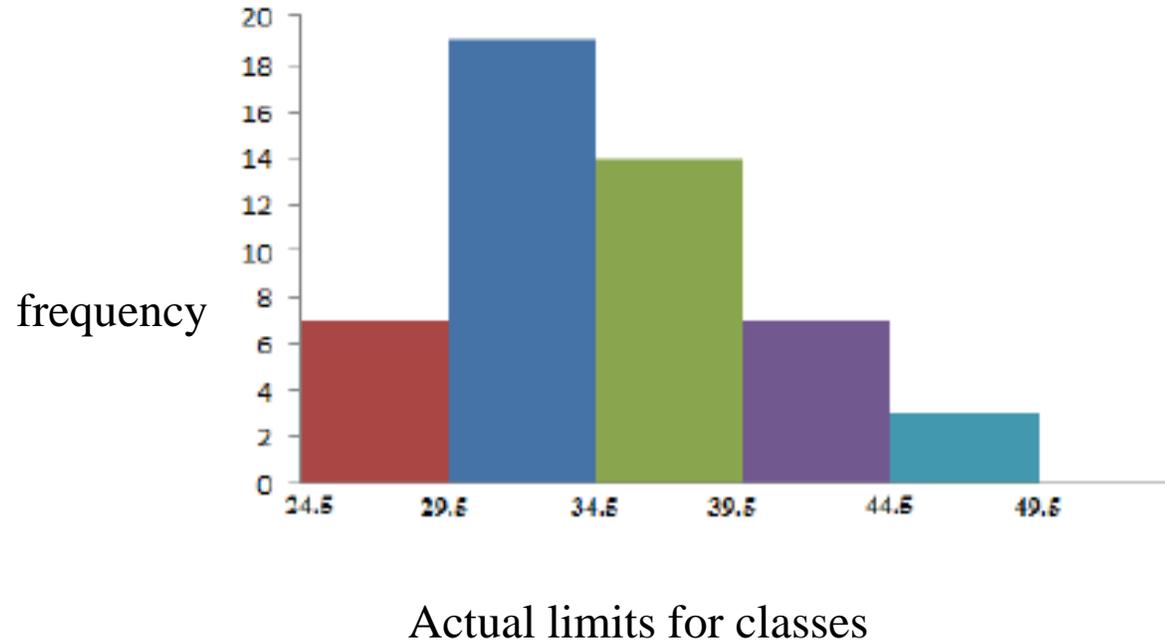
Actual limits for classes	Frequencies (f_i)	Relative Frequency	percentile Frequency
24.5 – 29.5	7	$7/50 = 0.14$	14
29.5 – 34.5	19	$19/50 = 0.38$	38
34.5 – 39.5	14	$14/50 = 0.28$	28
39.5 – 44.5	7	$7/50 = 0.14$	14
44.5 – 49.5	3	$3/50 = 0.06$	6
		$\Sigma = 1$	$\Sigma = 100$

Graphically represent the frequency distributions

تمثيل التوزيعات التكرارية بيانيا

1- طريقة المدرج التكراري (Frequency Histogram): حيث يتم تمثيل تكرار كل فئة بمستطيل حدود قاعدته يمثل الحدود الفعلية لتلك الفئة و ارتفاعه يتناسب مع تكرارها.

Returning to the previous example, the histogram distribution is represented by the frequency histogram method.



2- المضلع التكراري (Frequency polygon): و هو مضلع مغلق نحصل عليه بتصنيف الاضلاع العلوية للمستطيلات في المدرج التكراري و توصيل نقاط التصنيف مع بعضها و من ثم اغلاق المضلع مع المحور الافقي.

3- المنحني التكراري (Frequency Curve): و يرسم بتمديد خطوط المضلع التكراري لجعله منحنيا، و يستخدم في حالة البيانات الكبيرة من النوع المتصل مثل بيانات الزمن والوزن.

4- المضلع التكراري المتجمع (Accumulated Frequency Polygon): و يرسم بمد خطوط بين قيم التكرار المتجمع للفئات في مناطق نهايات الفئات الحقيقية.

Statistics

Lecture 5

Ex1:The following data shows the educational level of twenty female workers in a factory:

the educational level	uneducated	primary	Secondary 1	Secondary 2
The number	5	8	1	6

Represent the following using the pie-charts method.

Sol:

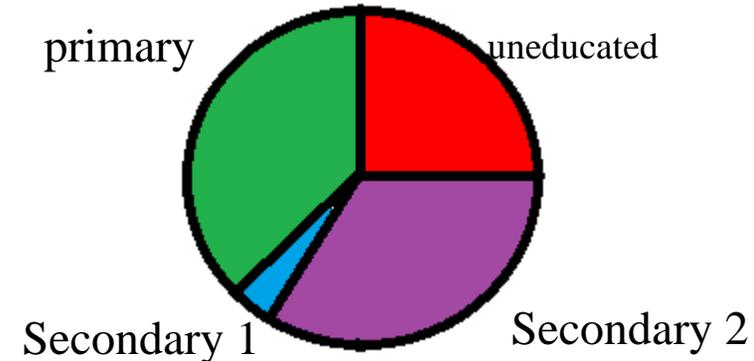
Total summation for team = $5+8+1+6=20$

Central angle for the uneducated = $(5/20)*360^\circ = 90^\circ$

Central angle for the primary = $(8/20)*360^\circ = 144^\circ$

Central angle for the Secondary 1 = $(1/20)*360^\circ = 18^\circ$

Central angle for the Secondary 2 = $(6/20)*360^\circ = 108^\circ$

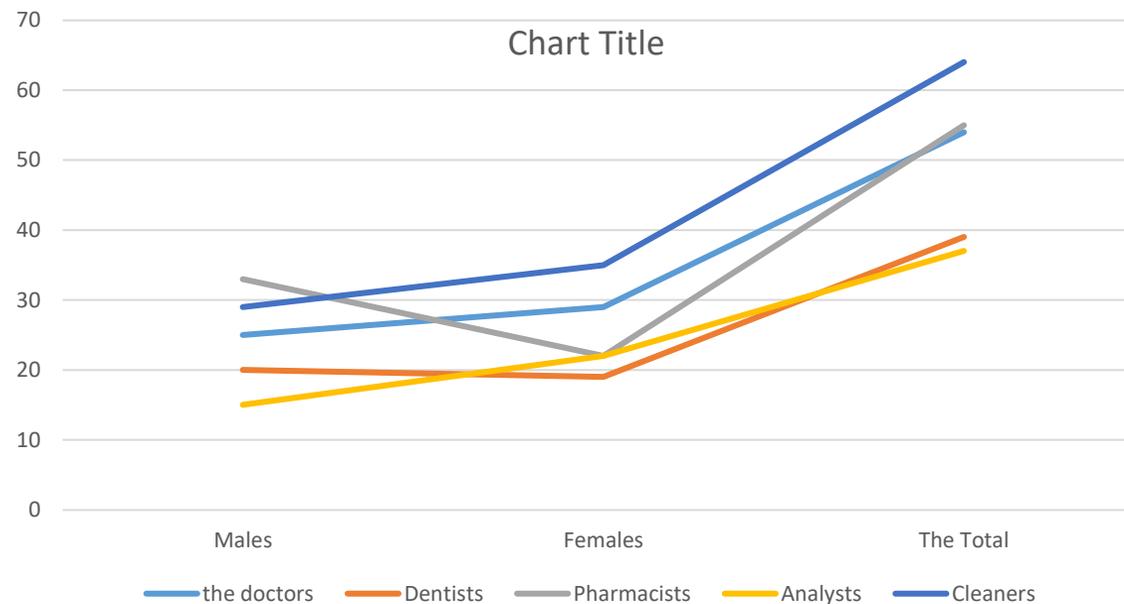


Ex2: The following data represent the number of male and female employees at Al-Hayat Hospital for all categories from doctors to cleaning workers:

Categories	the doctors	Dentists	Pharmacists	Analysts	Cleaners
Males	25	20	33	15	29
Females	29	19	22	22	35
The Total	54	39	55	37	64

Represent the following data in the table in the included line method.

Sol:



Ex3: In the process of building a building with bricks, which is transported from the factory to the work site by vehicles especially. The time required to transport each cargo was measured in minutes, and the results were as follows:

30	18	17	13	24	20	20	16	24	25
19	26	28	23	23	28	17	18	10	18

Arrange the following data in the frequency table with the classes.

Sol:

Assume that the number of classes= 7.

rang of the data = $30 - 10 = 20$

Class length = $\frac{\text{rang}}{\text{Number of classes}} = \frac{20}{7} = 2.85 = 3$

The minimum limit for the first class = 10

The actual minimum for the first class = $10 - 0.5 = 9.5$

The actual upper limit for the first class = $9.5 + 3 = 12.5$

The upper limit for the first class = $12.5 - 0.5 = 12$

The lower and upper limits of the rest of the classes are determined by adding the length of the category to each boundary.

Center of classes $(X_i) = \frac{10+12}{2} = 11$

Center of classes of the rest = Center of the previous category+ Class length

Collect of frequencies (The result must be equal to the original data count).

Class limits	Actual limits for classes	Centers of classes	Frequencies
10 - 12	9.5 – 12.5	11	1
13 – 15	12.5 – 15.5	14	1
16 – 18	15.5 – 18.5	17	6
19 – 21	18.5 – 21.5	20	3
22 – 24	21.5 – 24.5	23	4
25 – 27	24.5 – 27.5	26	2
28 - 30	27.5 – 30.5	19	3
			$\Sigma f_i = 20$

Ex4: The following data represents the number of customers of a commercial store in a city daily during two consecutive months:

(40 55 47 49 50 48 41 49 45 46 54 51 45 46 48 54 58 62 47 48 47 58 51 49 48 46 44 39 36 43 41 44 43 38 39 61 48 47 43 46 40 44 52 50 60 35 36 56 53 44 39 41 38 45 50 52 55 60 57 49)

Show this data in a frequency distribution table of 6 classes of equal length.

Sol:

Assume that the number of classes = 6.

rang of the data = $62 - 35 = 27$

Class length = $\frac{\text{rang}}{\text{Number of classes}} = \frac{27}{6} = 4.5 = 5$

The minimum limit for the first class = 35

The actual minimum for the first class = $35 - 0.5 = 34.5$

The actual upper limit for the first class = $34.5 + 5 = 39.5$

The upper limit for the first class = $39.5 - 0.5 = 39$

The lower and upper limits of the rest of the classes are determined by adding the length of the category to each boundary.

$$\text{Center of classes } (X_i) = \frac{35+39}{2} = 37$$

Center of classes of the rest = Center of the previous category + Class length

Collect of frequencies (The result must be equal to the original data count).

Class limits	Actual limits for classes	Centers of classes	Frequencies
35 - 39	34.5 – 39.5	37	8
40 – 44	39.5 – 44.5	42	12
45 – 49	44.5 – 49.5	47	20
50 – 54	49.5 – 54.5	52	10
55 – 59	54.5 – 59.5	57	6
60 - 64	59.5 – 64.5	62	4
			$\Sigma f_i = 60$

Ex5: For the example 4:

- 1- Calculate the up word combined frequency.
- 2- Calculate the down word combined frequency.
- 3- Calculate the distribution of percentile Frequency.

Sol:

1-

Actual limits for classes	Frequencies (f_i)	The up word combined frequency
Less than 34.5	0	0
34.5 – 39.5	8	8
39.5 – 44.5	12	20
44.5 – 49.5	20	40
49.5 – 54.5	10	50
54.5 – 59.5	6	56
59.5 – 64.5	4	The total number of data (n) = 60

2-

Actual limits for classes	Frequencies (f_i)	The down word combined frequency
34.5 – 39.5	8	60
39.5 – 44.5	12	52
44.5 – 49.5	20	40
49.5 – 54.5	10	20
54.5 – 59.5	6	10
59.5 – 64.5	4	4
More than 64.5	0	0

3-

Actual limits for classes	Frequencies (f_i)	Relative Frequency	percentile Frequency
34.5 – 39.5	8	$(8/60)=0.14$	14
39.5 – 44.5	12	$(12/60)=0.2$	20
44.5 – 49.5	20	$(20/60)=0.34$	34
49.5 – 54.5	10	$(10/60)=0.16$	16
54.5 – 59.5	6	$(6/60)=0.1$	10
59.5 – 64.5	4	$(4/60)=0.06$	6
		$\Sigma =1$	$\Sigma =100$

Statistics

Lecture 10

Measures of dispersion

In statistics, the measures of dispersion help to interpret the variability of data i.e. to know how much homogenous or heterogeneous the data is. In simple terms, it shows how squeezed or scattered the variable is.

في الإحصاء ، تساعد مقاييس التشتت في تفسير تباين البيانات ، أي معرفة مقدار البيانات المتجانسة أو غير المتجانسة. بعبارة بسيطة ، يوضح مدى ضغط أو تشتت المتغير.

Types of measures of dispersion:

1. **Range.** (المدى)
2. **Quartile deviation.** (الانحراف الربيعي)
3. **Average of the absolute deviations.** (متوسط الانحرافات المطلقة)
4. **Variance .** (التباين)
5. **Standard deviation.** (الانحراف المعياري)

In addition to **the coefficient of variation** that is used to compare two or more sets of data in terms of Dispersion.

بالإضافة إلى معامل الاختلاف الذي يستخدم لمقارنة مجموعتين أو أكثر من البيانات من حيث التشتت.

1. Range. (المدى)

The range is the difference between the largest and smallest value in the data, it is the period that In which Type equation here. the variable under study changes and is denoted by the symbol (R).

A. Non- classified data:

ويحسب المدى في حالة البيانات غير المبوبة كما يلي :

$R = \text{largest value} - \text{smallest value}$

$R = \text{اصغر قيمة} - \text{اكبر قيمة}$

B. Classified data:

أما في حالة البيانات المبوبة نعتبر أكبر قيمة هي الحد الأعلى للفئة الأخيرة وأصغر قيمة هي الحد الأدنى للفئة الأولى مع مراعاة أن تكون فئات الجدول مرتبة ترتيباً تصاعدياً ، أي يحسب المدى في حالة البيانات المبوبة كما يلي :

$R = \text{The upper limit for the last class} - \text{The minimum limit for the first class}$

$R = \text{الحد الادنى للفئة الاولى} - \text{الحد الاعلى للفئة الاخيرة}$

إذا كان المدى صغيراً فيعني ذلك أن البيانات منتشرة في فترة قصيرة أي قريبة من بعضها وتشتتها صغير ، أما إذا كان المدى كبيراً فيعني ذلك أن البيانات منتشرة في فترة طويلة أي متباعدة عن بعضها وتشتتها كبير.

Ex1: Using the term compare the dispersion of the three groups.

G1: 6, 6, 6, 6, 6.

G2: 4, 5, 6, 7, 8.

G3: 0, 1, 6, 10, 13.

Sol:

$$R(G1) = 6 - 6 = 0$$

$$R(G2) = 8 - 4 = 4$$

$$R(G3) = 13 - 0 = 13$$

بما أن مدى المجموعة الأولى يساوي صفرًا ، فيعني ذلك أنه لا يوجد اختلاف أو تشتت بين قيم هذه المجموعة ، أي أن كل القيم داخل هذه المجموعة متساوية.

ونلاحظ أن مدى المجموعة الثانية أصغر من مدى المجموعة الثالثة، وهذا يعني أن تشتت القيم داخل المجموعة الثانية أقل من تشتت القيم داخل المجموعة الثالثة.

H.W1: Compare the dispersion of the statistics subject's scores for three groups of students:

G1: 50, 62, 30, 70, 45, 80.

G2: 0, 20, 70, 55, 65, 82, 35, 77.

G3: 49, 54, 85, 50, 60, 63, 75, 65.

Ex2: Calculate the range for the following data showing weights in kilograms of 100 students:

Weight (kg)	Number of students
60 – 63	5
63 – 66	15
66 – 69	40
69 – 72	28
72 – 75	12

Sol:

$$\begin{aligned} R &= \text{The upper limit for the last class} - \text{The minimum limit for the first class} \\ &= 75 - 60 = 15. \end{aligned}$$

Range properties:

1. Range is a measure that is easy to calculate and simple in concept and significance.
2. It is concerned with only two values in the data and neglects the rest.
3. It depends in its calculation on the smallest and largest value only, and therefore it is considered a misleading measure, because it is when the largest or minor values, or both, are outliers, then the range is large while the group values are not divergent, for example if we have the following data:

55, 14, 11, 9, 8, 12, 10.

$$R = 55 - 8 = 47$$

فالمدى كبير ، ويشير إلى وجود تشتت كبير في المجموعة في حين أن القيم متقاربة ، ولذلك فهو مقياس مضلل ، وسبب ذلك هو اعتماده على القيم المتطرفة ، فلو حذفنا القيمة المتطرفة وهي القيمة 55 فنجد أن قيمة المدى تساوي :

$$R = 14 - 8 = 6$$

وهي قيمة صغيرة وواقعية .

2. Variance(التباين)

The variance is defined as the arithmetic mean of the squares of the deviations of values from their arithmetic mean (S^2).

A. Non- classified data:

- نحسب الوسط الحسابي للبيانات.
- نحسب انحراف كل قيمة عن الوسط الحسابي ،
- حيث : انحراف القيمة عن الوسط الحسابي = القيمة - الوسط الحسابي
- نوجد مربع انحراف كل قيمة عن الوسط الحسابي .
- نوجد مجموع مربعات انحرافات القيم عن الوسط الحسابي .
- نحسب قيمة التباين باستخدام القانون التالي :

$$S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n - 1}$$

Ex3: Calculate the variance for data representing 8 students' grades:

7, 6, 10, 9, 8, 5, 5, 6.

Sol:

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n} = \frac{56}{8} = 7$$

Value (X_i)	$(\bar{X} - X_i)$	$(\bar{X} - X_i)^2$
7	0	0
6	-1	1
10	3	9
9	2	4
8	1	1
5	-2	4
5	-2	4
6	-1	1
	Σ	24

$$S^2 = \frac{24}{7} = 3.43$$

B. Classified data:

في حالة البيانات المعروضة في جداول تكرارية نتبع الخطوات التالية :

- نحسب مراكز الفئات.
- نحسب الوسط الحسابي للتوزيع .
- نحسب انحراف كل مركز عن الوسط الحسابي .
- نوجد مربع انحراف كل مركز عن الوسط الحسابي .
- يضرب مربع انحراف كل مركز عن الوسط الحسابي لفئة في تكرار هذه الفئة .
- نجمع حواصل الضرب التي تحصلنا عليها في الخطوة السابقة .
- نحسب قيمة التباين باستخدام القانون التالي :

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

حيث : x_i : مركز الفئة ، f_i : تكرار الفئة ، \bar{x} : الوسط الحسابي

Ex4: Calculate the variance for the following data, which represent the distribution of 200 workers by overtime spend on the job at the factory every month:

Overtime (h)	Number of Workers
0 – 10	20
10 – 20	80
20 – 30	50
30 – 40	40
40 – 50	10
Σ	200

Sol:

$$S^2 = \frac{\sum_{i=1}^k f_i (x_i - \bar{X})^2}{\sum_{i=1}^k f_i - 1}$$

class	f_i	X_i	$X_i f_i$	$(X_i - \bar{x})$	$(X_i - \bar{x})^2$	$f_i(X_i - \bar{x})^2$
0 – 10	20	5	100	-17	289	5780
10 – 20	80	15	1200	-7	49	3920
20 – 30	50	25	1250	3	9	450
30 – 40	40	35	1400	13	169	6760
40 – 50	10	45	450	23	529	5290
Σ	200		4400			22200

$$\bar{X} = \frac{\sum_{i=1}^k x_i f_i}{\sum_{i=1}^k f_i} = \frac{4400}{200} = 22$$

$$S^2 = \frac{\sum_{i=1}^k f_i (x_i - \bar{x})^2}{\sum_{i=1}^k f_i - 1}$$
$$= \frac{22200}{200-1} = \frac{22200}{199} = 111.56 \text{ h}^2$$

ونلاحظ أن التباين وحداته هي مربع وحدات القياس الأصلية، وكثيراً ما تكون غير ذات معنى فمثلاً في هذا المثال التباين = 111.56 h² فهنا ساعة تربيع ليس لها أي معنى، وهذا هو عيب التباين ؛ ولذلك نُرجع الوحدات إلى أصلها بأخذ الجذر التربيعي للتباين ، ويسمى المقياس الجديد بالانحراف المعياري .

3. Standard deviation (الانحراف المعياري)

Standard deviation is the measure of dispersion of a set of data from its mean. It measures the absolute variability of a distribution; the higher the dispersion or variability, the greater is the standard deviation and greater will be the magnitude of the deviation of the value from their mean.

It is the positive square root of the mean of the squares of deviations of values from their mean arithmetic, that is, it is the positive square root of the variance, and it is denoted by a symbol (S).

هو الجذر التربيعي الموجب للوسط الحسابي لمربعات انحرافات القيم عن وسطها الحسابي ، أي: هو الجذر التربيعي الموجب للتباين.

A. Non- classified data:

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

B. Classified data:

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

For Ex :

$$S = \sqrt{3.43} = 1.85$$

For Ex :

$$S = \sqrt{111.56} = 10.56$$

إذا لم يكن الوسط الحسابي عدداً صحيحاً فإن حساب التباين ومن ثم الانحراف المعياري باستخدام الصيغ السابقة الذكر ، يصبح أمراً غير سهلٍ ، ولذلك اشتقت من الصيغة الأساسية للتباين والتي تعتمد على الانحرافات عن الوسط الحسابي ، صيغة أخرى تعتمد على القيم مباشرة ، وذلك لتسهيل العمليات الحسابية ، وبالطبع الصيغتان تعطيان نفس النتيجة تماماً ، وصيغ التباين التي تعتمد على القيم مباشرة هي :

A. Non- classified data:

$$S^2 = \frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i \right)^2}{n} \right]$$

B. classified data:

$$S^2 = \frac{1}{\sum_{i=1}^k f_i - 1} \left[\sum_{i=1}^k x_i^2 f_i - \frac{\left(\sum_{i=1}^k x_i f_i \right)^2}{\sum_{i=1}^k f_i} \right]$$

حيث : X_i : مركز الفئة ، f_i : تكرار الفئة ،
يجب الانتباه بأن هناك فرقاً بين المقدار $\sum_{i=1}^k x_i^2 f_i$ والمقدار $\left(\sum_{i=1}^k x_i f_i \right)^2$.

Ex5: Calculate the variance and standard deviation of the data mentioned in an example 3, so using the second variance formula (direct values formula).

Sol:

$$S^2 = \frac{1}{n - 1} \left[\sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i \right)^2}{n} \right]$$

X_i	7	6	10	9	8	5	5	6	$\Sigma X_i = 56$
X_i^2	49	36	100	81	64	25	25	36	$\Sigma X_i^2 = 416$

$$\begin{aligned}
 S^2 &= \frac{1}{8-1} \left[416 - \frac{(56)^2}{8} \right] \\
 &= \frac{1}{7} [416 - 392] = \frac{24}{7} = 3.43
 \end{aligned}$$

$$S = \sqrt{3.43} = 1.85$$

نلاحظ انها نفس النتيجة في المثال السابق.

Ex6: Calculate the variance and standard deviation of the data mentioned in an example 4, so using the second variance formula (direct values formula).

Sol:

$$S^2 = \frac{1}{\sum_{i=1}^k f_i - 1} \left[\sum_{i=1}^k X_i^2 f_i - \frac{\left(\sum_{i=1}^k X_i f_i \right)^2}{\sum_{i=1}^k f_i} \right]$$

class	f_i	X_i	$X_i f_i$	X_i^2	$f_i X_i^2$
0 – 10	20	5	100	25	500
10 – 20	80	15	1200	225	18000
20 – 30	50	25	1250	625	31250
30 – 40	40	35	1400	1225	49000
40 – 50	10	45	450	2025	20250
Σ	200		4400		119000

$$S^2 = \frac{1}{200-1} \left[119000 - \frac{(4400)^2}{200} \right]$$
$$= \frac{1}{199} [119000 - 96800] = 111.56$$

$$S = \sqrt{111.56} = 10.56$$

نلاحظ انها نفس النتيجة التي حصلنا عليها في المثال الرابع.

Properties of Standard Deviation:

- Standard deviation is the most important and widely used measure of dispersion.
- Enter into his calculation all the observed values without neglecting any value.
- It is characterized by its ability to treat algebra.
- It cannot be calculated for open frequency distributions.

Coefficient of variation (معامل الاختلاف):

كل مقاييس التشتت السابقة تعتمد على وحدات القياس، وبالتالي لا يمكن استعمالها لمقارنة تشتت توزيعين مختلفين في وحدات القياس كمقارنة تشتت الأطوال بتشتت الأوزان مثلاً ، ولذلك يجب التعامل مع مقياس نسبي لا يعتمد على الوحدات المستعملة ويسمى هذا المقياس معامل الاختلاف ويحسب كما يلي :

$$\text{Coefficient of variation} = \frac{S}{\bar{x}} * 100$$

$$100 \times \frac{\text{الانحراف المعياري}}{\text{الوسط الحسابي}} = \text{معامل الاختلاف}$$

وبالطبع كلما زاد التشتت أخذ معامل الاختلاف نسبة أكبر .

Ex7: Compare the dispersion of the students' lengths and their weights, using the following data, which represent the lengths and weights 100 students.

Lengths (cm)	Number of students
110 – 120	12
120 – 130	15
130 – 140	25
140 – 150	30
150 – 160	10
160 - 170	8

Weights (Kg)	Number of students
60 – 63	6
63 – 66	16
66 – 69	40
69 – 72	28
72 – 75	10

Sol:

نحسب معامل الاختلاف لكل من التوزيعين ثم نقوم بالمقارنة ، ولحساب معامل الاختلاف لكل توزيع يلزمنا حساب الوسط الحسابي والانحراف المعياري لكل منهما ، وذلك كما يلي :

A. Distribution of lengths:

نوضح العمليات الحسابية التي تلزمنا لحساب الوسط الحسابي والانحراف المعياري في الجدول التالي:

class	f_i	X_i	$X_i f_i$	$(X_i - \bar{x})$	$f_i(X_i - \bar{x})^2$
110 – 120	12	115	1380	-23.5	6627
120 – 130	15	125	1875	-13.5	2733.75
130 – 140	25	135	3375	-3.5	306.25
140 – 150	30	145	4350	6.5	1267.5
150 – 160	10	155	1550	16.5	2722.5
160 - 170	8	165	1320	26.5	5618
Σ	100		13850		19275

$$\bar{x} = \frac{\sum_{i=1}^k x_i f_i}{\sum_{i=1}^k f_i} = \frac{13850}{100} = 138.50$$

$$S^2 = \frac{1}{\sum_{i=1}^k f_i - 1} \left[\sum_{i=1}^k f_i (x_i - \bar{x})^2 \right]$$
$$= \frac{1}{100-1} 19275 = 194.7$$

$$S = \sqrt{194.7} = 13.95$$

$$\text{Coefficient of variation for lengths} = \frac{13.95}{138.5} * 100 = 10.07\%$$

B. Distribution of weights:

الجدول التالي يوضح لنا العمليات الحسابية اللازمة للحصول على الوسط الحسابي والانحراف المعياري للاوزان:

class	f_i	X_i	$X_i f_i$	$(X_i - \bar{x})$	$f_i(X_i - \bar{x})^2$
60 – 63	6	61.5	369	-6.6	261.36
63 – 66	16	64.5	1032	-3.6	207.36
66 – 69	40	67.5	2700	-0.6	14.4
69 – 72	28	70.5	1974	2.4	161.28
72 – 75	10	73.5	735	5.4	291.6
Σ	100		6810		936

$$\bar{x} = \frac{\sum_{i=1}^k x_i f_i}{\sum_{i=1}^k f_i} = 68.1$$

$$S^2 = \frac{1}{\sum_{i=1}^k f_i - 1} \left[\sum_{i=1}^k f_i (x_i - \bar{X})^2 \right]$$
$$= \frac{936}{99} = 9.45$$

$$S = \sqrt{9.45} = 3.07$$

$$\text{Coefficient of variation for weights} = \frac{3.07}{68.1} * 100 = 4.51\%$$

بما أن معامل الاختلاف للأطوال أكبر من معامل الاختلاف للأوزان ، إذن تشتت الأطوال أكبر من تشتت الأوزان.

H.W 2: If we had two sets of data:

set 1: 21, 15, 23, 16, 14, 12, 25.

set 2: 30, 42, 64, 60, 28, 70, 82, 56.

A- For each of these two groups, calculate the following:

Range, Variance, Standard deviation.

B- Compare the dispersion of the two groups using the coefficient of variation.

H.W 3: If you know the coefficient of variation = 25% and that the Variance = 25, that for the following data:

10, 18, 25, 15, X. Calculate the missing value X.

Statistics

Lecture 12

Probability Distribution

التوزيعات الاحتمالية

The benefit of studying probability distributions from the valuable implementation success is in a major field:

1. The possibility of determining or calculating a probability in certain evidence circumstances.
2. Appropriateness of the probability distribution of some of the data collected in the field. (Therefore, it is possible to extract information from this data after making a comparison with these distributions are compatible with it).

1. امكانية تحديد او احتساب احتمال ما في ظروف بيانات معينة.

2. امكانية ملائمة التوزيع الاحتمالي لبعض البيانات التي تم جمعها ميدانيا. وبالتالي امكانية استخراج المعلومات من هذه البيانات بعد عمل مقارنة مع هذه التوزيعات المتلائمة معها.

المتغيرات العشوائية Random variables

إن بعض الكميات تكون قيمتها العددية غير محددة تحديداً نهائياً، وإنما تتغير بتغير عوامل عشوائية ومن الامثلة على هذه الكميات اختلاف مقاومات الانضغاط للمكعبات المأخوذة من نفس الخلطة ومعالجتها بنفس الظروف او مقاومات الشد لقضبان حديد التسليح المصنوعة من نفس الفولاذ وبنفس الظروف. إن هذا الاختلاف بالرغم من تساوي فرص الحصول على توافق ظاهرياً هو بسبب ظاهرة التشتت وسببها عوامل عشوائية لم يكن بالإمكان حسابها او التنبؤ بها مسبقاً. لذا يمكن تعريف الكميات العشوائية: وهي الكميات التي تأخذ قيماً مختلفة في العمليات أو المحاولات المختلفة وتبقى مختلفة مهما حاولنا جعل ظروف وشروط حدوثها متجانسة تبعاً لتغيرات عشوائية خارجة عن الحساب. **إن اهم ما يجب معلومته بما يخص المتغيرات العشوائية هو جميع القيم العددية التي يمكن أن تأخذها تلك المتغيرات، بالإضافة الى معرفة الاحتمالات المقابلة لتلك القيم لحالة المتغيرات المدروسة.** إن معلومية احتمال القيم الممكنة للكميات العشوائية المدروسة يكفي لحل أي سؤال يرتبط بهذه الكمية العشوائية وتبيان فعاليتها. وتقسم المتغيرات العشوائية الى نوعين:

- أ- المتغيرات العشوائية المتقطعة او المنفصلة : ويكون فيها المتغير العشوائي يأخذ أي قيمة من قيم مجموعة متقطعة أو قابلة للعد مثل متغير عدد الطلبة الذكور في احدى المدارس.
- ب- المتغيرات العشوائية المتصلة : وهي المتغيرات العشوائية التي من الممكن أن تأخذ أي قيمة داخل مدى (مجال) معين من القيم مثل مقاييس درجة الحرارة أو الوزن أو الطول.

Depending on the type of the random variable, the probability distributions can be divided into two types:

بالاعتماد على انواع المتغير العشوائي ، يمكن تقسيم التوزيعات الاحتمالية الى نوعين:

A. Discrete Probability Distributions ((توزيعات احتمالية متقطعة (منفصلة))

- Binomial Dist. (توزيع ذي الحدين)
- Poisson Dist.(توزيع بواسون)

B. Continuous Probability Distributions (توزيعات احتمالية متصلة)

- Normal Distribution (التوزيع الطبيعي)
- T-Distribution (توزيع تي)
- Chi-Square Distribution (توزيع كاي تربيع)

A. Discrete Probability Distributions (توزيعات احتمالية متقطعة (منفصلة))

وهي التوزيعات الاحتمالية (دوال كتلة الاحتمال) لمتغيرات عشوائية متقطعة.

تعريف دالة كتلة الاحتمال للمتغير العشوائي المتقطع: إذا كان X متغيراً عشوائياً متقطعاً بحيث إن مجموعة القيم الممكنة له هي

$$X(S) = \{X_1, X_2, \dots, X_n\}$$

فإن دالة كتلة الاحتمال ويرمز لها $f(X)$ تأخذ القيم التالية:

$$f(X) = \begin{cases} P(X = X_i); & X_i \in X(S) \\ 0 & ; X_i \notin X(S) \end{cases}$$

ويجب ان تحقق دالة كتلة الاحتمال الشروط التالية:

- i. ان قيمة دالة كثافة الاحتمال عند قيمة X ما يجب أن لا تزيد عن 1 ولا تقل عن صفر.
- ii. ان مجموع قيم دوال كثافة الاحتمال لكل قيم المتغير X المتوقعة (ضمن القيم الممكنة ل X ان تأخذها أو فضاء العينة) يجب أن يساوي "1".

• Binomial Dist. (توزيع ذي الحدين)

وهو من اهم التوزيعات المتقطعة المستعملة في مجالات التطبيقات الهندسية لمتغيرات عشوائية منفصلة. حيث يمكن الحكم على نتائج تجربة ما ب "نعم" أو "لا". ويستخدم عند تحقيق الشروط التالية:

- ان نتيجة أي محاولة لتجربة ما تكون اما النجاح أو الفشل.
- كل محاولة نتيجتها مستقلة عن نتيجة المحاولة الاخرى. (أي لا تؤثر عليها مثل عملية السحب بإرجاع).
- احتمال النجاح في كل محاولة ثابت وليكن (P) وبالتالي فإن احتمال الفشل ثابت أيضاً وهو (q) ويحسب من (q=1-P).
- تجري التجربة بعدد معين من المحاولات أو المرات أو بمعنى آخر (يكون هنالك عدد n من المحاولات).

وعند اجراء تجربة ما وكان عدد المحاولات = n ، عدد النجاحات أو عدد مرات حصول الحادثة = x ، احتمال النجاح للتوزيع = P ، واحتمال الفشل للتوزيع = q فإن الاقتران الاحتمالي لمتغير ذي الحدين ويرمز له P(x;n,p) ويحسب كما يلي:

$$P(x; n, p) = C_x^n \cdot p^x \cdot q^{(n-x)} \quad ; X=0,1,2,\dots,n$$

ملاحظة: تم أخذ قيمة الصفر لأنه يمثل احتمال عدم الحصول او احتمال ولا مرة يحصل حادث النجاح.

ملاحظة: اذا كان X متغيراً ذا حدين (خاضع لنظرية ذات الحدين) فإن:

$$M = n \cdot p$$

$$S^2 = n \cdot p \cdot q$$

i. الوسط او التوقع لقيمة X الممكنة هو

ii. التباين لقيمة X هو

مثال(1): رُميت زهرة نرد منتظمة 5 مرات ، ما احتمال عدم ظهور الرقم 6 فيها؟ ما احتمال ظهور الرقم 6 ثلاث مرات؟

$$s = \{1,2,3,4,5,6\}$$

الحل: ان فضاء العينة لتجربة رمي زهرة نرد واحدة منتظمة لمرة واحدة هو

$$^{\wedge} \text{ احتمال ظهور الرقم 6 } P(6) = \frac{1}{6} \rightarrow ^{\wedge} \text{ احتمال عدم ظهور الرقم 6 } q(6) = 1 - \frac{1}{6} = \frac{5}{6}$$

$$^{\wedge} P(6) = \frac{1}{6}, q(6) = \frac{5}{6}, n=5, x=0,1,2,3,4,5$$

$$P\left(x; 5, \frac{1}{6}\right) = C_x^5 \cdot \left(\frac{1}{6}\right)^x \cdot \left(\frac{5}{6}\right)^{(5-x)}, x=0,1,2,3,4,5$$

وهي تمثل الصيغة العامة لهذه التجربة

$$^{\wedge} \text{ احتمال عدم ظهور الرقم 6 } P\left(x = 0; 5, \frac{1}{6}\right) = C_0^5 \cdot \left(\frac{1}{6}\right)^0 \cdot \left(\frac{5}{6}\right)^{(5-0)} = 1 \times 1 \times \left(\frac{5}{6}\right)^5 = 0.40187$$

$$^{\wedge} \text{ احتمال عدم ظهور الرقم 6 ثلاث مرات } P\left(x = 3; 5, \frac{1}{6}\right) = C_3^5 \cdot \left(\frac{1}{6}\right)^3 \cdot \left(\frac{5}{6}\right)^{(5-3)} = 10 \times \frac{1}{216} \times \left(\frac{25}{36}\right) = 0.03215$$

$$M = n \cdot p = 5 \cdot \frac{1}{6} = \frac{5}{6}$$

ولحساب القيمة المتوقعة للمتغير X نحسب الوسط

$$S^2 = n \cdot p \cdot q = 5 \cdot \frac{1}{6} \cdot \frac{5}{6} = \frac{25}{36}$$

ولحساب قيمة التباين المتوقعة للمتغير X نحسب التباين

مثال(2): يُراد توزيع عدد من الكراسي الخاصة بالطلبة الأعسرين (مستخدمي اليد اليسرى) في كل قاعة دراسية ، فإذا عُلم ان 10% من جميع الطلبة هم أعسرين فما هو احتمال إن صفاً من 20 طالب: (أ)- لا يحتوي على أي طالب أعسر؟ (ب)- اثنان منهم أعسران؟ (ج)- أكثر من اثنين منهم أعسرين؟

الحل: نفرض ان المتغير X يمثل عدد الطلبة الأعسرين

$$x=0,1,2,\dots,20, \quad n=20$$

$$q(L)=1 - 0.1=0.9 \text{ احتمال أن لا يكون الطالب أعسر} \quad , \quad P(L)=0.1 \text{ احتمال أن يكون الطالب أعسر}$$

من المعلومات أعلاه نلاحظ تحقق شروط نظرية ذي الحدين، لذا فإن الصيغة العامة للتجربة هي

$$P(x; 20,0.1) = C_x^{20} \cdot (0.1)^x \cdot (0.9)^{(20-x)}, \quad x=0,1,2,\dots,20$$

$$P(x = 0; 20,0.1) = C_0^{20} \cdot (0.1)^0 \cdot (0.9)^{(20-0)} = 1 \times 1 \times 0.12157 = 0.12157 \text{ (أ)}$$

$$P(x = 2; 20,0.1) = C_2^{20} \cdot (0.1)^2 \cdot (0.9)^{(20-2)} = 190 \times 0.01 \times 0.15009$$

$$=0.28518$$

• Poisson Dist.(توزيع بواسون)

يعتبر توزيع بواسون من التوزيعات الاحتمالية المتقطعة التي تستخدم للحوادث نادرة الحصول أو الوقوع، أو المتغيرات التي تحدث في أزمنة عشوائية معلومة. مثل عدد الأخطاء المطبعية في صفحة كتاب ما، عدد الزلازل السنوية الحاصلة في موقع ما، حدوث فيضان بتصريف معين في نهر ما... الخ.

وتسمى التجارب التي تعطي عدد نجاحات معين في فترة زمنية أو فترة قياس كالطول أو الحجم أو المساحة... الخ بالتجربة البواسونية، وهي تجربة احصائية تمتلك الخصائص التالية:

- نتائج التجربة مُعرفة بالنجاح أو الفشل والتجارب مستقلة الواحدة عن الأخرى.
 - الاحتمال ثابت وذو قيمة صغيرة بشكل عام (تقترب من 0) وإن عدد المحاولات n يكون كبيراً عادةً (يقترب من ∞) ولكن المهم ان مقدار حاصل الضرب بينهما يكون مقدار ثابت يساوي المتوسط أو المعدل (λ) بحيث أن ($\lambda > 0$).
- $$\lambda = n \cdot P$$

ومن خصائص توزيع بواسون الاحتمالي:

- المعدل أو المتوسط للتوزيع ($\lambda = n \cdot P$) ، والتباين يساوي (λ)
- الانحراف المعياري يساوي ($\sqrt{\lambda}$)

$$P(x, \lambda) = \frac{e^{-\lambda} \cdot \lambda^x}{x!}, \quad x=0,1,2,\dots$$

ان الاقتران الاحتمالي لتوزيع بواسون هو:

حيث ان:

λ : معدل عدد النجاحات (حصول الحدث) ضمن الفترة الزمنية أو فترة القياس (الطول، المساحة، الحجم) وغيرها.

x : عدد النجاحات الفعلي أو عدد مرات القياس المطلوبة وهي تمثل متغير بواسون العشوائي.

e : أساس اللوغاريتم الطبيعي وتساوي تقريباً ($e=2.718281828$).

مثال(1): معدل عدد حوادث السيارات على طريق صحراوي ما هو 5 (حادث\أسبوع). (أ) ما احتمال عدم حدوث أي حادث على ذلك الطريق في اسبوع معين؟
 (ب) ما احتمال حدوث 4 حوادث أو أقل في اسبوع معين؟ (ج) ما احتمال أن يزيد عدد الحوادث عن حادثين خلال أسبوعين؟

$$\lambda=5, x=0,1,2,\dots, P(x, \lambda) = \frac{e^{-\lambda} \cdot \lambda^x}{x!}$$

الحل:

$$P(x, \lambda = 5) = \frac{e^{-5} \cdot 5^x}{x!}$$

الصيغة الأساسية لحل المسألة

$$(أ) \text{ عدم حصول أي حادث } P(x = 0, \lambda = 5) = \frac{e^{-5} \cdot 5^0}{0!} = 0.00674$$

$$(ب) \text{ حدوث 4 حوادث أو أقل } P(x \leq 4, \lambda = 5) = P(0,5) + P(1,5) + P(2,5) + P(3,5) + P(4,5)$$

$$= 0.00674 + \frac{e^{-5} \cdot 5^1}{1!} + \frac{e^{-5} \cdot 5^2}{2!} + \frac{e^{-5} \cdot 5^3}{3!} + \frac{e^{-5} \cdot 5^4}{4!} = 0.4405$$

زيادة عدد الحوادث عن 2 خلال إسبوعين (هنا نلاحظ أن معدل حصول الحوادث اختلف بدلاً من الفترة اسبوع الى اسبوعين فهنا ستتغير الصيغة (ت)
 الأساسية للسؤال حسب المعدل الجديد).

$$\lambda = 5(\text{اسبوع\حادث}) \times 2 = 10(\text{اسبوعين\حادث})$$

$$^* P(x, \lambda = 10) = \frac{e^{-10} \cdot 10^x}{x!}, x = 0,1,2, \dots$$

$$P(x > 2, \lambda = 10) = 1 - \{P(0,10) + P(1,10) + P(2,10)\} = 1 - \left\{ \frac{e^{-10} \cdot 10^0}{0!} + \frac{e^{-10} \cdot 10^1}{1!} + \frac{e^{-10} \cdot 10^2}{2!} \right\} = 0.99723$$

Approximation of the binomial distribution of the Poisson distribution

تقريب توزيع ثنائي الحدين لتوزيع بواسون

إذا اعتبرنا n او عدد المحاولات كبيرة بشكل كافٍ بحيث $(n \geq 50)$ و احتمال النجاح P صغير $(P < 0.1)$ بحيث يبقى حاصل ضرب قيمتهما $(n.P)$ مقدار ثابت وليكن λ ، عندها يمكن تقريب توزيع ذي الحدين الى توزيع بواسون بالمقدار $(\lambda = n.P)$.

مثال (2): في اختبار مُضاف مُحسن للخلطات الخرسانية، اذا افترضنا انه ضار للخلطة باحتمال 0.001 وانه تم استخدامه فعليا في انتاج 2000 خلطة خرسانية، استخدم تقريب بواسون للإجابة على: (أ) ما احتمال ان تتضرر أكثر من خلطتين نتيجة استخدام ذلك المُضاف؟ (ب) ما احتمال تضرر خلطتان على الأقل؟

Solution: $n=2000$, $P=0.001$ (احتمال حصول الضرر)

نلاحظ ان عدد المحاولات كبير جداً $(n > 50)$ واحتمال حصول الحدث P أقل من 0.1 لهذا يمكن استخدام تقريب بواسون بذي الحدين أو بمعنى آخر

$$\lambda = n.P = 0.001 \times 2000 = 2$$

$$(أ) P(x > 2, \lambda = 2) = 1 - \{P(0,2) + P(1,2) + P(2,2)\} = 1 - \left\{ \frac{e^{-2} \cdot 2^0}{0!} + \frac{e^{-2} \cdot 2^1}{1!} + \frac{e^{-2} \cdot 2^2}{2!} \right\} = 0.32332$$

$$(ب) P(x \geq 2, \lambda = 2) = 1 - P(x < 2, \lambda = 2) = 1 - \{P(0,2) + P(1,2)\} = 0.59399$$

Not:

الفرق الأساسي بين توزيعي ذي الحدين وتوزيع بواسون هو إن n أو عدد التجارب الكلي في توزيع ذي الحدين يكون معلوماً ولهذا فإن قيم المتغير العشوائي x تكون $(x=0,1,2,\dots,n)$ أما في توزيع بواسون فإن n أو عدد التجارب غير معلوم أو مبهم وفي هذه الحالة فإن المتغير العشوائي x يأخذ المدى $(x=0,1,2,\dots)$. ولهذا في حالة إحتساب احتمالات أكبر من قيمة ما و لتكن a في توزيع بواسون فإننا نضطر طرح احتمال تلك القيمة والتي أقل منها من (1) وهو الأحتمال الكلي لفضاء العينة أو بمعنى آخر $P(x > a, \lambda) = 1 - P(x \leq a, \lambda)$ وذلك لأننا لا نعرف إلى أي حد سينتهي إليه المتغير العشوائي x .

مثال(2): إذا كان 3% من انتاج أحد معامل الطابوق غير مطابق للمواصفات القياسية وقد تم سحب عينة من الانتاج بشكل عشوائي. ما هو احتمال بأن العينة تحتوي بالضبط على طابوقتين فاشلتين؟ استخدم توزيعين مختلفين للحل ثم قارن النتيجةين؟

الحل: (أ) باستخدام توزيع ذي الحدين ($P=0.03$, $q=0.97$, $n=100$)

$$\therefore P(x; 100,0.03) = C_x^{100} \cdot (0.03)^x \cdot (0.97)^{100-x} \quad x = 0,1, \dots, 100$$

$$\therefore P(x = 2; 100,0.03) = \frac{100!}{2! \times (100 - 2)!} \times (0.03)^2 \times (0.97)^{100-2} = 4950 \times 0.0009 \times 0.0505 = 0.225153$$

(ب) باستخدام توزيع بواسون ($\lambda = n \cdot P = 0.03 \times 100 = 3$)

$$P(x, \lambda = 3) = \frac{e^{-3} \cdot 3^x}{x!}, \quad x = 0,1,2, \dots$$

$$\therefore P(x = 2, \lambda = 3) = \frac{e^{-3} \cdot 3^2}{2!} = 0.22404$$

مثال(3): اذا كانت احدى البدالات تتلقى مكالمات هاتفية بمعدل مكالمتين لكل دقيقة. أوجد احتمال أن تتلقى البدالة 3 مكالمات على الأقل؟

Solution: ($\lambda=2$ call/min)

$$P(x, 2) = \frac{e^{-2} \cdot 2^x}{x!}, \quad x = 0, 1, 2, \dots$$

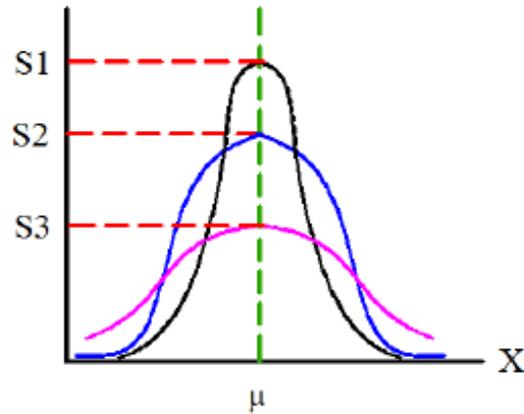
$$P(x \geq 3, \lambda = 2) = 1 - \left\{ \frac{e^{-2} \cdot 2^0}{0!} + \frac{e^{-2} \cdot 2^1}{1!} + \frac{e^{-2} \cdot 2^2}{2!} \right\} = 1 - (0.67667) = 0.3233$$

التوزيعات الاحتمالية المتصلة (Continuous Probability Distributions)

1-1- التوزيع الطبيعي (Normal Distribution):

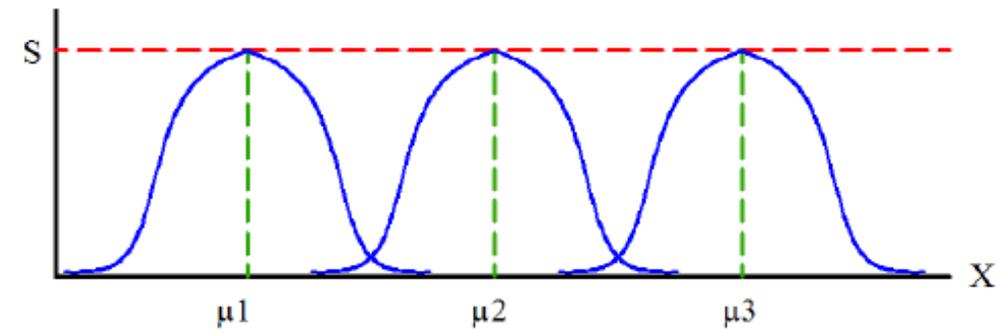
وهو من أشهر التوزيعات الاحتمالية المتصلة وأكثرها استخداماً وتطبيقاً ومن الأمثلة عليه التوزيعات البايومترية (توزيعات الطول والوزن) وتوزيعات اخطاء المشاهدات (الفروق بين القيم الحقيقية والقيم المُشاهدة) وتم تطوير هذا التوزيع من قبل العالم غاوس سنة 1809 ولهذا يطلق عليه في بعض الأحيان اسم (توزيع غاوس).

ان منحنى التوزيع الطبيعي الذي يأخذ شكل الجرس هو منحنى توزيع نظري لدالة كثافة (كتلة الاحتمال) للمتغير العشوائي الذي يخضع لهذا التوزيع. هنالك ما لانهاية من منحنيات التوزيع الطبيعي التي قد تختلف عن بعضها البعض حسب قيمة كل من الوسط الحسابي (μ) والانحراف المعياري (S)، فمثلاً قد تتفق المنحنيات بالانحراف المعياري لكنها تختلف بالوسط الحسابي **كما في الشكل (1)** أو بالعكس قد تختلف بالانحراف المعياري وتتفق بالوسط **كما في الشكل (2)**.



شكل (2)

منحنيات توزيع طبيعي لها نفس الوسط الحسابي مع اختلاف الانحراف المعياري



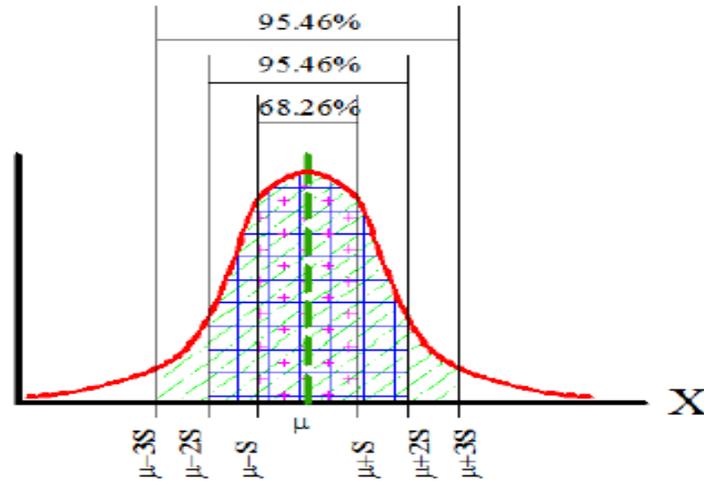
شكل (1)

منحنيات توزيع طبيعي لها نفس الانحراف المعياري مع اختلاف الوسط الحسابي

خصائص التوزيع الطبيعي:

- منحنى التوزيع الطبيعي متماثل حول الوسط μ (إتوائه يساوي صفراً) ، وشكله بشكل الجرس.
- له قيمة منوالية (منوال) واحد ينطبق على الوسط μ ويقترب طرفاه من قيمتي $(-\infty, \infty)$ دون أن تقطع المحور الافقي.
- قيمة الانحراف المعياري تحدد عرض المنحنى الطبيعي، مثلاً قيمة انحراف معياري أعلى تنتج منحنى أعرض وأكثر تسطحاً.
- الاحتمالات للمتغير الطبيعي العشوائي تُعطى بالمساحة تحت منحنى التوزيع الطبيعي. والمساحة الكلية تحت منحنى التوزيع الطبيعي للفترة $(-\infty \leq x \leq +\infty)$ تساوي (1) و تتوزع مناصفةً 0.5 للنصف الأيمن و 0.5 للنصف الأيسر.
- قيمة الاحتمال عند قيمة مفردة محددة للمتغير العشوائي الطبيعي x تساوي (صفر) ولأي قيمة وذلك لان المساحة تحت المنحنى والمحددة بقيمة واحدة تكون عبارة عن خط مستقيم ولهذا فلهذا فلحساب الاحتمال يجب ان تحدد فترة ما للمتغير العشوائي ومن خلالها تحسب المساحة تحت المنحنى ضمن هذه الفترة والتي سوف تساوي قيمة الاحتمال ضمن تلك الفترة.

$$P(x=0) = 0 , P(x=1) = 0 , P(x=100) = 0 , P(0 < x < a) = ? , P(x > a) = ?$$



المساحة تحت المنحنى الطبيعي تتوزع كما يلي (ثابتة لجميع منحنيات التوزيع الطبيعي)

- 1- 68.26% من المساحة تحت المنحنى تكون ضمن حدود $(\mu - S \leq x \leq \mu + S)$
 - 2- 95.46% من المساحة تحت المنحنى تكون ضمن حدود $(\mu - 2S \leq x \leq \mu + 2S)$
 - 3- 99.74% من المساحة تحت المنحنى تكون ضمن حدود $(\mu - 3S \leq x \leq \mu + 3S)$
 - 4- 95% من المساحة تحت المنحنى تكون ضمن حدود $(\mu - 1.96S \leq x \leq \mu + 1.96S)$
 - 5- 99% من المساحة تحت المنحنى تكون ضمن حدود $(\mu - 2.58S \leq x \leq \mu + 2.58S)$
- يُقال للمتغير العشوائي المتصل x بأنه موزع طبيعياً اذا كانت دالة كثافة (كتلة) الاحتمال له

مُعطاة بالصيغة :

$$F(x) = \frac{1}{S\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{S}\right)^2} , \text{ للفترة } (-\infty \leq x \leq +\infty)$$

ويرمز لهذا التوزيع بالرمز التالي $N(\mu, S^2)$

التوزيع الطبيعي المعياري (القياسي):

بسبب ان شكل منحنى التوزيع الطبيعي يعتمد على قيمتي الوسط μ والانحراف المعياري S فلهذا هنالك مالا نهاية من المنحنيات الخاصة بذلك النوع من التوزيع وبالتالي فان قيم الاحتمالات المعتمدة على قيمة المساحة المحصورة تحت هذه المنحنيات سوف تتغير تبعاً لتلك القيم وبالتالي فلأغراض المقارنة بين قيم الاحتمالات للتوزيعات الطبيعية تم إقتراح توزيع طبيعي معياري موحد يمكنه التعبير عن جميع منحنيات التوزيعات الطبيعية: وهو التوزيع الطبيعي الذي وسطه ($\mu=0$) وانحرافه المعياري ($S=1$) وله دالة كثافة (كتلة) احتمال معرفة كما يلي:

$$F(z) = \frac{1}{1 \times \sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{z-0}{1}\right)^2} = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(z)^2}, \quad \text{للفترة } (-\infty \leq z \leq +\infty)$$

هنا لاستخراج المساحة تحت منحنى التوزيع الطبيعي المعياري لاستخراج الاحتمال ضمن الفترة المطلوبة والتي سوف تمثل حدود التكامل باستخدام صيغة التكامل التالية:

$$P(-\infty \leq z \leq +\infty) = \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(z)^2} . dz$$

أما لحساب الاحتمال لفترة ما فيتم بتبديل حدود التكامل أعلاه بقيم تلك الفترة فنستخرج قيمة المساحة تحت المنحنى ضمن تلك الفترة والتي تساوي قيمة الاحتمال لتلك الفترة.

ملاحظة: لتسهيل الحسابات والعمل بهذا النوع من التوزيع يتم استبدال التكامل بجداول خاصة بالتوزيع الطبيعي المعياري واعتماداً على قيمة المتغير الطبيعي العشوائي (Z) ، وعلى قيم الحدود للفترة المطلوب حساب الاحتمال لها.

ملاحظة: للتحويل من حالة التوزيع الطبيعي الى حالة التوزيع الطبيعي المعياري نقوم بتحويل المتغير العشوائي المعياري x الى المتغير العشوائي Z من خلال العلاقة التالية:

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

مثال (1): إذا كان x متغير عشوائي طبيعي خاضع للتوزيع $N(3,4)$ فما هي قيمة الاحتمال عند وقوع x بين القيمتين 3 و 5 ؟

Solution: $z = \frac{x-\mu}{s}$, $z_1 = \frac{3-3}{\sqrt{4}} = 0$, $z_2 = \frac{5-3}{\sqrt{4}} = 1$ → $P(0 \leq z \leq 1) = ?$ وبتطبيق هذه الحدود تصبح المسألة

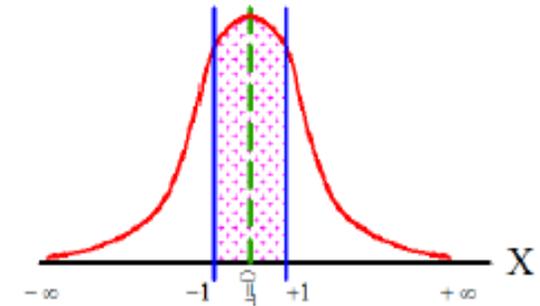
ومن الجداول الخاصة بالتوزيع الطبيعي المعياري → $P(0 \leq z \leq 1) = 0.3413$

مثال (2): إذا كان المتغير العشوائي الطبيعي x خاضع للتوزيع $N(\mu, S^2)$ ، إحسب كلٍ من $P(\mu-S \leq x \leq \mu+S)$ ، $P(\mu-2S \leq x \leq \mu+2S)$ و $P(\mu-3S \leq x \leq \mu+3S)$ ؟

Solution:

(1) For $P(\mu-S \leq x \leq \mu+S)$, $z_1 = \frac{(\mu-S)-\mu}{s} = -1$, $z_2 = \frac{(\mu+S)-\mu}{s} = 1$, → $P(-1 \leq z \leq 1) = 2 \times P(0 \leq z \leq 1) = 2 \times 0.3413 = 0.6826 = 68.26\%$ (باستخدام الجداول)

(2) و (3) H.W



ملاحظات لاستخدام الجداول:

1. في حالة وقوع قيمة المتغير العشوائي بين قيمتين في الجدول فيتم استخدام التقريب الخطي وتشابه المثلثات لإيجاد قيمة الاحتمال الصحيح.
2. يجب قبل استخدام الجدول ملاحظة الرسم الموجود في أعلى الصفحة وملاحظة المساحة المظللة تحت المنحني فهي تمثل حدود التكامل التي على غرارها أو وفقاً لها تم استخراج قيم الجدول.
3. في حالة كانت قيمة المتغير Z خارج قيم الجدول أكبر أو أقل فيتم اعتماد طريقة التكامل لاستخراج المساحة تحت المنحني وبالتالي قيمة الاحتمال المطلوبة.

مثال: للتوزيع الطبيعي المعياري $Z \sim N(0,1)$ أوجد قيمة K ، اذا كانت:

- 1) $P(Z \leq K) = 0.5$
- 2) $P(Z \geq K) = 0.2451$
- 3) $P(Z \leq K) = 0.8289$
- 4) $P(-K \leq Z \leq K) = 0.901$
- 5) $P(0 \leq Z \leq K) = 0.492$
- 6) $P(Z > K) = 0.0548$
- 7) $P(|Z| < K) = 0.901$
- 8) $P(|Z| > K) = 0.099$
- 9) $P(-1.24 < Z < K) = 0.8$

الحل:

$$1. P(Z \leq K) = 0.5$$

نلاحظ هنا بخلاف المثال السابق، القيمة المعطاة هي قيمة المساحة تحت المنحنى بينما المجهول هو إيجاد القيمة على المحور الأفقي، وبما أن جدول Z متماثل حول الصفر والمساحة الكلية تساوي الواحد فهذا يعني أن المساحة يمين الصفر تساوي 0.5 ويسار الصفر تساوي 0.5، أي أن قيمة K تساوي صفر.

$$2. P(Z \geq K) = 0.2451$$

نلاحظ أن قيمة المساحة أصغر من 0.5، وعلامة الاحتمال أكبر من ، أي أن قيمة K ستكون موجبة :

$$P(Z \geq K) = 0.2451 \rightarrow 1 - P(Z < K) = 0.2451 \rightarrow P(Z < K) = 0.7549$$

ومن خلال الجدول نجد أنه عند المساحة 0.7549 قيمة K تساوي 0.69 .

$$3. P(Z \leq K) = 0.8289$$

نلاحظ أن قيمة المساحة أكبر من 0.5، وعلامة الاحتمال أصغر من ، أي أن قيمة K ستكون موجبة ومن خلال الجدول مباشرة نجد أن قيمة K عند المساحة 0.8289 تساوي 0.95 .

$$P(-K \leq Z \leq K) = 0.901 \quad .4$$

$$P(-K \leq Z \leq K) = P(-K < Z < 0) + P(0 < Z < K)$$

من خلال خاصية التماثل في جدول التوزيع الطبيعي المعياري، نجد أن:

$$P(-K < Z < 0) = P(0 < Z < K)$$

وبذلك يكون لكل منهما نصف المساحة أي أن :

$$P(0 \leq Z \leq K) = 0.4505$$

$$P(0 \leq Z \leq K) = P(Z \leq K) - P(Z \leq 0) = 0.4505$$

$$P(Z \leq K) = P(Z \leq 0) + 0.4505 = 0.5 + 0.4505 = 0.9505$$

من خلال جدول التوزيع وعند المساحة 0.9505 نجد أن K تساوي 1.65

$$P(0 \leq Z \leq K) = 0.492 \quad .5$$

$$P(0 \leq Z \leq K) = P(Z \leq K) - P(Z \leq 0) = 0.492$$
 بنفس الطريقة السابقة:

$$P(Z \leq K) = 0.492 + 0.5 = 0.992$$

من خلال الجدول ، نجد أن قيمة K تساوي 2.41 .

$$P(Z > K) = 0.0548 \quad .6$$

$$P(Z > K) = 1 - P(Z < K) = 0.0548 \rightarrow P(Z < K) = 1 - 0.0548 = 0.9452$$

من خلال الجدول نجد أن قيمة K تساوي 1.6 .

$$P(|Z| < K) = 0.901 \quad .7$$

لاحظ أن علامة الاحتمال أصغر من وبالتالي نستطيع تحويل القيمة المطلقة الى:

$$P(-K \leq Z \leq K) = 0.901$$

وكما في الفقرة 4 تكون قيمة K 1.65 .

8 . $P(|Z| > K) = 0.099$: لاحظ أن علامة الاحتمال أكبر من أي عكس الفقرة السابقة ويتم التعامل معها كالتالي:

$$P(|Z| > K) = P(Z > K) + P(Z < -K) = 1 - P(-K \leq Z \leq K) = 0.099$$

$$\rightarrow P(-K \leq Z \leq K) = 1 - 0.099 = 0.901$$

ومن خلال الفقرة السابقة تكون قيمة K تساوي 1.65

$$P(-1.24 < Z < K) = 0.8 \quad .9$$

$$P(Z < K) - P(Z < -1.24) = 0.8 \rightarrow P(Z < K) = 0.8 + 0.1075 = 0.9075$$

ومن خلال الجدول وعند المساحة 0.9075 نجد أن قيمة K تساوي تقريبا 1.325 .

Good luck

Statistics

Lecture 13

Probability Distribution

التوزيعات الاحتمالية

The benefit of studying probability distributions from the valuable implementation success is in a major field:

1. The possibility of determining or calculating a probability in certain evidence circumstances.
2. Appropriateness of the probability distribution of some of the data collected in the field. (Therefore, it is possible to extract information from this data after making a comparison with these distributions are compatible with it).

1. امكانية تحديد او احتساب احتمال ما في ظروف بيانات معينة.

2. امكانية ملائمة التوزيع الاحتمالي لبعض البيانات التي تم جمعها ميدانيا. وبالتالي امكانية استخراج المعلومات من هذه البيانات بعد عمل مقارنة مع هذه التوزيعات المتلائمة معها.

المتغيرات العشوائية Random variables

إن بعض الكميات تكون قيمتها العددية غير محددة تحديداً نهائياً، وإنما تتغير بتغير عوامل عشوائية ومن الامثلة على هذه الكميات اختلاف مقاومات الانضغاط للمكعبات المأخوذة من نفس الخلطة ومعالجتها بنفس الظروف او مقاومات الشد لقضبان حديد التسليح المصنوعة من نفس الفولاذ وبنفس الظروف. إن هذا الاختلاف بالرغم من تساوي فرص الحصول على توافق ظاهرياً هو بسبب ظاهرة التشتت وسببها عوامل عشوائية لم يكن بالإمكان حسابها او التنبؤ بها مسبقاً. لذا يمكن تعريف الكميات العشوائية: وهي الكميات التي تأخذ قيماً مختلفة في العمليات أو المحاولات المختلفة وتبقى مختلفة مهما حاولنا جعل ظروف وشروط حدوثها متجانسة تبعاً لتغيرات عشوائية خارجة عن الحساب. **إن اهم ما يجب معرفته بما يخص المتغيرات العشوائية هو جميع القيم العددية التي يمكن أن تأخذها تلك المتغيرات، بالإضافة الى معرفة الاحتمالات المقابلة لتلك القيم لحالة المتغيرات المدروسة.** إن معلومية احتمال القيم الممكنة للكميات العشوائية المدروسة يكفي لحل أي سؤال يرتبط بهذه الكمية العشوائية وتبيان فعاليتها. وتقسّم المتغيرات العشوائية الى نوعين:

- أ- المتغيرات العشوائية المتقطعة او المنفصلة : ويكون فيها المتغير العشوائي يأخذ أي قيمة من قيم مجموعة متقطعة أو قابلة للعد مثل متغير عدد الطلبة الذكور في احدى المدارس.
- ب- المتغيرات العشوائية المتصلة : وهي المتغيرات العشوائية التي من الممكن أن تأخذ أي قيمة داخل مدى (مجال) معين من القيم مثل مقاييس درجة الحرارة أو الوزن أو الطول.

Depending on the type of the random variable, the probability distributions can be divided into two types:

بالاعتماد على انواع المتغير العشوائي ، يمكن تقسيم التوزيعات الاحتمالية الى نوعين:

A. Discrete Probability Distributions ((توزيعات احتمالية متقطعة (منفصلة))

- Binomial Dist. (توزيع ذي الحدين)
- Poisson Dist.(توزيع بواسون)

B. Continuous Probability Distributions (توزيعات احتمالية متصلة)

- Normal Distribution (التوزيع الطبيعي)
- T-Distribution (توزيع تي)
- Chi-Square Distribution (توزيع كاي تربيع)

A. Discrete Probability Distributions (توزيعات احتمالية متقطعة (منفصلة))

وهي التوزيعات الاحتمالية (دوال كتلة الاحتمال) لمتغيرات عشوائية متقطعة.

تعريف دالة كتلة الاحتمال للمتغير العشوائي المتقطع: إذا كان X متغيراً عشوائياً متقطعاً بحيث إن مجموعة القيم الممكنة له هي

$$X(S) = \{X_1, X_2, \dots, X_n\}$$

فإن دالة كتلة الاحتمال ويرمز لها $f(X)$ تأخذ القيم التالية:

$$f(X) = \begin{cases} P(X = X_i); & X_i \in X(S) \\ 0 & ; X_i \notin X(S) \end{cases}$$

ويجب ان تحقق دالة كتلة الاحتمال الشروط التالية:

- i. ان قيمة دالة كثافة الاحتمال عند قيمة X ما يجب أن لا تزيد عن 1 ولا تقل عن صفر.
- ii. ان مجموع قيم دوال كثافة الاحتمال لكل قيم المتغير X المتوقعة (ضمن القيم الممكنة ل X ان تأخذها أو فضاء العينة) يجب أن يساوي "1".

• Binomial Dist. (توزيع ذي الحدين)

وهو من اهم التوزيعات المتقطعة المستعملة في مجالات التطبيقات الهندسية لمتغيرات عشوائية منفصلة. حيث يمكن الحكم على نتائج تجربة ما ب "نعم" أو "لا". ويستخدم عند تحقيق الشروط التالية:

- ان نتيجة أي محاولة لتجربة ما تكون اما النجاح أو الفشل.
- كل محاولة نتيجتها مستقلة عن نتيجة المحاولة الاخرى. (أي لا تؤثر عليها مثل عملية السحب بإرجاع).
- احتمال النجاح في كل محاولة ثابت وليكن (P) وبالتالي فإن احتمال الفشل ثابت أيضاً وهو (q) ويحسب من (q=1-P).
- تجري التجربة بعدد معين من المحاولات أو المرات أو بمعنى آخر (يكون هنالك عدد n من المحاولات).

وعند اجراء تجربة ما وكان عدد المحاولات = n ، عدد النجاحات أو عدد مرات حصول الحادثة = x ، احتمال النجاح للتوزيع = P ، واحتمال الفشل للتوزيع = q فإن الاقتران الاحتمالي لمتغير ذي الحدين ويرمز له P(x;n,p) ويحسب كما يلي:

$$P(x; n, p) = C_x^n \cdot p^x \cdot q^{(n-x)} \quad ; X=0,1,2,\dots,n$$

ملاحظة: تم أخذ قيمة الصفر لأنه يمثل احتمال عدم الحصول او احتمال ولا مرة يحصل حادث النجاح.

ملاحظة: اذا كان X متغيراً ذا حدين (خاضع لنظرية ذات الحدين) فإن:

$$M = n \cdot p$$

$$S^2 = n \cdot p \cdot q$$

i. الوسط او التوقع لقيمة X الممكنة هو

ii. التباين لقيمة X هو

مثال(1): رُميت زهرة نرد منتظمة 5 مرات ، ما احتمال عدم ظهور الرقم 6 فيها؟ ما احتمال ظهور الرقم 6 ثلاث مرات؟

$$s = \{1,2,3,4,5,6\}$$

الحل: ان فضاء العينة لتجربة رمي زهرة نرد واحدة منتظمة لمرة واحدة هو

$$^{\wedge} \text{ احتمال ظهور الرقم 6 } P(6) = \frac{1}{6} \rightarrow ^{\wedge} \text{ احتمال عدم ظهور الرقم 6 } q(6) = 1 - \frac{1}{6} = \frac{5}{6}$$

$$^{\wedge} P(6) = \frac{1}{6}, q(6) = \frac{5}{6}, n=5, x=0,1,2,3,4,5$$

$$P\left(x; 5, \frac{1}{6}\right) = C_x^5 \cdot \left(\frac{1}{6}\right)^x \cdot \left(\frac{5}{6}\right)^{(5-x)}, x=0,1,2,3,4,5$$

وهي تمثل الصيغة العامة لهذه التجربة

$$^{\wedge} \text{ احتمال عدم ظهور الرقم 6 } P\left(x = 0; 5, \frac{1}{6}\right) = C_0^5 \cdot \left(\frac{1}{6}\right)^0 \cdot \left(\frac{5}{6}\right)^{(5-0)} = 1 \times 1 \times \left(\frac{5}{6}\right)^5 = 0.40187$$

$$^{\wedge} \text{ احتمال عدم ظهور الرقم 6 ثلاث مرات } P\left(x = 3; 5, \frac{1}{6}\right) = C_3^5 \cdot \left(\frac{1}{6}\right)^3 \cdot \left(\frac{5}{6}\right)^{(5-3)} = 10 \times \frac{1}{216} \times \left(\frac{25}{36}\right) = 0.03215$$

$$M = n \cdot p = 5 \cdot \frac{1}{6} = \frac{5}{6}$$

ولحساب القيمة المتوقعة للمتغير X نحسب الوسط

$$S^2 = n \cdot p \cdot q = 5 \cdot \frac{1}{6} \cdot \frac{5}{6} = \frac{25}{36}$$

ولحساب قيمة التباين المتوقعة للمتغير X نحسب التباين

مثال(2): يُراد توزيع عدد من الكراسي الخاصة بالطلبة الأعسرين (مستخدمي اليد اليسرى) في كل قاعة دراسية ، فإذا عُلم ان 10% من جميع الطلبة هم أعسرين فما هو احتمال إن صفاً من 20 طالب: (أ)- لا يحتوي على أي طالب أعسر؟ (ب)- اثنان منهم أعسران؟ (ج)- أكثر من اثنين منهم أعسرين؟

الحل: نفرض ان المتغير X يمثل عدد الطلبة الأعسرين

$$x=0,1,2,\dots,20, \quad n=20$$

$$q(L)=1 - 0.1=0.9 \text{ احتمال أن لا يكون الطالب أعسر } , \quad P(L)=0.1 \text{ احتمال أن يكون الطالب أعسر } .$$

من المعلومات أعلاه نلاحظ تحقق شروط نظرية ذي الحدين، لذا فإن الصيغة العامة للتجربة هي

$$P(x; 20,0.1) = C_x^{20} \cdot (0.1)^x \cdot (0.9)^{(20-x)}, \quad x=0,1,2,\dots,20$$

$$P(x = 0; 20,0.1) = C_0^{20} \cdot (0.1)^0 \cdot (0.9)^{(20-0)} = 1 \times 1 \times 0.12157 = 0.12157 \text{ (أ)}$$

$$P(x = 2; 20,0.1) = C_2^{20} \cdot (0.1)^2 \cdot (0.9)^{(20-2)} = 190 \times 0.01 \times 0.15009$$

$$=0.28518$$

• Poisson Dist.(توزيع بواسون)

يعتبر توزيع بواسون من التوزيعات الاحتمالية المتقطعة التي تستخدم للحوادث نادرة الحصول أو الوقوع، أو المتغيرات التي تحدث في أزمنة عشوائية معلومة. مثل عدد الأخطاء المطبعية في صفحة كتاب ما، عدد الزلازل السنوية الحاصلة في موقع ما، حدوث فيضان بتصريف معين في نهر ما... الخ.

وتسمى التجارب التي تعطي عدد نجاحات معين في فترة زمنية أو فترة قياس كالطول أو الحجم أو المساحة... الخ بالتجربة البواسونية، وهي تجربة احصائية تمتلك الخصائص التالية:

- نتائج التجربة مُعرفة بالنجاح أو الفشل والتجارب مستقلة الواحدة عن الأخرى.
 - الاحتمال ثابت وذو قيمة صغيرة بشكل عام (تقترب من 0) وإن عدد المحاولات n يكون كبيراً عادةً (يقترب من ∞) ولكن المهم ان مقدار حاصل الضرب بينهما يكون مقدار ثابت يساوي المتوسط أو المعدل (λ) بحيث أن ($\lambda > 0$).
- $$\lambda = n \cdot P$$

ومن خصائص توزيع بواسون الاحتمالي:

- المعدل أو المتوسط للتوزيع ($\lambda = n \cdot P$) ، والتباين يساوي (λ)
- الانحراف المعياري يساوي ($\sqrt{\lambda}$)

$$P(x, \lambda) = \frac{e^{-\lambda} \cdot \lambda^x}{x!}, \quad x=0,1,2,\dots$$

ان الاقتران الاحتمالي لتوزيع بواسون هو:

حيث ان:

λ : معدل عدد النجاحات (حصول الحدث) ضمن الفترة الزمنية أو فترة القياس (الطول، المساحة، الحجم) وغيرها.

x : عدد النجاحات الفعلي أو عدد مرات القياس المطلوبة وهي تمثل متغير بواسون العشوائي.

e : أساس اللوغاريتم الطبيعي وتساوي تقريباً ($e=2.718281828$).

مثال(1): معدل عدد حوادث السيارات على طريق صحراوي ما هو 5 (حادث\أسبوع). (أ) ما احتمال عدم حدوث أي حادث على ذلك الطريق في اسبوع معين؟
 (ب) ما احتمال حدوث 4 حوادث أو أقل في اسبوع معين؟ (ج) ما احتمال أن يزيد عدد الحوادث عن حادثين خلال أسبوعين؟

$$\lambda=5 , x=0,1,2,\dots , P(x, \lambda) = \frac{e^{-\lambda} \cdot \lambda^x}{x!}$$

الحل:

$$P(x, \lambda = 5) = \frac{e^{-5} \cdot 5^x}{x!}$$

الصيغة الأساسية لحل المسألة

$$(أ) P(x = 0, \lambda = 5) = \frac{e^{-5} \cdot 5^0}{0!} = 0.00674$$

$$(ب) P(x \leq 4, \lambda = 5) = P(0,5) + P(1,5) + P(2,5) + P(3,5) + P(4,5)$$

$$= 0.00674 + \frac{e^{-5} \cdot 5^1}{1!} + \frac{e^{-5} \cdot 5^2}{2!} + \frac{e^{-5} \cdot 5^3}{3!} + \frac{e^{-5} \cdot 5^4}{4!} = 0.4405$$

زيادة عدد الحوادث عن 2 خلال إسبوعين (هنا نلاحظ ان معدل حصول الحوادث اختلف بدلاً من الفترة اسبوع الى اسبوعين فهنا ستتغير الصيغة (ت)
 الأساسية للسؤال حسب المعدل الجديد).

$$\lambda = 5(\text{اسبوع\حادث}) \times 2 = 10(\text{اسبوعين\حادث})$$

$$^* P(x, \lambda = 10) = \frac{e^{-10} \cdot 10^x}{x!} , x = 0,1,2, \dots$$

$$P(x > 2, \lambda = 10) = 1 - \{P(0,10) + P(1,10) + P(2,10)\} = 1 - \left\{ \frac{e^{-10} \cdot 10^0}{0!} + \frac{e^{-10} \cdot 10^1}{1!} + \frac{e^{-10} \cdot 10^2}{2!} \right\} = 0.99723$$

Approximation of the binomial distribution of the Poisson distribution

تقريب توزيع ثنائي الحدين لتوزيع بواسون

إذا اعتبرنا n او عدد المحاولات كبيرة بشكل كافٍ بحيث $(n \geq 50)$ و احتمال النجاح P صغير $(P < 0.1)$ بحيث يبقى حاصل ضرب قيمتهما $(n.P)$ مقدار ثابت وليكن λ ، عندها يمكن تقريب توزيع ذي الحدين الى توزيع بواسون بالمقدار $(\lambda = n.P)$.

مثال (2): في اختبار مُضاف مُحسن للخلطات الخرسانية، اذا افترضنا انه ضار للخلطة باحتمال 0.001 وانه تم استخدامه فعليا في انتاج 2000 خلطة خرسانية، استخدم تقريب بواسون للإجابة على: (أ) ما احتمال ان تتضرر أكثر من خلطتين نتيجة استخدام ذلك المُضاف؟ (ب) ما احتمال تضرر خلطتان على الأقل؟

Solution: $n=2000$, $P=0.001$ (احتمال حصول الضرر)

نلاحظ إن عدد المحاولات كبير جداً $(n > 50)$ واحتمال حصول الحدث P أقل من 0.1 لهذا يمكن استخدام تقريب بواسون بذي الحدين أو بمعنى آخر

$$\lambda = n.P = 0.001 \times 2000 = 2$$

$$(أ) P(x > 2, \lambda = 2) = 1 - \{P(0,2) + P(1,2) + P(2,2)\} = 1 - \left\{ \frac{e^{-2} \cdot 2^0}{0!} + \frac{e^{-2} \cdot 2^1}{1!} + \frac{e^{-2} \cdot 2^2}{2!} \right\} = 0.32332$$

$$(ب) P(x \geq 2, \lambda = 2) = 1 - P(x < 2, \lambda = 2) = 1 - \{P(0,2) + P(1,2)\} = 0.59399$$

Not:

الفرق الأساسي بين توزيعي ذي الحدين وتوزيع بواسون هو إن n أو عدد التجارب الكلي في توزيع ذي الحدين يكون معلوماً ولهذا فإن قيم المتغير العشوائي x تكون $(x=0,1,2,\dots,n)$ أما في توزيع بواسون فإن n أو عدد التجارب غير معلوم أو مبهم وفي هذه الحالة فإن المتغير العشوائي x يأخذ المدى $(x=0,1,2,\dots)$. ولهذا في حالة إحتساب احتمالات أكبر من قيمة ما و لتكن a في توزيع بواسون فإننا نضطر طرح احتمال تلك القيمة والتي أقل منها من (1) وهو الأحتمال الكلي لفضاء العينة أو بمعنى آخر $P(x > a, \lambda) = 1 - P(x \leq a, \lambda)$ وذلك لأننا لا نعرف إلى أي حد سينتهي إليه المتغير العشوائي x .

مثال(2): إذا كان 3% من انتاج أحد معامل الطابوق غير مطابق للمواصفات القياسية وقد تم سحب عينة من الانتاج بشكل عشوائي. ما هو احتمال بأن العينة تحتوي بالضبط على طابوقتين فاشلتين؟ استخدم توزيعين مختلفين للحل ثم قارن النتيجةين؟

الحل: (أ) باستخدام توزيع ذي الحدين ($P=0.03$, $q=0.97$, $n=100$)

$$\therefore P(x; 100,0.03) = C_x^{100} \cdot (0.03)^x \cdot (0.97)^{100-x} \quad x = 0,1, \dots, 100$$

$$\therefore P(x = 2; 100,0.03) = \frac{100!}{2! \times (100 - 2)!} \times (0.03)^2 \times (0.97)^{100-2} = 4950 \times 0.0009 \times 0.0505 = 0.225153$$

(ت) باستخدام توزيع بواسون ($\lambda = n \cdot P = 0.03 \times 100 = 3$)

$$P(x, \lambda = 3) = \frac{e^{-3} \cdot 3^x}{x!}, \quad x = 0,1,2, \dots$$

$$\therefore P(x = 2, \lambda = 3) = \frac{e^{-3} \cdot 3^2}{2!} = 0.22404$$

مثال(3): اذا كانت احدى البدالات تتلقى مكالمات هاتفية بمعدل مكالمتين لكل دقيقة. أوجد احتمال أن تتلقى البدالة 3 مكالمات على الأقل؟

Solution: ($\lambda=2$ call/min)

$$P(x, 2) = \frac{e^{-2} \cdot 2^x}{x!}, \quad x = 0, 1, 2, \dots$$

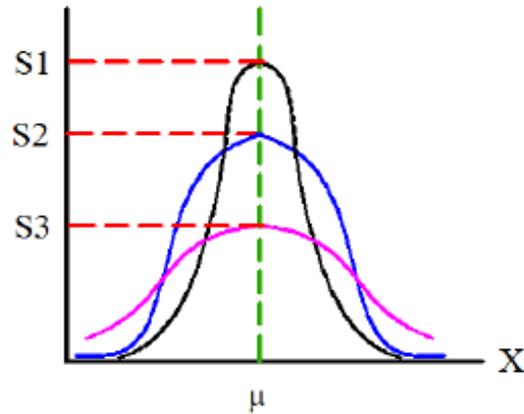
$$P(x \geq 3, \lambda = 2) = 1 - \left\{ \frac{e^{-2} \cdot 2^0}{0!} + \frac{e^{-2} \cdot 2^1}{1!} + \frac{e^{-2} \cdot 2^2}{2!} \right\} = 1 - (0.67667) = 0.3233$$

التوزيعات الاحتمالية المتصلة (Continuous Probability Distributions)

-1-1 التوزيع الطبيعي (Normal Distribution):

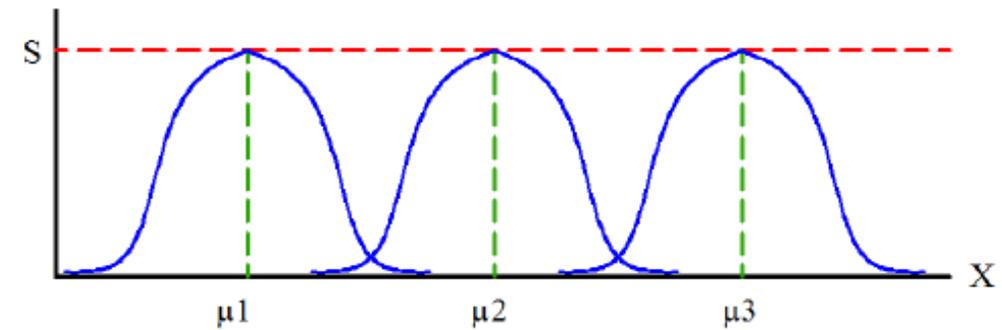
وهو من أشهر التوزيعات الاحتمالية المتصلة وأكثرها استخداماً وتطبيقاً ومن الأمثلة عليه التوزيعات البايومترية (توزيعات الطول والوزن) وتوزيعات اخطاء المشاهدات (الفروق بين القيم الحقيقية والقيم المُشاهدة) وتم تطوير هذا التوزيع من قبل العالم غاوس سنة 1809 ولهذا يطلق عليه في بعض الأحيان اسم (توزيع غاوس).

ان منحنى التوزيع الطبيعي الذي يأخذ شكل الجرس هو منحنى توزيع نظري لدالة كثافة (كتلة الاحتمال) للمتغير العشوائي الذي يخضع لهذا التوزيع. هنالك ما لانهاية من منحنيات التوزيع الطبيعي التي قد تختلف عن بعضها البعض حسب قيمة كل من الوسط الحسابي (μ) والانحراف المعياري (S)، فمثلاً قد تتفق المنحنيات بالانحراف المعياري لكنها تختلف بالوسط الحسابي **كما في الشكل (1)** أو بالعكس قد تختلف بالانحراف المعياري وتتفق بالوسط **كما في الشكل (2)**.



شكل (2)

منحنيات توزيع طبيعي لها نفس الوسط الحسابي مع اختلاف الانحراف المعياري



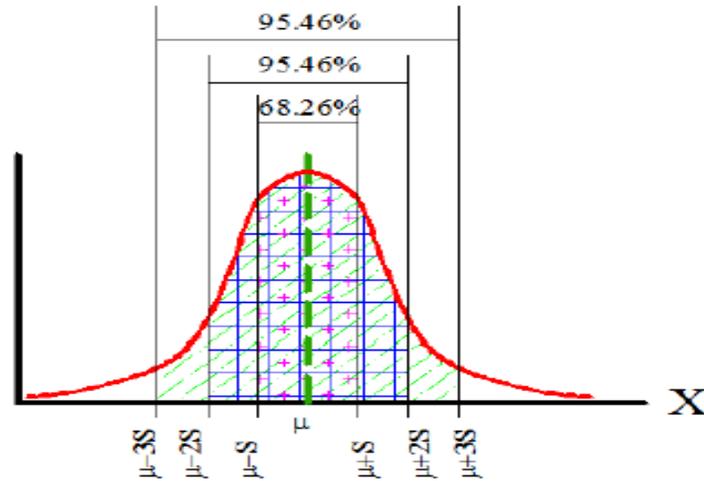
شكل (1)

منحنيات توزيع طبيعي لها نفس الانحراف المعياري مع اختلاف الوسط الحسابي

خصائص التوزيع الطبيعي:

- منحنى التوزيع الطبيعي متماثل حول الوسط μ (إتوائه يساوي صفراً) ، وشكله بشكل الجرس.
- له قيمة منوالية (منوال) واحد ينطبق على الوسط μ ويقترب طرفاه من قيمتي $(-\infty, \infty)$ دون أن تقطع المحور الافقي.
- قيمة الانحراف المعياري تحدد عرض المنحنى الطبيعي، مثلاً قيمة انحراف معياري أعلى تنتج منحنى أعرض وأكثر تسطحاً.
- الاحتمالات للمتغير الطبيعي العشوائي تُعطى بالمساحة تحت منحنى التوزيع الطبيعي. والمساحة الكلية تحت منحنى التوزيع الطبيعي للفترة $(-\infty \leq x \leq +\infty)$ تساوي (1) و تتوزع مناصفةً 0.5 للنصف الأيمن و 0.5 للنصف الأيسر.
- قيمة الاحتمال عند قيمة مفردة محددة للمتغير العشوائي الطبيعي x تساوي (صفر) ولأي قيمة وذلك لان المساحة تحت المنحنى والمحددة بقيمة واحدة تكون عبارة عن خط مستقيم ولهذا فلهذا فلحساب الاحتمال يجب ان تحدد فترة ما للمتغير العشوائي ومن خلالها تحسب المساحة تحت المنحنى ضمن هذه الفترة والتي سوف تساوي قيمة الاحتمال ضمن تلك الفترة.

$$P(x=0) = 0 , P(x=1) = 0 , P(x=100) = 0 , P(0 < x < a) = ? , P(x > a) = ?$$



المساحة تحت المنحنى الطبيعي تتوزع كما يلي (ثابتة لجميع منحنيات التوزيع الطبيعي)

- 1- 68.26% من المساحة تحت المنحنى تكون ضمن حدود $(\mu - S \leq x \leq \mu + S)$
 - 2- 95.46% من المساحة تحت المنحنى تكون ضمن حدود $(\mu - 2S \leq x \leq \mu + 2S)$
 - 3- 99.74% من المساحة تحت المنحنى تكون ضمن حدود $(\mu - 3S \leq x \leq \mu + 3S)$
 - 4- 95% من المساحة تحت المنحنى تكون ضمن حدود $(\mu - 1.96S \leq x \leq \mu + 1.96S)$
 - 5- 99% من المساحة تحت المنحنى تكون ضمن حدود $(\mu - 2.58S \leq x \leq \mu + 2.58S)$
- يُقال للمتغير العشوائي المتصل x بأنه موزع طبيعياً اذا كانت دالة كثافة (كتلة) الاحتمال له

مُعطاة بالصيغة :

$$F(x) = \frac{1}{S\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{S}\right)^2} , \text{ للفترة } (-\infty \leq x \leq +\infty)$$

ويرمز لهذا التوزيع بالرمز التالي $N(\mu, S^2)$

التوزيع الطبيعي المعياري (القياسي):

بسبب ان شكل منحنى التوزيع الطبيعي يعتمد على قيمتي الوسط μ والانحراف المعياري S فلهذا هنالك مالا نهاية من المنحنيات الخاصة بذلك النوع من التوزيع وبالتالي فان قيم الاحتمالات المعتمدة على قيمة المساحة المحصورة تحت هذه المنحنيات سوف تتغير تبعاً لتلك القيم وبالتالي فلأغراض المقارنة بين قيم الاحتمالات للتوزيعات الطبيعية تم إقتراح توزيع طبيعي معياري موحد يمكنه التعبير عن جميع منحنيات التوزيعات الطبيعية: وهو التوزيع الطبيعي الذي وسطه ($\mu=0$) وانحرافه المعياري ($S=1$) وله دالة كثافة (كتلة) احتمال معرفة كما يلي:

$$F(z) = \frac{1}{1 \times \sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{z-0}{1}\right)^2} = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(z)^2} \quad , \quad \text{للفترة } (-\infty \leq z \leq +\infty)$$

هنا لاستخراج المساحة تحت منحنى التوزيع الطبيعي المعياري لاستخراج الاحتمال ضمن الفترة المطلوبة والتي سوف تمثل حدود التكامل باستخدام صيغة التكامل التالية:

$$P(-\infty \leq z \leq +\infty) = \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(z)^2} . dz$$

أما لحساب الاحتمال لفترة ما فيتم بتبديل حدود التكامل أعلاه بقيم تلك الفترة فنستخرج قيمة المساحة تحت المنحنى ضمن تلك الفترة والتي تساوي قيمة الاحتمال لتلك الفترة.

ملاحظة: لتسهيل الحسابات والعمل بهذا النوع من التوزيع يتم استبدال التكامل بجداول خاصة بالتوزيع الطبيعي المعياري واعتماداً على قيمة المتغير الطبيعي العشوائي (Z) ، وعلى قيم الحدود للفترة المطلوب حساب الاحتمال لها.

ملاحظة: للتحويل من حالة التوزيع الطبيعي الى حالة التوزيع الطبيعي المعياري نقوم بتحويل المتغير العشوائي المعياري X الى المتغير العشوائي Z من خلال العلاقة التالية:

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

مثال (1): إذا كان x متغير عشوائي طبيعي خاضع للتوزيع $N(3,4)$ فما هي قيمة الاحتمال عند وقوع x بين القيمتين 3 و 5 ؟

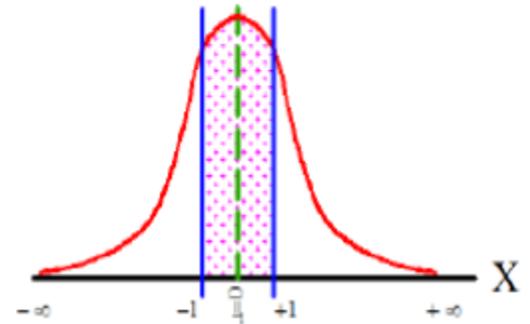
Solution: $z = \frac{x-\mu}{s}$, $z_1 = \frac{3-3}{\sqrt{4}} = 0$, $z_2 = \frac{5-3}{\sqrt{4}} = 1$ → وبتطبيق هذه الحدود تصبح المسألة $P(0 \leq z \leq 1) = ?$

ومن الجداول الخاصة بالتوزيع الطبيعي المعياري → $P(0 \leq z \leq 1) = 0.3413$

مثال (2): إذا كان المتغير العشوائي الطبيعي x خاضع للتوزيع $N(\mu, S^2)$ ، احسب كلٍ من $P(\mu-S \leq x \leq \mu+S)$ ، $P(\mu-2S \leq x \leq \mu+2S)$ و $P(\mu-3S \leq x \leq \mu+3S)$ ؟

Solution:

(1) For $P(\mu-S \leq x \leq \mu+S)$, $z_1 = \frac{(\mu-S)-\mu}{s} = -1$, $z_2 = \frac{(\mu+S)-\mu}{s} = 1$, → $P(-1 \leq z \leq 1) = 2 \times P(0 \leq z \leq 1) = 2 \times 0.3413 = 0.6826 = 68.26\%$ (باستخدام الجداول)



Statistics

Lecture 6

Measures of central tendency

مقاييس النزعة المركزية

The basic idea behind measures of central tendency is to represent a large set of data with a single value. This value is usually in the middle of the data, that is, in its center, and it is the value that other values tend to and clustered around. Thus, this property that most data possesses is called the central tendency property, and I term measures that measure the average value around which the cluster gets, measures of central tendency or averages.

The average value is taken as a representative of the whole group, on the basis that it is a non-extreme value, but rather a value that is gathered around it. Most of the values are, and therefore, first than others in representing data. Knowing the middle value of the data is useful for us. Study the characteristics of the frequency distribution and compare the different repetitive distributions for the same phenomenon.

الفكرة الأساسية التي تعتمد عليها مقاييس النزعة المركزية هي تمثيل مجموعة كبيرة من البيانات بقيمة واحدة ، وهذه القيمة عادة تكون في وسط البيانات أي في مركزها وهي القيمة التي تميل إليها بقية القيم وتتجمع حولها ، وبالتالي سميت هذه الخاصية التي تتمتع بها معظم البيانات بخاصية النزعة المركزية ، وأطلق على المقاييس التي تقيس هذه القيمة المتوسطة التي يحصل حولها التجمع ، مقاييس النزعة المركزية أو المتوسطات. وتؤخذ قيمة المتوسط كممثل للمجموعة كلها ، على أساس أنها قيمة غير متطرفة بل هي قيمة تتجمع حولها أغلبية القيم ، وبالتالي هي أولى من غيرها في تمثيل البيانات. ومعرفة القيمة الوسطى للبيانات تفيدنا في دراسة خصائص التوزيع التكراري والمقارنة بين التوزيعات التكرارية المختلفة لنفس الظاهرة.

One of the main averages that we are going to study the following:

1. Mean:

The mean is defined as the value that represents the center of gravity of the data, i.e. The anchor point then the equilibrium occurs, if we have a set of data and match it with weights of equal weight on a stepped board, it will balance if suspended or fixed from its center of gravity and center of gravity this is what represents the mean value of these values.

$$\text{Mean} = \frac{\text{Sum of values}}{\text{Number of values}}$$

We denote the mean of the sample by the symbol (\bar{X}) , and we denote the mean of the community by the symbol μ .

عرف الوسط الحسابي بأنه القيمة التي تمثل مركز ثقل البيانات أي نقطة الارتكاز التي يحصل عندها التوازن ، فإذا كان لدينا مجموعة من البيانات ومثلناها بأثقال متساوية الوزن على لوح مدرج ، فسنجد أن هذا اللوح سيتزن إذا علق أو ثبت من مركز ثقله ومركز الثقل هذا هو الذي يمثل قيمة الوسط الحسابي لهذه القيم.

In the following, we will show how to apply the mean formula in the case of ungrouped data, i.e. Not shown in the frequency distributions tables, and in the case of grouped data, i.e. shown in Frequency distributions tables.

1. Calculate the arithmetic mean of unclassified data :

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

(\bar{X}) : الوسط الحسابي

x : رمز المتغير

n : عدد القيم

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$$

Ex1: The following data represent 6 scores in the statistics subject, (4, 6, 3, 7, 9, 1) calculate the arithmetic mean

of these scores.

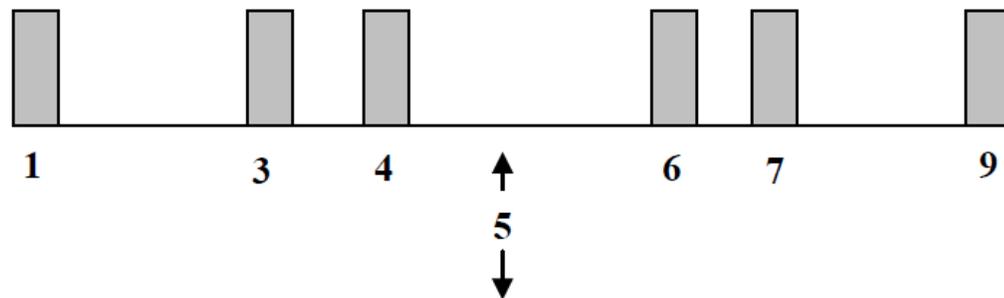
Sol:

$x_1= 4, x_2= 6, x_3= 3, \dots\dots\dots x_6= 1.$

$$\bar{X} = \frac{\sum_{i=1}^n x_i}{n}$$

$$= \frac{4+6+3+7+9+1}{6} = 5$$

This means that the center of gravity, i.e., the focal point at which the balance occurs, is degree 5 this is as shown in the figure.



Ex2: The following data shows the value of imports (in millions of Iraqi dinars) for a port in years from 1976 to 1980, according to bulletins issued by the Bureau of Statistics and Census.

the years	1976	1977	1978	1979	1980
Import value	381	292	293	333	391

Calculate the arithmetic mean of the value of this port's imports in these five years.

Sol:

$$\begin{aligned}\bar{X} &= \frac{\sum_{i=1}^n x_i}{n} \\ &= \frac{381+272+293+333+391}{5} = \frac{1670}{5} = 334\end{aligned}$$

2. Calculation of the arithmetic mean of the classified data:

A. If the data is classified in a frequency distribution table to represent each category of the table only one value, in this case in order to calculate the arithmetic mean, we perform the following steps:

1. نحسب المجموع الكلي للقيم التابعة لكل فئة بضرب قيمة x_i التي تمثلها الفئة في تكرار الفئة f_i ، اي نحسب $(x_i f_i)$ لكل الفئات.

2. نحسب المجموع الكلي لجميع القيم وذلك بجمع المجاميع الخاصة بكل الفئات والتي تحصلنا عليها في الخطوة السابقة، أي نحسب:

حيث $k =$ عدد الفترات أو عدد القيم المختلفة.

$$\left(\sum_{i=1}^k x_i f_i \right)$$

3. نحسب العدد الكلي للقيم المشاهدة وهو يساوي المجموع الكلي للتكرارات $\left(\sum_{i=1}^k f_i \right)$.

4. نحسب الوسط الحسابي بقسمة المجموع الكلي للقيم $\left(\sum_{i=1}^k x_i f_i \right)$ على العدد الكلي للقيم $\left(\sum_{i=1}^k f_i \right)$.

B. If the data is grouped in a repeating table so that each of the table categories is more representative
Of a single value.

إذا كانت البيانات مبوبة في جدول تكراري بحيث كل فئة من فئات الجدول تمثل أكثر من قيمة واحدة ، ففي هذه الحالة لا نستطيع معرفة القيم المشاهدة التابعة لكل فئة ، والذي نعرفه هو عددها فقط والمتمثل في تكرار الفئة ، ولحساب المجموع الكلي للقيم التابعة لكل فئة سنفترض أن القيم المشاهدة موزعة حول مركز الفئة داخل كل فئة من فئات الجدول توزيعاً عادلاً ، وبالتالي تكون قيمة مركز الفئة هي القيمة الافتراضية لجميع القيم داخل الفئة.

ولكي نحسب الوسط الحسابي في هذه الحالة نجري الخطوات التالية :

1. نحسب مركز كل فئة ونرمز للمركز بالرمز x_i .
2. نحسب المجموع الكلي للقيم التابعة لكل فئة بضرب القيمة x_i (حيث x_i مركز الفئة) في تكرار الفئة f_i ، أي نحسب $(x_i f_i)$ لكل الفئات.
3. نحسب المجموع الكلي لجميع القيم وذلك بجمع المجاميع الخاصة بكل الفئات والتي تحصلنا عليها في الخطوة السابقة، أي نحسب $(\sum_{i=1}^k x_i f_i)$
4. نحسب العدد الكلي للقيم المشاهدة وهو يساوي المجموع الكلي للتكرارات $(\sum_{i=1}^k f_i)$.
5. نحسب الوسط الحسابي بقسمة المجموع الكلي للقيم $(\sum_{i=1}^k x_i f_i)$ على العدد الكلي للقيم $(\sum_{i=1}^k f_i)$.

وهكذا تكون صيغة المتوسط الحسابي للبيانات المبوبة كما يلي :

$$\bar{X} = \frac{\sum_{i=1}^k x_i f_i}{\sum_{i=1}^k f_i}$$

حيث :

x : القيمة التي تمثلها الفئة (عندما الفئة تمثل قيمة واحدة فقط) ،

وهي عبارة عن مركز الفئة (عندما الفئة تمثل اكثر من قيمة واحدة).

f: تكرار الفئة.

Ex1: The following data represent the results of the compression test for fifty concrete cubes in units (MPa).

(29 30 32 39 38 38 44 28 33 31 35 37 42 49 25 34 30 31 35 37 40 26 32 33 33 39 44 45 26 31 34 31 36 36 41 27 30 31 32 37 40 39 25 34 31 30 38 35 43 48).

Calculate the mean for the two cases:

1. Unclassified data.
2. Data classified in a frequency table with categories.

Sol:

$$1. \quad \bar{X} = \frac{\sum_{i=1}^n x_i}{n} = \frac{29+30+32+\dots+48}{50} = \frac{1744}{50} = 34.88$$

2.

Class limits حدود الفئات	Centers of classes مراكز الفئات	Frequencies التكرارات	The arithmetic mean of the actual values within the classes
25 - 29	27	7	26.571
30 - 34	32	19	31.737
35 - 39	37	14	37.071
40 - 44	42	7	42
45 - 49	47	3	47.333

$$\bar{X} = \frac{\sum_{i=1}^k x_i f_i}{\sum_{i=1}^k f_i}$$

$$\bar{X} = \frac{27 * 7 + 32 * 19 + 37 * 14 + 42 * 7 + 47 * 3}{50} = 35$$

- نلاحظ هنا الفرق بين النتيجةين وسبب واضح وهو وجود تقريب في حالة تبويب البيانات لاننا استخدمنا بديل عن البيانات الواقعية وهو مراكز الفئات بدلا من الفئات نفسها.

Ex2: The data shown in the following frequency table represent the scores of 50 students:

the scores	Number of students
30 – 39	4
40 – 49	6
50 – 59	8
60 – 69	12
70 – 79	9
80 – 89	7
90 – 99	4

Calculate the mean of these scores.

Sol:

$$\bar{X} = \frac{\sum_{i=1}^k x_i f_i}{\sum_{i=1}^k f_i}$$

the scores	Number of students f_i	Centers of classes مراكز الفئات x_i	$x_i f_i$
30 – 39	4	34.5	138
40 – 49	6	44.5	267
50 – 59	8	54.5	436
60 – 69	12	64.5	774
70 – 79	9	74.5	670.5
80 – 89	7	84.5	591.5
90 – 99	4	94.5	378
SUM=	50		3255

$$\bar{X} = \frac{3255}{50} = 65.1$$

Statistics

Lecture 7

Properties of the arithmetic mean:

ان لقيمة الوسط الحسابي صفات تؤهله بأن يكون ممثلاً عن القيم للبيانات الأخرى المحسوب منها. من مزايا الوسط الحسابي انه يأخذ بنظر الاعتبار كل القيم ومن عيوبه انه يتأثر بالقيم الشاذة و المتطرفة و يصعب حسابه في حالة البيانات الوصفية و في حالة الجداول التكرارية المفتوحة. ومن اهم خصائصه:

1. أهم مقاييس النزعة المركزية وأكثرها استعمالاً نظراً لسهولة حسابه وإمكانية التعامل معه رياضياً ، ولذلك له أهمية قصوى في التحليل الإحصائي إذ إنه يدخل في حساب كثير من المقاييس الأخرى.
2. يأخذ بعين الاعتبار جميع القيم المشاهدة ، ونستطيع أن نحصل على مجموع القيم المشاهدة إذا عرفنا قيمة الوسط الحسابي كما يلي :
مجموع القيم = الوسط الحسابي × عدد القيم
3. الوسط الحسابي هو قيمة نظرية وليس بالضروري أن تكون واحدة من القيم التي يمكن أن يأخذها المتغير محل الدراسة.
4. يتأثر بوجود قيم متطرفة في البيانات وينحاز لها ، فمثلاً إذا تبرع خمسة عشر فرداً لعمل خيري بالمبالغ التالية :
10 , 15 , 15 , 20 , 15 , 15 , 12 , 20 , 10 , 10 , 18 , 20 , 10 , 10 , 1000 .

نلاحظ أن الوسط الحسابي للتبرعات كبير ولا يمثل ما دفعه معظم الأفراد ، فهنا انحاز الوسط الحسابي للقيمة المتطرفة 1000 ، وبالتالي يعدُّ الوسط الحسابي في هذه الحالة مقياساً مضللاً . ولكن لو استبعدنا القيمة المتطرفة فسنلاحظ أن الوسط الحسابي سيكون واقعياً .

$$\bar{X} = \frac{1200}{15} = 80$$

5. لا يمكن حسابه في حالة البيانات النوعية (الوصفية) .

6. لا يمكن إيجاده من الرسم (أي بيانياً) .

7. لا يمكن حسابه في حالة جداول التوزيعات التكرارية المفتوحة ؛ لأننا لا نستطيع حساب مراكز الفئات المفتوحة.

8. مجموع انحرافات القيم عن وسطها الحسابي يساوي صفرأ ، حيث انحراف القيمة عن الوسط الحسابي المقصود به القيمة مطروحاً منها الوسط الحسابي ، وبالتالي فهذه الخاصية تعني أن:

$$\sum_{i=1}^n (X_i - \bar{X}) = 0$$

For example:

إذا كانت لدينا البيانات التالية: 6, 8, 5, 11, 10 فإن الوسط الحسابي لهذه البيانات هو:

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n} = \frac{6+8+5+11+10}{5} = 8$$

فسنجد أن مجموع انحرافات هذه القيم عن وسطها الحسابي يساوي صفرأ ، كما هو مبين فيما يلي :

الانحرافات ($X_i - \bar{X}$)	القيمة (X_i)
$6 - 8 = -2$	6
$8 - 8 = 0$	8
$5 - 8 = -3$	5
$11 - 8 = 3$	11
$10 - 8 = 2$	10
صفر	المجموع

9. مجموع مربعات انحرافات القيم عن وسطها الحسابي أقل ما يمكن.

2. The Median (الوسيط):

The median is the value in the middle of the data after arranging it in ascending or descending order (M).

الوسيط هو القيمة الموجودة في منتصف البيانات بعد ترتيبها تصاعدياً أو تنازلياً ، أي هو القيمة الوسطى التي يكون نصف البيانات أقل منها والنصف الآخر أكبر منها.

The median is calculated as follows:

A. Non-classified data:

يحسب الوسيط في حالة البيانات غير المبوبة باتباع الخطوات التالية :

أ – نرتب القيم المشاهدة تصاعدياً ، أي نبدأ بأصغر قيمة ثم نرتب ما بعدها الأكبر فالأكبر .

ب – نحدد ترتيب الوسيط بين القيم وذلك كما يلي :

ترتيب الوسيط هو $\frac{n+1}{2}$ حيث n عدد القيم.

$$\text{Median arrangement} = \frac{n+1}{2}$$

ج – نحدد قيمة الوسيط وتكون هي القيمة التي ترتيبها $\frac{n+1}{2}$ عندما يكون عدد القيم n عدداً فردياً، أما إذا كانت قيمة n عدداً زوجياً سيقع الترتيب بين قيمتين ويكون الوسيط هو معدل هاتين القيمتين.

Ex1: Find the median value of the following data:

13, 17, 15, 18, 20, 19, 16.

Sol:

نرتب القيم ترتيباً تصاعدياً:

13, 15, 16, 17, 18, 19, 20

$n = 7$ (عدد فردي)

$$\text{Median arrangement} = \frac{n+1}{2} = \frac{7+1}{2} = 4$$

أي أن قيمة الوسيط هي القيمة الرابعة في البيانات بعد ترتيب البيانات تصاعدياً وبالتالي فإن :

$$\text{الوسيط} = 17$$

Ex2: Find the median value of the following data:

3, 8, 9, 6, 5, 11.

Sol:

نرتب القيم ترتيباً تصاعدياً:

3, 5, 6, 8, 9, 11.

n = 6 (عدد زوجي)

$$\text{Median arrangement} = \frac{n+1}{2} = \frac{6+1}{2} = 3.5$$

قيمة الوسيط تقع بين القيمتين الثالثة والرابعة ، إذاً الوسيط يساوي المعدل لهاتين القيمتين أي:

$$M = \frac{6+8}{2} = 7$$

B. classified data:

لحساب الوسيط في حالة البيانات المعروضة في جداول تكرارية بحيث تكون فئات الجدول مرتبة ترتيباً تصاعدياً أو تنازلياً، نتبع الخطوات التالية:

1. نحدد ترتيب الوسيط ويساوي في البيانات المبوبة

$$\text{Median arrangement} = \frac{\sum_{i=1}^k f_i}{2}$$

2. نحدد الفئة الوسيطة وهي الفئة التي تحتوي على الوسيط ويتم تحديدها بالاستعانة بالتكرار المتجمع الصاعد ، فتكون الفئة الوسيطة هي أول فئة

يكون تكرارها المتجمع الصاعد أكبر من أو يساوي ترتيب الوسيط .

3. نقوم بالتعويض في القانون التالي للحصول على قيمة الوسيط :

$$m = L + \left(\frac{\frac{\sum_{i=1}^k f_i}{2} - F}{f_m} \right) \times c$$

حيث:

L: الحد الأدنى للفئة الوسيطة.

F: التكرار المتجمع الصاعد للفئة السابقة للفئة الوسيطة .

f_m : تكرار الفئة الوسيطة .

C: طول الفئة الوسيطة .

والقاعدة الرياضية التي يعتمد عليها قانون الوسيط هو افتراض أن القيم التابعة لكل فئة موزعة حول مركزها توزيعاً عادلاً ، وبالتالي سيكون هناك تناسب بين بعد الوسيط عن الحد الأدنى للفئة الوسيطة والذي نرسم له بالرمز (L) وطول الفئة الوسيطة (C) وبين عدد القيم السابقة للوسيط داخل الفئة الوسيطة (ترتيب الوسيط F) والعدد الكلي للقيم المشاهدة التابعة للفئة الوسيطة (f_m)، أي أن هناك تناسباً بين المسافات والتكرارات .

Ex3: Find the median value of the following data:

Class limits	f_i
50 – 54	10
54 – 58	30
58 – 62	90
62 – 66	60
66 – 70	20

Sol:

$$\text{Median arrangement} = \frac{\sum_{i=1}^k f_i}{2} = \frac{210}{2} = 105$$

نحدد الفئة الوسيطة أي نحدد الفئة التي تحتوي على الوسيط وهو القيمة التي ترتيبها 105 ويتم ذلك بالاستعانة بالترتيب المتجمع الصاعد الموضح في الجدول التالي:

Class limits	f_i	The up word combined frequency
50 – 54	10	10
54 – 58	30	40
58 – 62	90	130
62 – 66	60	190
66 – 70	20	210

→ الفئة الوسيطة

حيث إن الفئة الثانية تكرارها المتجمع يساوي 40 ، والفئة الثالثة تكرارها المتجمع يساوي 130 ، فيعني ذلك أن التسعين قيمة الموجودة في الفئة الثالثة ترتيباتها تبدأ من الترتيب 41 إلى الترتيب 130 ، وحيث إن ترتيب الوسيط يساوي 105 ، إذن يجب أن يكون ضمن قيم الفئة الثالثة ، وبالتالي فإن الفئة الثالثة هي الفئة الوسيطة.

وبعد تحديد الفئة الوسيطة نطبق قانون الوسيط ، حيث :

$$f_m = 90 \quad C=4 \quad L =58 \quad F=40$$

$$m = L + \left(\frac{\frac{\sum_{i=1}^k f_j}{2} - F}{f_m} \right) \times c$$
$$= 58 + \frac{(105 - 40)}{90} \times 4$$
$$= 58 + 2.89 = 60.89$$

ويعني ذلك ان 50% الفئات تكون اقل من 60.89 والنسبة المتبقية تكون اكثر من 60.89 .

Statistics

Lecture 8

Measures of dispersion

In statistics, the measures of dispersion help to interpret the variability of data i.e. to know how much homogenous or heterogeneous the data is. In simple terms, it shows how squeezed or scattered the variable is.

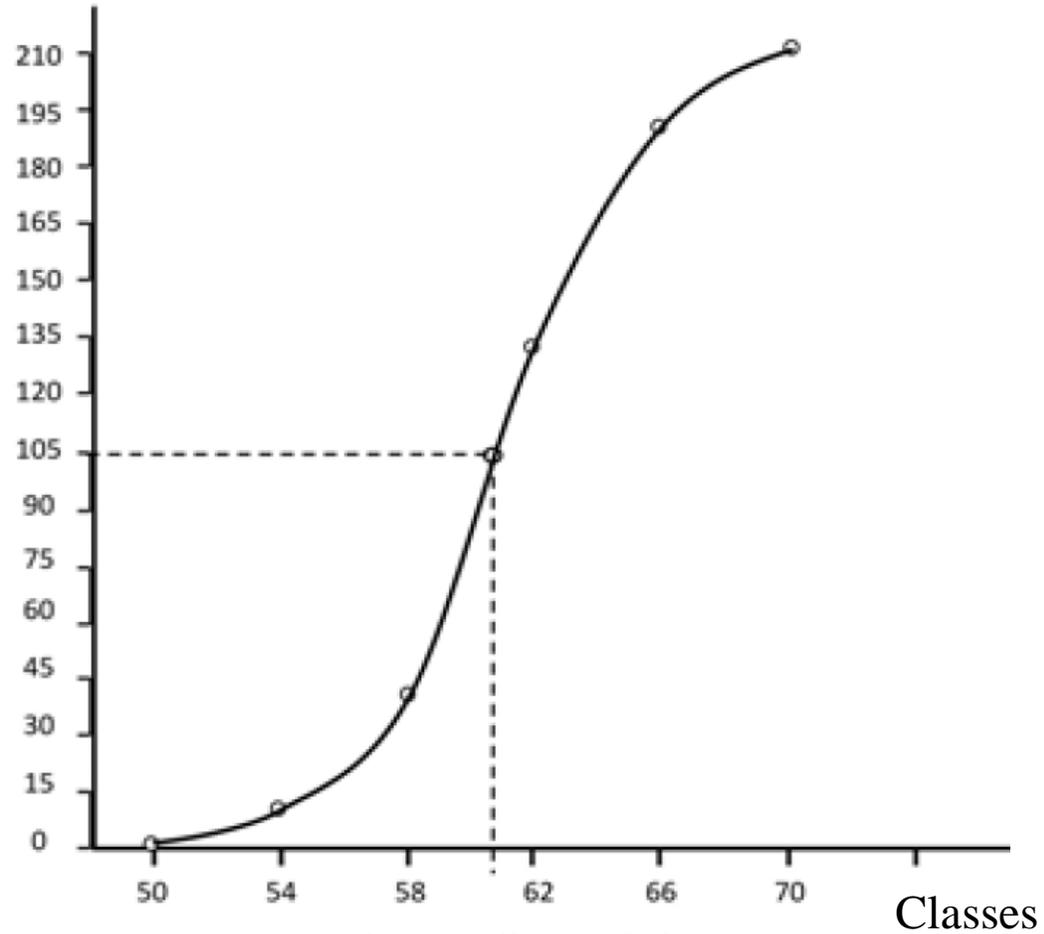
في الإحصاء ، تساعد مقاييس التشتت في تفسير تباين البيانات ، أي معرفة مقدار البيانات المتجانسة أو غير المتجانسة. بعبارة بسيطة ، يوضح مدى ضغط أو تشتت المتغير.

Sol:

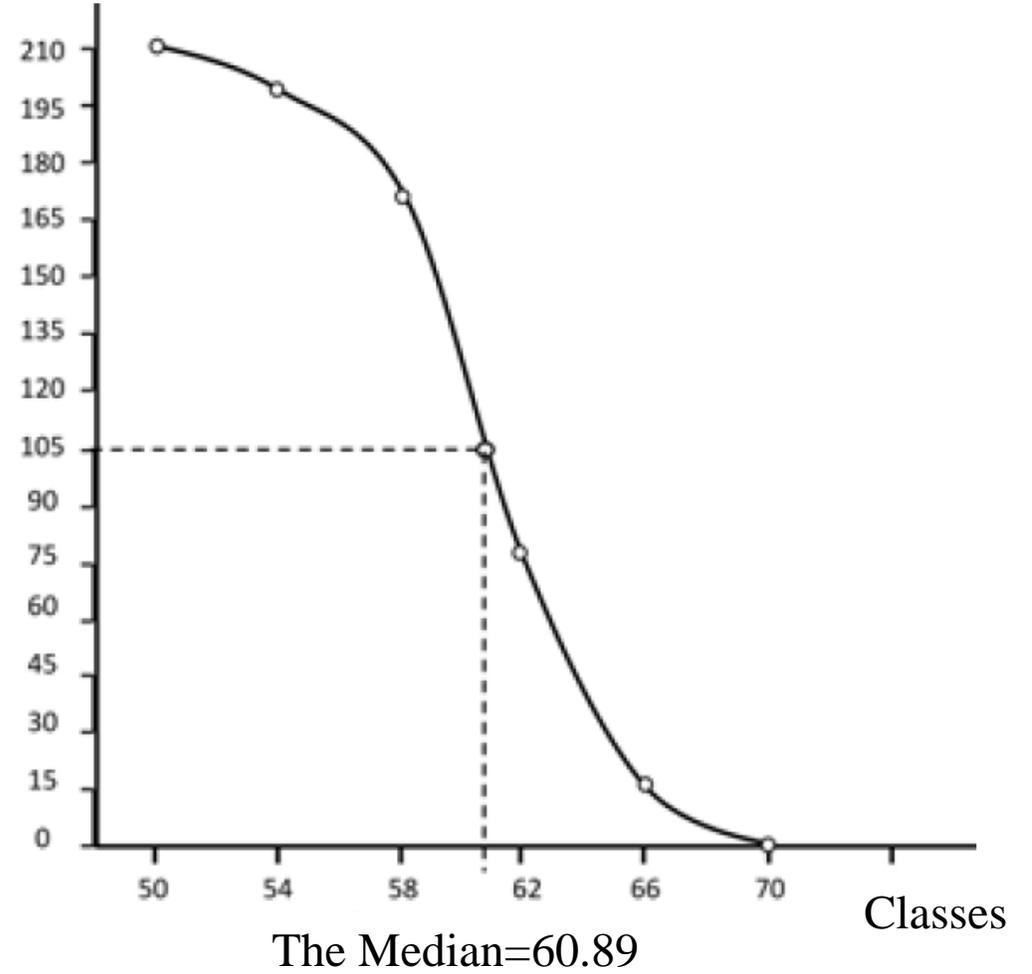
نكوّن أولاً جدول التكرار المتجمع الصاعد أو جدول التكرار المتجمع الهابط لهذه البيانات وبعد رسم المنحنى المتجمع وتحديد ترتيب الوسيط الذي يساوي 105 على المحور الرأسي ، نرسم من هذه النقطة خطاً أفقياً يوازي محور السينات حتى يلاقى المنحنى المتجمع ومن نقطة التلاقي نسقط عموداً ليلقي محور السينات فنحصل على قيمة الوسيط.

Classes	The up word combined frequency	Classes	The down word combined frequency
Less than 50	0	More than 50	210
Less than 54	10	More than 50	200
Less than 58	40	More than 50	170
Less than 62	130	More than 50	80
Less than 66	190	More than 50	20
Less than 70	210	More than 50	0

The up word combined frequency



The down word combined frequency



Median properties:

- سهل التعريف وسهل الحساب .
- يعتمد على القيمة الوسطى فقط إذا كانت (n) فردية وعلى القيمتين الوسيطيتين إذا كانت زوجية، ويهمل بقية القيم .
- لا يتأثر بوجود قيم متطرفة .

Ex2: Find the median value of the following data:

(2 , 7 , 6 , 4 , 10 , 1 , 9).

Sol:

1 , 2 , 4 , 6 , 7 , 9 , 10

نرتب الاعداد ترتيبا تصاعديا

$$n = 7$$

$$\text{Median arrangement} = \frac{n+1}{2} = \frac{7+1}{2} = 4$$

The Median = 6

لو استبدلنا بدل 10 رقم اخر يعتبر متطرف مثلا 80 فإن ترتيب الوسيط يبقى نفسه وهو الرقم الرابع وقيمته نفسها 6 .

- يمكن إيجاد الوسيط للبيانات النوعية (الوصفية) بشرط إمكانية ترتيبها تصاعدياً أو تنازلياً ويكون عددها فردياً .
فمثلاً إذا كان لدينا البيانات النوعية التالية والتي تمثل تقديرات 9 طلبة في مادة الإحصاء :
مقبول، جيد، ضعيف، ممتاز، جيد جداً، مقبول، مقبول، جيد، مقبول الترتيب التصاعدي للبيانات:
ضعيف،مقبول،مقبول،مقبول،مقبول،جيد،جيد، جيد جداً، ممتاز

$$\text{Median arrangement} = \frac{9+1}{2} = 5$$

إذن الوسيط هو التقدير الخامس بعد ترتيب البيانات تصاعدياً، أي أن الوسيط هو تقدير مقبول.

- يمكن حسابه من جداول التوزيعات التكرارية المفتوحة.
- يمكن إيجاد الوسيط بيانياً.

The Mode (M_0):

المنوال: هو القيمة أو الصفة الأكثر شيوعاً في البيانات ، أي القيمة أو الصفة التي لها أكبر تكرار ، أي التي تتكرر أكثر من غيرها من القيم أو الصفات .

The mode is calculated as follows:

A. Non-classified data:

في هذه الحالة لا توجد أي عمليات حسابية لإيجاد المنوال ، كل ما يتطلبه إيجاد المنوال هو معرفة القيمة التي تتكرر أكثر من غيرها.

Ex2: Find the mode value of the following data:

(11 , 12 , 9 , 10 , 9 , 11 , 13 , 9)

Sol:

The mode of these values is equal 9 ,because it is the value that is repeated the most.

إذا كان في البيانات منوال واحد فتسمى بيانات وحيدة المنوال ، وإذا وجد في البيانات منوالان فتسمى بيانات ثنائية المنوال ، وإذا وجد في البيانات أكثر من منوالين فتسمى بيانات عديدة المنوال ، وأحياناً لا توجد في البيانات قيمة أو صفة تتكرر أكثر من غيرها من القيم فتسمى بيانات عديمة المنوال .

Ex3: Find the mode for the data reported in each of the following groups:

Group 1: 18 , 14 , 15 , 19 , 20 , 24.

Group 2: 12 , 10 , 12 , 16 , 10 , 15 , 14.

Group 3: 12 , 12 , 15 , 13 , 15 , 13.

Group 4: 12 , 15 , 19 , 13 , 16 , 14 , 15 , 15.

Group 5: 10 , 9 , 15 , 12 , 5 , 11 , 5 , 10 , 12.

Sol:

For Group 1: There is no mode (بيانات عديدة المنوال)

For Group 2: 12 , 10 (بيانات ثنائية المنوال)

For Group 3: 12 , 13 , 15 (بيانات عديدة المنوال)

For Group 4: 15 (بيانات أحادية المنوال)

For Group 5: 5 , 12 , 10 (بيانات عديدة المنوال)

B. Classified data:

- لحساب المنوال لبيانات مبوبة في جداول تكرارية منتظمة أي فئاتها متساوية الطول نتبع الخطوات التالية :
- نحدد أولاً الفئة التي تحتوي المنوال ويطلق عليها الفئة المنوالية وهي الفئة المقابلة لأكبر تكرار.
 - نحسب قيمة المنوال باستخدام القانون التالي :

$$M = L + \frac{\Delta_1}{\Delta_1 + \Delta_2} \times C$$

حيث:

L: الحد الأدنى للفئة المنوالية (أو الحد الأدنى الحقيقي في حالة البيانات التي تمثل متغيراً منفصلاً)

Δ_1 = تكرار الفئة المنوالية - تكرار الفئة السابقة لها.

Δ_2 = تكرار الفئة المنوالية - تكرار الفئة اللاحقة لها.

C: طول الفئة المنوالية .

ملاحظة : إذا كان جدول التوزيع التكراري غير منتظم يجب تعديل تكراراته قبل تطبيق خطوات إيجاد المنوال التي أشرنا إليها .

Ex4: Find the mode value of the following data:

Class limits	f_i
50 – 54	10
54 – 58	30
58 – 62	90
62 – 66	60
66 – 70	20
Σ	210

Sol:

الفئة المنوالية هي الفئة (58 - 62) لأنها هي الفئة التي لها أكبر تكرار ونجد أن:

$$L=58 \quad C=4 \quad \Delta_1 = 90 - 30 = 60 \quad \Delta_2 = 90 - 60 = 30$$

$$M = L + \frac{\Delta_1}{\Delta_1 + \Delta_2} \times C$$

$$M = 58 + \frac{60}{60+30} * 4 = 60.67$$

Determine the Mode Graphically:

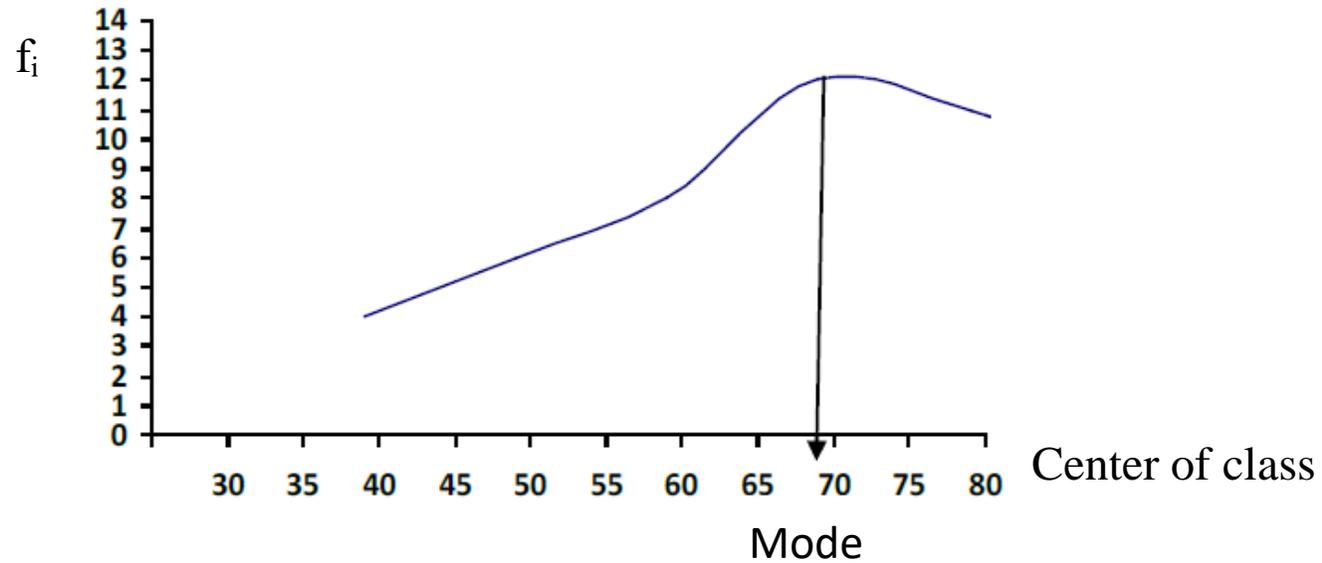
يمكن تحديد قيمة المنوال بيانياً باستخدام المنحنى التكراري أو المضلع فتكون قيمة المنوال هي القيمة المقابلة لقيمة المنحنى ، لأن القمة تمثل أكبر تكرار ، وحيث إن المنحنى يكون ممهداً باليد فغالباً ستكون قيمة المنوال التي نتحصل عليها بهذه الطريقة غير دقيقة ، كما يمكن تحديد قيمة المنوال باستخدام المدرج التكراري.

Ex5: Determine the mode Graphically of the following data:

Class limits	f_i
30 – 39	4
40 – 49	6
50 – 59	8
60 – 69	12
70 – 79	11
80 – 89	9

Sol:

Class limits	f_i	Center of class	Point (x , y)
30 – 39	4	34.5	(43.5 , 4)
40 – 49	6	44.5	(44.5 , 6)
50 – 59	8	54.5	(54.5 , 8)
60 – 69	12	64.5	(64.5 , 12)
70 – 79	11	74.5	(74.5 , 11)
80 – 89	9	84.5	(84.5 , 9)



Mode properties:

- أسهل مقاييس النزعة المركزية في حسابه .
- لا يتأثر بوجود قيم متطرفة .
- يمكن حسابه في حالة التوزيعات التكرارية المفتوحة بشرط ألا تكون الفئة المفتوحة هي للفئة المنوالية.
- يمكن إيجاد المنوال للبيانات النوعية.
- ليس له معنى إذا كانت البيانات قليلة العدد وقد لا يوجد أصلاً ، أما في حالة البيانات كثيرة العدد فله معنى معقول وله أهمية كبيرة وخاصة في عملية التسويق ، فمثلاً شركات تسويق الأحذية في مدينة ما لا تهتم بالوسط الحسابي أو بالوسيط بل تهتم بالمقياس الأكثر شيوعاً وهو المنوال .
- يمكن إيجاد المنوال بيانياً .
- قد لا يكون للبيانات منوالٌ وقد تحتوي على منوالين أو أكثر .
- يتأثر كثيراً بطريقة اختيار الفئات التكرارية للتوزيع ، فإذا غيرنا تقسيم الفئات لنفس التوزيع فيحدث تغيراً في التكرارات، وفي الغالب يحدث تغيراً في موقع الفئة المنوالية ، ولذلك نحصل على قيم مختلفة للمنوال .

H.W:

The following data shows the number of teachers in 25 primary schools in a city:

(30 , 24 , 23 , 27 , 28 , 31 , 21 , 25 , 18 , 29 , 29 , 20 , 27 , 28 , 26 , 26 , 24 , 23 , 27 , 28 , 22 , 26 , 26 , 28 , 30)

1. Create a frequency distribution table for this data, using 5 classes of equal length.
2. Find the up word combined frequency for this data.
3. Find the down word combined frequency for this data.
4. Find the arithmetic mean.
5. Find the median value for this data.
6. Determine the Median Graphically for this data.
7. Find the mode value for this data.
8. Determine the mode Graphically for this data.

Statistics

Lecture 9

- **H.W:**

The following data shows the number of teachers in 25 primary schools in a city:

(30 , 24 , 23 , 27 , 28 , 31 , 21 , 25 , 18 , 29 , 29 , 20 , 27 , 28 , 26 , 26 , 24 , 23 , 27 , 28 , 22 , 26 , 26 , 28 ,
30)

- 1.Create a frequency distribution table for this data, using 5 classes of equal length.
- 2.Find the up word combined frequency for this data.
- 3.Find the down word combined frequency for this data.
- 4.Find the arithmetic mean.
- 5.Find the median value for this data.
- 6.Determine the Median Graphically for this data.
- 7.Find the mode value for this data.
- 8.Determine the mode Graphically for this data.

Sol:

1.

Number of classes = 5.

Rang = $31 - 18 = 13$

Class length = $\frac{13}{5} = 2.6 = 3$

The minimum limit for the first class = 18

The actual minimum for the first class = $18 - 0.5 = 17.5$

The actual upper limit for the first class = $17.5 + 3 = 20.5$

The actual upper limit for the first class = $20.5 - 0.5 = 20$

$X_i = \frac{18+20}{2} = 19$

Table:

Class limits	Actual limits for classes	Centers of classes	Frequencies
18 – 20	17.5 – 20.5	19	2
21 – 23	20.5 – 23.5	22	4
24 – 26	23.5 – 26.5	25	7
27 – 29	26.5 – 29.5	28	9
30 - 32	29.5 – 32.5	31	3
		Σ	25

2.

Actual limits for classes	Frequencies	The up word combined frequency
Less than 17.5	0	فئة غير موجودة 0
17.5 – 20.5	2	2
20.5 – 23.5	4	6
23.5 – 26.5	7	13
26.5 – 29.5	9	22
29.5 – 32.5	3	25
		The total number of data (n) = 25



3.

Actual limits for classes	Frequencies	The down word combined frequency
17.5 – 20.5	2	25
20.5 – 23.5	4	23
23.5 – 26.5	7	19
26.5 – 29.5	9	12
29.5 – 32.5	3	3
More than 32.5	0	0

4.

$$\bar{X} = \frac{\sum_{i=1}^k x_i f_i}{\sum_{i=1}^k f_i}$$

$$\bar{X} = \frac{19*2+22*4+25*7+28*9+31*3}{25} = 25.84$$

5.

$$\text{Median arrangement} = \frac{\sum_{i=1}^k f_i}{2} = \frac{25}{2} = 12.5$$

$$m = L + \left(\frac{\frac{\sum_{i=1}^k f_i}{2} - F}{f_m} \right) \times c$$

$$L = 24$$

$$F = 6$$

$$f_m = 7$$

$$C = 3$$

$$\frac{\sum_{i=1}^k f_i}{2} = 12.5$$

$$m = 26.78$$

6.

نرسم رسماً بيانياً ونجد الوسيط حسب الطريقة المشروحة

7.

$$M = L + \frac{\Delta_1}{\Delta_1 + \Delta_2} \times C$$

$$L = 27$$

$$\Delta_1 = 9 - 7 = 2$$

$$\Delta_2 = 9 - 3 = 6$$

$$C = 3$$

$$M = 27.75$$

8.

نرسم رسماً بيانياً ونجد الوسيط حسب الطريقة المشروحة