

Estimates of location

Mean

Weighted mean

Median

Weighted median

Trimmed mean

Robust

outliers

Key Terms for Data Types

Continuous:

Data that can take on any value in an interval

Discrete

Data that can take on only integer values, such as counts.

Categorical

Data that can take on only a specific set of values representing a set of possible categories.

Binary

A special case of categorical data with just two categories of values (0/1 True/False)

ordinal

Categorical data that has an explicit ordering

columns - features

rows - records

Page No.

Date

Data → Sample → Sampling distribution → hypothesis
→ Conclusion.

Rectangular Data
Excel dataset

2 chapters → IA-I

Exploratory Data Analysis

Q.1 What are the elements of structured data

1. Continuous:
The data that can take any value in an interval like your weight.

2. Discrete

The data that can take only integers
Discrete data that cannot be converted into points

3. Categorical

Data that can take only a specific set of values representing a set of possible categories.
College having different streams.
(IT, AIDS, comps)

4. Binary

A categorical data with just two categories of values

Ex: Switch ON/OFF, PASS /FAIL

5. Ordinal.

A categorical data that has an explicit ordering like quick, rank. A group of students having roll nos from 1 to 60.

Q.3. What are the types of structured data?

The answer there are two types

Numerical, Categorical

Types of data in machine learning

Structured Data

Numerical

Categorical

discrete

continuous

Count of

Time duration

occurrence of events

* Purpose: The purpose of data analysis and predictive modelling, the data type is important to help determine the type of visual display, data analysis or statistical model. The data type for a variable determines how software such as python/R will handle computations for that variable. Knowing that data is categorical can act as a signal telling software how statistical procedures such as producing a chart or fitting a model should behave. 5. Storage and indexing can be optimized

Q. 3 What is Rectangular data?

Ans : 3) The typical frame of reference for an analysis in data science is a rectangular data object like a spreadsheet or a database table.

1. Dataframe:

Rectangular data is the basic data structure for statistical and machine learning models

2. Feature:

The column in the table is called as a feature.

It can also be called as Attribute, input, Predictor, and variable.

3. Outcome

Ex: Mr. Narendra Modi was elected as the next Prime minister of India.

Outcome can also be called as dependent Variable, Response, target or output.

4. Records

A row in a table is called as records.

Summarizing Point → Rectangular data is a 2dimensional matrix with rows indicating cases and columns indicating variables

	Processor	RAM/Rom	Camera	Display	Battery	Price
--	-----------	---------	--------	---------	---------	-------

Iphone 14 pro max	1.25	1.5	1.25	1.5	3.5	5
-------------------	------	-----	------	-----	-----	---

Samsung S23 Ultra	1.50	1.25	1	1	1.5	5
-------------------	------	------	---	---	-----	---

Pixel 7 Pro	2	3.5	1.75	2.5	3	3
-------------	---	-----	------	-----	---	---

Oneplus 11	2.5	3	3	5	2.5	2.75
------------	-----	---	---	---	-----	------

Asus ROG 6	1	1.25	2.5	2.5	2	3
------------	---	------	-----	-----	---	---

Overall	ASUS	Samsung / Asus	Samsung	Samsung	Samsung	Oneplus 11
---------	------	-------------------	---------	---------	---------	------------

Q.4 Non Rectangular Data Structures:

1. Time Series Data

2. Spatial Data Structures

which are used for mapping and location analysis.

the object representation : The primary focus is an object Ex. house, office, college and its spatial co-ordinates. The field view focuses on small unit of space and value of a relevant matrix.

3. Graph Data Structures:

Ex: Distribution Hubs connected by roads are an example of physical network. Graphs are used for network optimization.

Q.5 What are the estimates of location:

1. mean.

2. weighted mean.

3. median

4. weighted median

5. Trimmed median

6. Robust

7. outliers.

overall

Estimates of Location

weight w_{xi} date x_{i1}

1 17 17

Estimates of location:

1 16 16

2 12 24 mean = 13.53

2 14 28 median = 14

1 16 16 trimmed = 14.36

1 18 18 mean

3 10 30 outlier = 2

2 11 22

2 12 24

2 13 26

2 15 30

1 16 16

2 11 22

4 2 8

1 20 20

$$\sum w_{xi} = 27$$

$$\sum w_i x_i = 317$$

median: 2 10 11 11 12 12 13 14 15 16 16 16 17 18 20

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

Trimmed Removing outliers

mean:

$$\text{mean} = \frac{10 + 11 + 11 + 12 + 12 + 13 + 14 + 15 + 16 + 16 + 17 + 18 + 20}{14}$$

$$14.36$$

Weighted mean $\frac{\sum w_i x_i}{\sum w_i} = \frac{317}{27} = 11.74$

Estimates of variability:

1. Deviation

The difference between the observed value and the estimate of location.

2. Variance

The sum of squared deviations from mean divided by $n-1$ where n is the no. of data values also called as Mean Squared Error.

3. Standard Deviation

Square root of variance

4. Mean Absolute Deviation

The mean of the absolute values of the deviation from the mean

5. Median Absolute Deviation

The median of the absolute values of the deviation from the median.

6. Range:

The difference between the largest and the smallest value of a dataset.

7. Order Statistics

Matrix based on the data values sorted from smallest to largest

8. Percentiles

The value such that P percent of the values take on this value or less and $100-P$ percent take on this value or more.

9. Inter Quartile Range

The difference between the 75^{th} Percentile and the 25^{th} Percentile.

Calculate MAD, variance, standard deviation

Date	$x - \bar{x}$	Taking mod	Taking square
17	$17 - 15.5 = 1.5$	1.5	2.25
16	$16 - 15.5 = 0.5$	0.5	0.25
12	$12 - 15.5 = -3.5$	3.5	12.25
14	$14 - 15.5 = -1.5$	1.5	2.25
16	$16 - 15.5 = 0.5$	0.5	0.25
18	$18 - 15.5 = 2.5$	2.5	6.25
\bar{x}	$\overline{15.5}$	10	23.5

$$\frac{10}{6} = 1.66 \leftarrow \text{Mean Absolute Deviation}$$

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

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

$$\text{Variance} = 23.5 = \frac{23.5}{5} = 4.7$$

$$\text{Standard Deviation} = \sqrt{4.7} = 2.16$$

$n-1$ = Degree of Freedom

In statistics, there is always some discussion of why we have $n-1$ in the denominator, and in the variance instead of n we use $n-1$. $n-1$ is called as Degrees of Freedom

whether to consider n or $n-1$ will not have a major impact by analyzing large datasets however if you want to estimate about a population based on the sample set then n or $n-1$ will have a greater impact. If you use the denominators as ' n ' in the variance formula you will underestimate the true value of the variance and the standard deviation in the population. This is referred to as biased estimate. However if you divide $n-1$ instead of n , the standard deviation becomes an unbiased estimate.

Q. Fix question in IA-I

Show that Standard deviation is always greater than mean absolute deviation and why mean absolute deviation is always greater than median absolute deviation

$$SD > \text{mean Absolute Deviation} > \text{median Absolute Deviation}$$

Q. The runs scored in a cricket match by 11 players are as follows : 7, 16, 121, 51, 101, 81, 116, 9, 11, 6

1. Standard deviation
2. Mean absolute deviation
3. Median Absolute deviation

No. of observations = 11.

$$\text{Score} \quad |x - \bar{x}| \quad (x - \bar{x})^2$$

7

31.18

972.1924

16

22.18

491.9524

121

82.82

6859.1524

51

12.82

164.3524

101

62.82

3946.3524

81

42.82

1833.5524

1

37.18

1382.3524

16

22.18

491.9524

9

29.18

851.4724

11

27.18

738.7524

6

32.18

1035.5524

Summation

420

402.54

18767.61

$$\bar{x} = \frac{420}{11} = 38.18$$

① Standard deviation = $\sqrt{\text{Variance}}$

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

$n = \text{no. of observations}$

Page No.	
Date	

$$= \sqrt{\frac{18767.61}{n-1}}$$

$$= \sqrt{\frac{18767.61}{10}}$$

$$= 43.32$$

② Mean absolute deviation = $\frac{\sum |x_i - \bar{x}|}{n} = \frac{402.54}{11} = 36.59$

Arranging all the scores in ascending order, we get.

1, 6, 7, 9, 11, 16, 16, 51, 81, 101, 121
median.

$$\text{ita. } |(x_i - \text{median})|$$

$$15$$

$$10$$

$$9$$

$$7$$

$$5$$

$$0$$

$$0$$

$$35$$

$$65$$

$$85$$

$$105$$

$$121$$

Now arranging new marks in ascending order,

0, 0, 5, 7, 9, 10, 15, 35, 65, 85, 105

Median.

median Absolute deviation. = 10

By observing all those values, we can conclude

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

Estimates based on Percentiles.

Different approach to estimating dispersion is based on looking at the spread of the sorted statistics based on sorted rank data are referred to as ordered statistics. The most basic measure is the range (The difference between the largest value and the smallest value).

The maximum and the minimum values themselves are useful to know and helpful in identifying outliers.

Therefore, range is most sensitive to outliers. To avoid the sensitivity to the outliers we can look at the range of the data after dropping values from each end. This types of estimates are based on differences between percentiles.

25^{th} percentile is called as quartile 1

50^{th} percentile is called as quartile 2 / median

75^{th} percentile is called as quartile 3

The difference between 75^{th} percentile & 25^{th} percentile is called as inter-quartile range.

- Q. The scores in a test are 40, 45, 49, 53, 61, 68, 71, 79, 85, 91. What will be the percentile for these scores.

1. 40

2. 45

3. 49

4. 53

5. 61

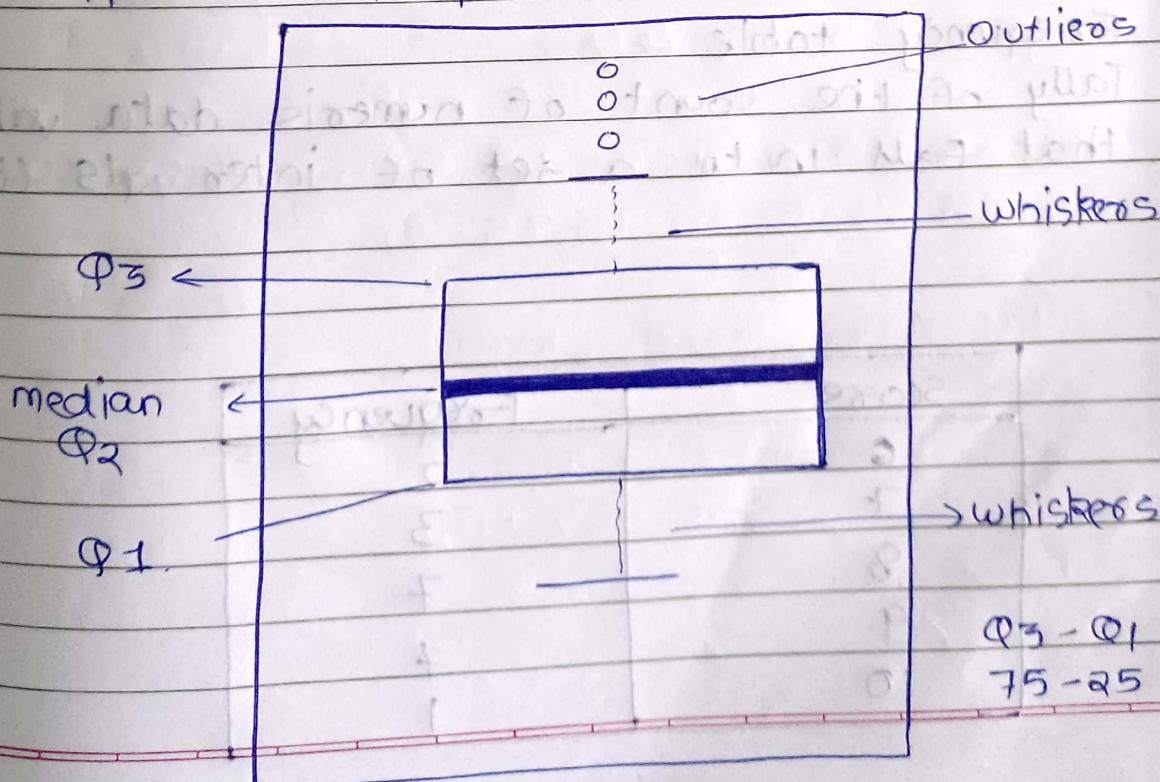
6. 68

7. 71 \rightarrow 70 percentile

8. 79 \rightarrow 80 percentile

9. 85 \rightarrow 90 percentile

10. 91 \rightarrow 100 percentile



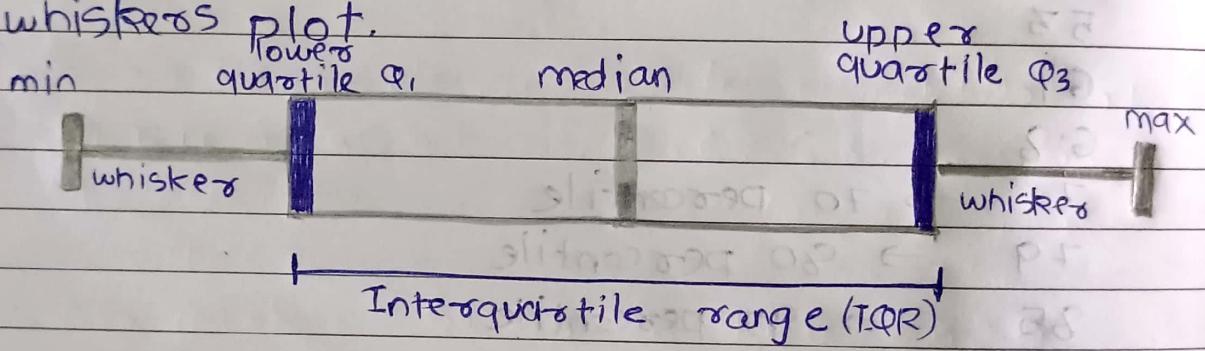
Data distribution

Each of the estimates that we have covered sums up the data in a singular single no to describe the location or variability of data

It is useful to explore how the data is distributed overall.

1. Box plot:

A plot introduced by John Tukey as a quick way to visualize the distribution of data based on percentile. It is also called as box and whiskers plot.



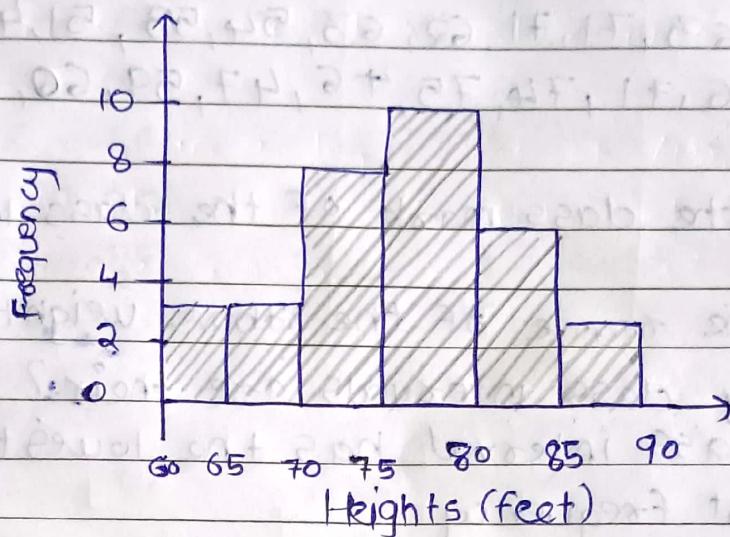
2. Frequency table.

Tally of the count of numeric data values that fall into a set of intervals (bins)

Score	Frequency
6	2
7	3
8	7
9	7
10	1

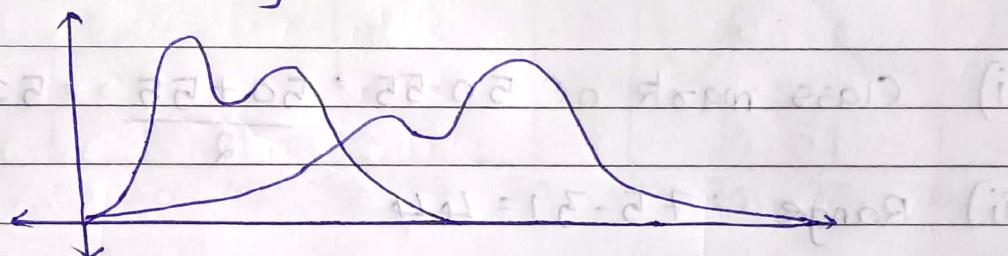
3. Histogram.

A plot of the frequency table with the bins on the x-axis and the count on the y-axis.



4. Density plot.

A smoothed version of the histogram, based on known density estimate.



Q. Construct a frequency distribution table for the following weights of 30 oranges using the equal class intervals. One of them is 40-45 (45 included). The weights are 31, 41, 46, 33, 44, 51, 56, 63, 71, 71, 62, 63, 54, 53, 51, 43, 36, 38, 54, 56, 66, 71, 74, 75, 46, 47, 59, 60, 61, 63

What is the class mark of the class interval 50-55?

What is the range of the above weights?

How many class intervals are there?

Which class interval has the lowest and the highest frequency?

Soln: 31, 33, 36, 38, 41, 43, 44, 46, 46, 47, 51, 51, 53, 54, 54, 56, 56, 59, 60, 61, 62, 63, 63, 63, 63, 71, 71, 71, 74, 75

i) Class mark of 50-55 = $\frac{50+55}{2} = 52.5$

ii) Range = 75 - 31 = 44

iii) 30-35, 35-40, 40-45, 45-50, 50-55, 55-60, 60-65, 65-70, 70-75, 75-80

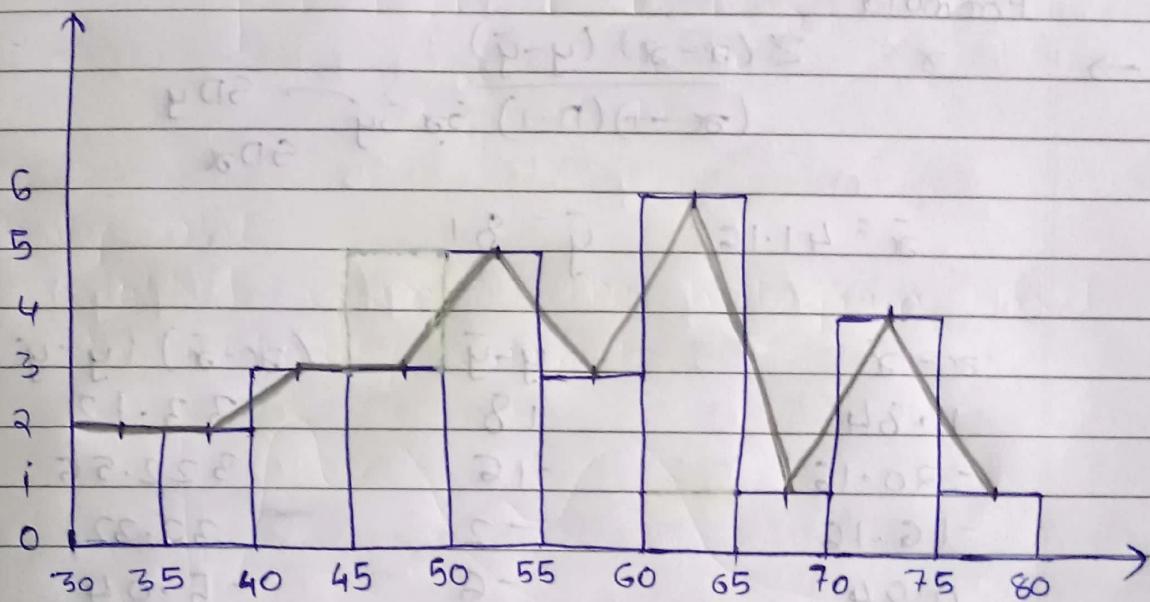
Total class interval = 10

frequency

30-35	2
35-40	2
40-45	3
45-50	3
50-55	5
55-60	3
60-65	6
65-70	1
70-75	4
75-80	1

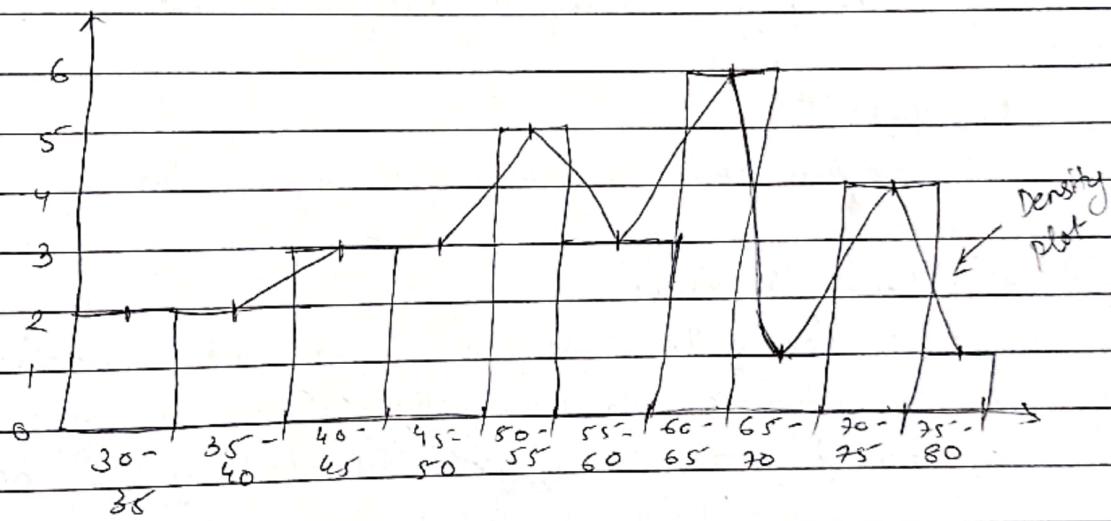
(6) → highest.

1 → lowest



(V)

	Frequency
30 - 35	2
35 - 40	2
40 - 45	3
45 - 50	3
50 - 55	5
55 - 60	3
60 - 65	6
65 - 70	1
70 - 75	4
75 - 80	1



- Q A survey of 36 students of a class was done to find out the mode of transport used by them while commuting to the college. The collected data is shown in the table given below. Represent the data in terms of bar graph.

Mode of transport	No. of student
Cycle	16
Bus	10
Car	10

→



Exploring Binary & Categorical Data.

1. Mode:-

The most commonly occurring value or category.

2. Expected value :-

When the categories can be associated with a numeric value, this gives an average value based on categories probability of occurrence.

3. Bar charts

The frequency for each category plotted as bars.

4. Pie charts

The frequency or proportion for each category plotted on wedges as pie.

a. The company offers 3 subscriptions plans. The basic version is priced at Rs 2500. Standard Version is priced at Rs 7000. And the pro version is priced at Rs 15,000. The probability of students signing for the course are 20%. for basic, 50% for standard and the remaining for pro version. Calculate the expected value.

$$\rightarrow \text{Expected value} = 2500 \times 0.2 + 7000 \times 0.5 + 15000 \times 0.3 \\ = 8500$$

Correlation

Exploratory data analysis (EDA) in many modeling projects (whether in data science or research) involves examining correlation among predictor variables like prediction and a target variable. Variables x & y are said to be highly correlated if high values of x go with high values of y , and low values of x go with low values of y .

If high values of x go with low values of y , the variables are -vely correlated.

Correlation Coefficient

A metric that measures the extent to which numeric variables are associated with one another ranges from -1 to 1.

Correlation Matrix

A table where the variables are shown on both rows & columns, & the cell values are the correlation b/w the variables.

Scatter plot

A plot in which the x-axis is the value of one variable and the y-axis is the value of another variable.

19/8/28

(*) Correlation Coefficient

Mod. 1 sum 16 Q. Find the value of the correlation coefficient from the given data

Subject	Age (x)	Glucose (y)	$(x - \bar{x})$	$(y - \bar{y})$
1	43	99	+1.84	18
2	21	65	-20.16	-16
3	25	79	-16.16	-2
4	42	75	0.84	-6
5	57	87	15.84	+6
6	59	81	17.84	0
	41.16	81		

Soln:- $\rho = \frac{\sum (x - \bar{x})(y - \bar{y})}{(n-1)S_x S_y}$

=

$$\sigma = (x - \bar{x})(y - \bar{y})$$

33.12

322.56

32.32

-5.04

95.04

0

478

$$(x - \bar{x})^2$$

1849

441

625

1764

3249

3481

$$(y - \bar{y})^2$$

9801

4225

6241

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

5.

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

$$= \sqrt{\frac{1240.83}{5}} = \sqrt{248.166} = 15.75$$

$$(x - \bar{x})^2$$

3.3855

406.4256

261.1456

0.6561

250.905

318.2656

1240.83

$$(y - \bar{y})^2$$

324

256

4

36

36

0

656

$$\sigma_y = \sqrt{\frac{\sum (y - \bar{y})^2}{n-1}} = \sqrt{\frac{656}{5}} = 11.45$$

$$\gamma = \frac{478}{5 \times 15.75 \times 11.45} \\ = +0.5299 \quad (-1.00 \text{ to } 1.00)$$

~~(S.M.)~~

Q Calculate the correlation coefficient for the following data.

x	y	$(x - \bar{x})$	$(y - \bar{y})$	$(x - \bar{x})^2 (y - \bar{y})^2$
4	5	-6	-4	36
8	10	-2	0	4
12	15	2	5	4
16	20	6	10	36
10	7.5			80

$$(x - \bar{x})(y - \bar{y})$$

$$+ 30 - 15$$

$$- 5$$

$$10 - 15$$

$$5 - 45$$

$$-25 - 20$$

$$\gamma = \frac{46}{5 \times 15.75 \times 8} = 0.5939$$

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

$$= \sqrt{\frac{80}{3}}$$

$$= 5.163$$

$$S_y = \sqrt{\frac{\sum (y - \bar{y})^2}{n-1}}$$

$$= \sqrt{\frac{50}{2}}$$

$$= \sqrt{25} = 5 \approx 0.4$$

$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$		
36	25		
4	0		
4	25		
36	.		
<hr/> 80	<hr/> 51		

$$SD_x = \sqrt{\text{Variance}}$$

$$SD_y = \sqrt{\text{Variance}}$$

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

$$= \sqrt{\frac{\sum (y_i - \bar{y})^2}{n-1}}$$

$$= \sqrt{\frac{80}{3}} = \sqrt{\frac{51}{3}}$$

$$= 5.16 \quad = 5.04941231$$

X

$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})(y_i + \bar{y})$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$
-6	-2.75	16.5	36
-2	2.25	-4.5	4
2	7.25	14.5	4
6	-6.75	-40.5	36
<hr/>	<hr/> -14	<hr/> 80	<hr/> 110.75

$$\sqrt{\frac{110.75}{3}}$$

$$SD_y = 6.0759 \quad SD_x = 5.16$$

$$\rho_{xy} = \frac{46}{3 \times 5.16 \times 6.0759} = \frac{46}{94.054932} = 0.489$$

$$\begin{array}{r}
 21 & 4 \\
 4 & 5 \\
 8 & 10 \\
 12 & 15 \\
 16 & 20 \\
 \hline
 40 & 30
 \end{array}$$

$$\bar{x} = \frac{40}{4} = 10 \quad \bar{y} = \frac{30}{4} = 7.5 \quad \bar{y} = \frac{50}{4} = 12.5$$

$$(x - \bar{x})(y - \bar{y}) \quad (x - \bar{x})(y - \bar{y})$$

$$\begin{array}{r}
 -6 & -7.5 & 45 \\
 -2 & -2.5 & 5 \\
 2 & 2.5 & 5 \\
 6 & 7.5 & 45 \\
 \hline
 & & 100
 \end{array}$$

$$\sigma^2 = \frac{100}{3 \times 5.16 \times 6.4549} = \frac{100}{99.921852} = 1.0007$$

$$(x - \bar{x})^2 \quad (y - \bar{y})^2$$

$$\begin{array}{r}
 36 & 56.25 \\
 4 & 6.25 \\
 4 & 0.25 \\
 36 & 56.25 \\
 \hline
 80 & 125
 \end{array}
 \quad SD_x = \sqrt{\frac{80}{3}} \quad SD_y = \sqrt{\frac{125}{3}}$$

$$SD_x = 5.16 \quad SD_y = 6.4549$$

Exploring 200 more variables:

- Estimators like mean and variance look at variables one at a time (univariate analysis)
- Correlation analysis compares 2 variables (bivariate analysis)
- Estimators which includes more than 2 variables are called as multivariate analysis.

1. Multivariate analysis

a. Contingency Table:

A tally of counts between two or more categorical variables

b. Hexagonal Binning

A plot of 2 variables with the records binned into hexagons

c. Contour plots.

A plot showing the density of 2 numeric variables like a topographical map.

d. Violin plots

Similar to the box plot, but showing the density estimate.

IA QB

- Chpt 1:
- difference between structured & unstructured data
 - what is the estimate of location
 - what is relation between median and robust estimator
 - what are the estimates of variability / dispersion
 - ^{show that} Standard Deviation $>$ mean Absolute Deviation
 $>$ median Absolute Deviation
 - what are the different tools used for data distribution.
 - what is expected value (sum)
^{theory}
 - what is correlation
 why it is called as bivariate analysis.
 - Numericals based on variance, SD, mean, mean absolute deviation, median absolute deviation
 - Numericals based on correlation coefficient.

Chpt 2:

- what are the types of samplings
 - methods of samplings
 - What is sample bias.
 - Data Quality v/s Data quantity. [Theory Q]
^[5 m]
 - What is sampling procedure
 - What is selection biased
- Answers includes
- a. data snooping
 - b. vast search Effect
- Why outliers / extreme observations tend to follow by mean / Regression to the Mean
 - What is Central Limit Theorem
 - What is standard Error
 - Resampling v/s Bootstrapping
 - What is normal Distribution
^(Gaussian) Also called as: Pg: 70
^{graph}

1. Demonstrate Q-Q plots for normal distribution and long tailed distribution
2. What is Student's t-distribution
3. What is binomial distribution
4. What is Chi-square test
5. What is F distribution
6. What is Poisson distribution
7. What is Weibull Distribution

$$\frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

Chapter 2: Random Sampling and Sample Bias

Random Sampling and Sample Bias

Sample: A subset from a larger dataset / population

Definition: A subset from a larger dataset / population

\bar{X} Population mean

\bar{x} Sample mean

N no. of observations of Population mean

n no. of observations of Sample mean

Random Sampling.

It is a process in which each available member of the population being sampled has an equal chance of being chosen for the sample at each draw.

or

The equal chance of being selected in the sample in each draw this is also called as simple Random Sampling

→ Sampling can be done with replacement or without replacement

Knowledge: Self Selection Sampling bias:

(sample)

Statistical bias: It refers to measurement or sampling errors that arise systematic and predict by the measurement or sampling process.

→ An important distinction should be made b/w errors due to random chance and errors due to bias.

\bar{a} or \bar{u}

\bar{x}

Sample mean v/s Population mean

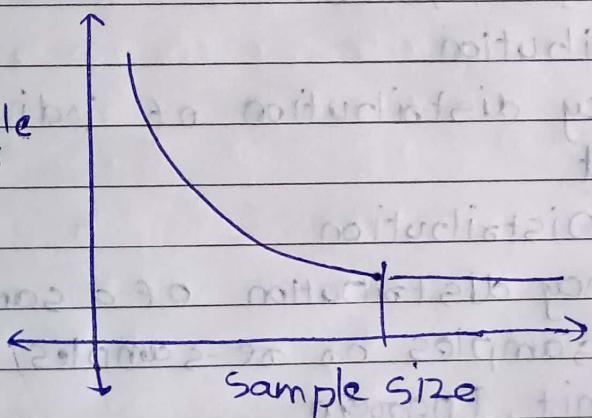
Selection bias refers to the practice of selectively choosing data consciously or unconsciously in a way that leads to a conclusion that is misleading.

Data snooping

Extensively hunting for data in search of something interesting.

Vast Search Effect:

Bias or non-reproducibility resulting from typical repeated data modelling or modeling data with large numbers of predictor variables.



Regression to the mean:

It refers to a phenomenon involving successive measurements on a given variable; extreme observations tend to be followed by the more controlled ones.

- 2 Attaching focus and meaning to the extreme value can lead to a form of selection biased.
- 3 Regression to the mean, meaning to "go back", is distinct from the statistical modelling method of linear regression, in which a linear relationship is estimated between dependent and independent variables.

Sampling Distribution of a Statistic

It refers to a distribution of some sample statistic over many samples drawn from the same population.

1. Sample statistic

A metric calculated for a sample of data drawn from a larger population

2. Data distribution

The frequency distribution of individual values in a dataset

3. Sampling Distribution

The frequency distribution of a sample statistic over many samples or re-samples

4. Central limit Theorem

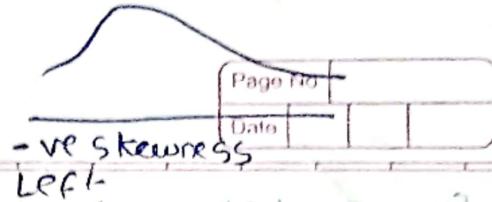
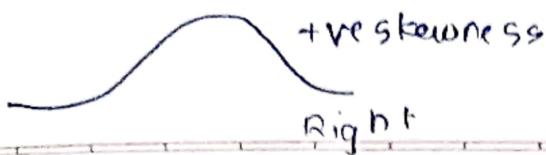
The tendency of the sampling distribution to take on a normal shape as sample size rises

5. Standard Error

The variability of sample statistic over many samples

~~trap~~

The variability between individual values is called standard deviation



(5m) Central Limit Theorem

The means drawn from multiple samples will resemble the familiar bell shaped normal curve, even if the source population is not normally distributed, provided that the sample size is large enough.

The Central Limit Theorem receives a lot of attention in traditional statistics text because it underlies the machinery of hypothesis test and confidence intervals, which themselves consumes half the space in such text. ~~many hardback~~
Data scientist should be aware of this.

~~Standard Error~~

The standard error is a single metric that sums up the variability in the sampling distribution for a statistic.

$$SE = \frac{s}{\sqrt{n}}$$

Standard deviation of sample values
Sample size

As the sample size increases, standard error decreases.

Collecting new samples to estimate the standard error is not very feasible. Instead of drawing new samples you can use bootstrap resamples.

The Bootstrap:

The effective way to estimate the sample distribution of a statistic is to draw additional samples with replacement, from the sample itself and recalculate the statistic or model for each resample. This procedure is called the Bootstrap.

Standard Error sums.

A certain property investment company with the an international presence, workers have a mean hourly wage of 12 dollars with a population standard deviation of 3 dollars. Given a sample size of 30, estimate and interpret the standard error of the sample mean.

$$SE = \frac{\sigma}{\sqrt{n}}$$

SD of population

$$= \frac{3}{\sqrt{30}}$$

$$= 0.5477$$

Explanation: If we were to draw several samples of 30 from the workers population and construct a sampling distribution of the sample means, we would end up with the mean of 12 dollars and a standard error of \$0.55.

Data quality vs Data quantity.
Example,

Assume that we have increased the sample size to 80 what will be the standard error.

$$SE = \frac{\sigma}{\sqrt{n}} = \frac{3}{\sqrt{80}} = 0.3354$$

Assume that we have decrease the sample size to 15

$$SE = \frac{\sigma}{\sqrt{n}} = \frac{3}{\sqrt{15}} = 0.77$$

Find the SE of estimate of mean weight of high school football players using the data given of weights of high school players

Players	Weights (in pounds)	$\bar{x} - \bar{\bar{x}}$	$(\bar{x} - \bar{\bar{x}})^2$
Player 1	150	-31.6	998.56
Player 2	203	21.4	457.96
Player 3	176	-5.6	31.36
Player 4	190	8.4	70.56
Player 5	168	-13.6	184.96
Player 6	193	11.4	129.96
Player 7	189	7.4	54.76
Player 8	178	-3.6	12.96
Player 9	197	15.4	237.16
Player 10	172	-9.6	92.16

$$\bar{x} = 181.6$$

$$\sum (\bar{x} - \bar{\bar{x}})^2 \\ = 2270.4$$

$$\sigma = \sqrt{\frac{\sum (\bar{x} - \bar{\bar{x}})^2}{N}} = 15.06$$

$$\sqrt{\frac{\sum (\bar{x} - \bar{\bar{x}})^2}{n-1}}$$

$$\sigma = 15.88$$

C30 we will consider sample.

Page No.	
Date	

$$\frac{\sigma_x}{\sqrt{n}} = \frac{15.88}{\sqrt{10}} = 5.02 \rightarrow S.E$$

σ_x sample Standard Deviation
 σ population Standard Deviation

- Q. Find the Standard Error of the estimate for the average no. of children in a household in your city by using the data collected from a sample of household in your city.

Household no.	no. of children	$(x - \bar{x})$	$(x - \bar{x})^2$
1	2	-0.25	0.0625
2	3	0.75	0.5625
3	1	-1.25	1.5625
4	0	0.25	0.0625
5	5	2.75	7.5625
6	2	-0.25	0.0625
7	1	-1.25	1.5625
8	4	1.75	3.0625
	2.25		$\sum(x - \bar{x})^2 = 19.5$

$$S.E = \frac{\sigma_x}{\sqrt{n}} = \frac{19.5}{\sqrt{8}}$$

$$\sigma_x = \sqrt{\frac{\sum(x - \bar{x})^2}{n-1}} = \sqrt{\frac{19.5}{7}} = 1.66$$

$$\sigma_x = 1.66 \approx 1.67$$



$$SE = \frac{\sigma}{\sqrt{n}} = \frac{1.67}{\sqrt{8}} = 0.59$$

Bootstrap esiditie e la efficienza

- i) Bootstrap Sample :
A sample taken with replacement from an observed data set.
 - ii) Resampling :
The process of taking repeated samples from observed data.

Household	Sample 1	Sample 2	3	4	5
1 element	1	2	3	1	2
2 individuals	2	2	4	2	2
3 individuals	2	2	4	4	2
4	2	3	1	3	2
5	4	4	2	2	2

The algorithm for bootstrap resampling of a mean for a sample size of n is as follows:

Step 1 Draw a sample value, record it and then replace it.

Step 2 Repeat n times for the string(s)

Step 3. Record the mean of the n flu sampled values

Step 4 Repeat steps 1, 2, 3 ~~A + R~~ times
iteration

Step 5: Use the R results to 1) calculate the standard deviation

- b. Produce the histogram or box plot
c. Find the confidence interval.

Resampling vs Bootstrapping

1. The Bootstrap is a tool for assessing the variability of sample statistic
2. The Bootstrap can be applied in a similar fashion in wide variety of circumstances
3. Bootstrap allows us to estimate sampling distribution for statistics where no mathematical approximation has been developed.
4. When applied to predictive models aggregating multiple bootstrap samples predictions outperform the use of single model.

Confidence Interval:

1. Confidence level

The percentage of confidence intervals, constructed in a same way from the same population that are expected to contain the statistic of interest.

2. Interval end points.

The top and bottom of the confidence interval

Note: The mean, median, variance, standard deviation for a sample or gives us or called as point estimate. Presenting an estimate not as a single number but as a range is called as confidence interval.

\bar{x} - sample mean
 μ - population mean

Page No.	
Date	

confidence intervals always comes with a coverage levels expressed as percentages e.g. 90%, 95%, 99%

x is a normally distributed variable with mean of 30 ($\mu = 30$) and standard deviation as 4
Find probability of a dataset which is greater than $P(\bar{x} < 40)$

$$P(\bar{x} > 21)$$

$$P(30 < \bar{x} < 35)$$

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

$$\textcircled{1} \quad Z = \frac{40 - 30}{4} \quad \textcircled{2} \quad Z = \frac{21 - 30}{4}$$

$$Z = \frac{10}{4}$$

$$Z = \frac{-9}{4}$$

$$Z = 2.5$$

$$0.9938$$

$$99.38\%$$

$$0.01222$$

$$\bar{x} >$$

we will minus the answer from 1 - 0.01222

$$0.98778$$

$$0.8$$

$$98.78\%$$

$$\textcircled{3} \quad P(30 < z < 35)$$

$$z = \frac{30 - 30}{4}$$

$$z = 0$$

$$z = \frac{35 - 30}{4}$$

$$z = 1.25$$

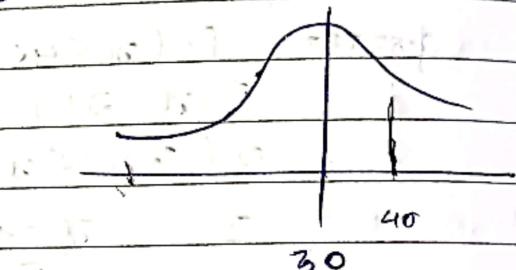
$$0 < p < 1.25$$

$$0.8943$$

$$0.8943 - 0.5000$$

$$0.3943$$

$$0.3943\%$$



Practice Problem:

13

A radar unit is measure speed of cars on a highway. The speed are normally distributed in a mean of 90 km/hr and standard deviation of 10 km/hr. What is the probability that a car picked randomly is travelling at more than 100 km/hr?

$$z > 100$$

$$z = \frac{x - \mu}{\sigma} = \frac{100 - 90}{10} = \frac{10}{10} = 1$$

$$0.8413$$

$$1 - 0.8413$$

$$0.1587$$

$$0.1587$$

$$15.87$$

Z-table - Normal Distribution Table

Page No.	
Date	

Player's Question:

Confidence Range:

95%

$$1.64 \times SE = 1.64 \times 5.022 = 8.23$$

181.

Confidence Interval

$$173.37 < p < 189.83$$

Children's Question:

99%

2.33 → From Table.

$$2.33 \times SE = 2.33 \times 2.25 = 2.33 \times 0.59$$

$$= 5.2425$$

1.37 to 1.92

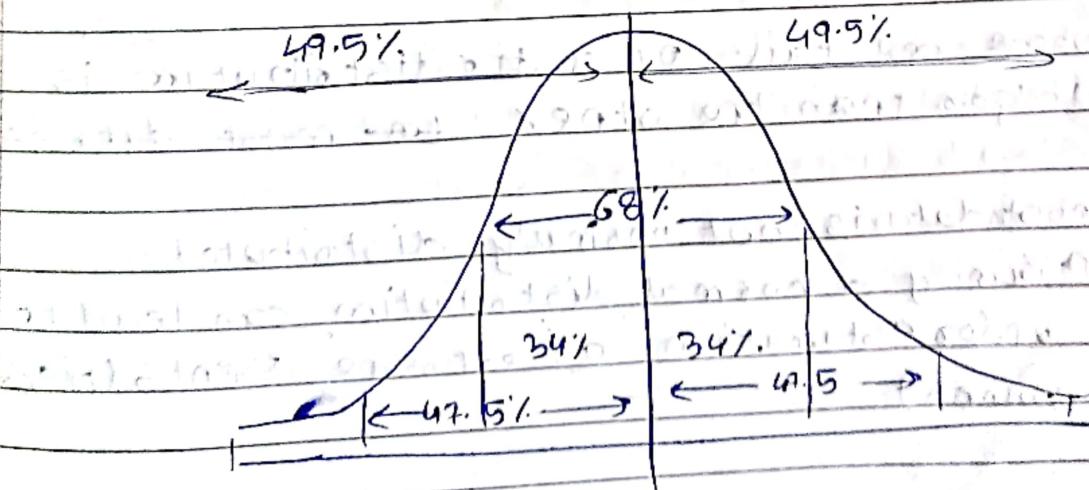
$\leftarrow p \leftarrow 7.4925$ or 1.92

$$0.88 \leftarrow \frac{t_0}{p} \leftarrow 3.62$$

Significant level and significance of table
Intermediate deviation is 99% of confidence interval

The bell shaped normal distribution states

that 68% of the data lies within one standard deviation of the mean and 95% lies within two standard deviation and 99% data lies within 3 standard deviation



Normal Distribution:

1. Error

The difference between a data point and predicted value

2. Standardize

Subtract the mean and divide by the standard deviation

3. Z-score

The result of standardizing an individual data point

4. Standard Normal

A normal distribution with mean 0 and standard deviation 1

5. Q-Q Plots

A plot to visualize how close a sample distribution is to a specified distribution

Long Tail Distribution:

1. Tail

The long narrow portion of a frequency distribution where relatively extreme values occurs at low frequency

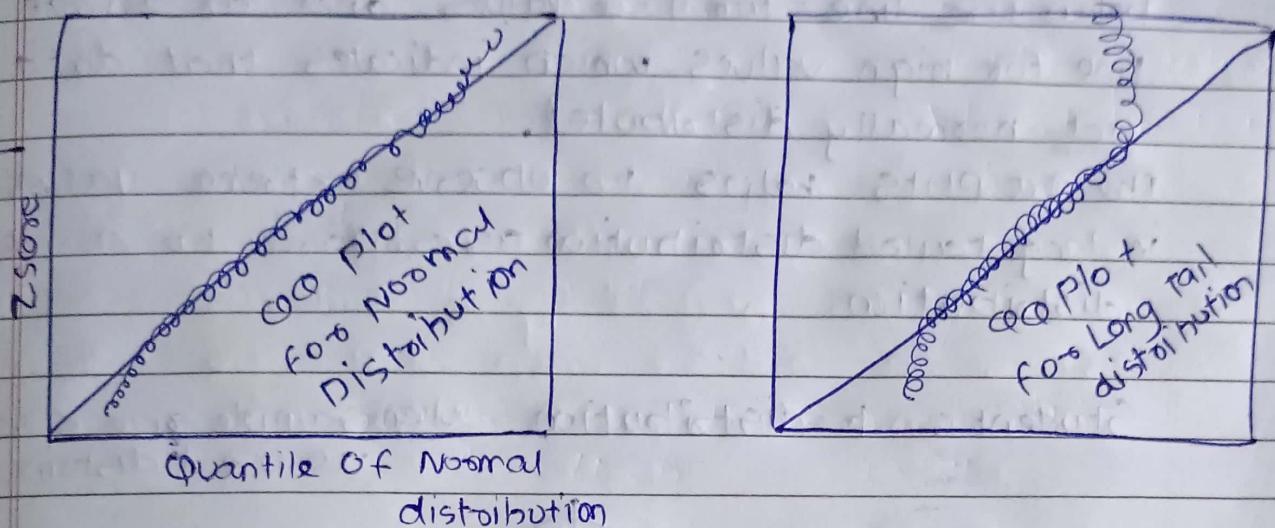
2. Skew

where one tail of a distribution is longer than the other (most data is)

1) most data is not normally distributed

2) Assuming a normal distribution can lead to underestimation of extreme events (black swans)

Student's t distribution



Student

The t distribution is a normally shaped distributed distribution ~~except~~ except that it is thicker and longer on the tails.

A Q-Q Plot is used to visually determine how close a sample is to a specified distribution.

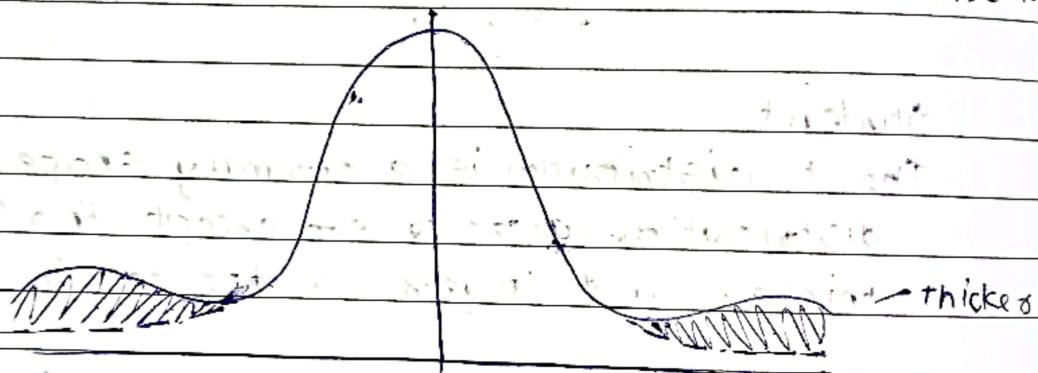
In normal distribution, Q-Q Plot orders the z-scores from low to high and plots each value's z-score on the y-axis; the x-axis is the corresponding Quantile of a normal distribution for that value's rank. Since the data is normalized, the units correspond to the no. of standard deviations away from the mean.

If the points roughly fall on the diagonal line then the sample distribution can be considered closed to normal

In long tailed distribution, the points are far below the line, for low values and far above the line for high values, which indicates that data is not normally distributed.

The QQ plots helps to observe extreme values in long tailed distribution as compared to normal distribution.

Student's t distribution. when sample size < 30
we use t distribution



- The t distribution is normally shaped distribution but thicker and longer on the tails
- It is used extensively in depicting distributions of sample statistics.
- The larger the sample, the more normally shaped t distribution becomes

The t distribution is widely used as a reference basis for the distribution of sample means, difference plus two sample means, regression parameters

$$\text{Variance } \frac{\sum (x - \bar{x})^2}{N}$$

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

Page No.
Date
Sample

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

Practice Problems

4 If die is rolled, find variance and standard deviation

	$(x - \bar{x})$	$(x - \bar{x})^2$
1	-2.5	6.25
2	-1.5	2.25
3	-0.5	0.25
4	0.5	0.25
5	1.5	2.25
6	2.5	6.25
Σ		17.5
21		

$$\bar{x} = \frac{21}{6} = 3.5$$

$$\text{Variance} = \frac{\sum (x - \bar{x})^2}{N} = \frac{17.5}{6}$$

$$\boxed{\text{Variance} = 2.917}$$

$$SD = \sqrt{\text{Variance}}$$

$$= \sqrt{2.917}$$

$$\boxed{SD = 1.707}$$

IA-I

Sums \rightarrow binomial \approx t, chi square distribution, normal distribution

$$t\text{-score} = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

t distribution sums:

If the sample mean and expected mean value of the marks obtained by 15 students in a class test is 290 and 300 resp. what is the t-score if the SD of the marks is 50?

$$t = \frac{\bar{x} - \bar{y}}{s}$$

$$t = \frac{290 - 300}{50} \\ t = \frac{-10}{50} \\ t = -\frac{1}{5}$$

$$t = \frac{-10}{12,909}$$

$$t = -0.7745$$

3. If the sample mean and expected mean value of height obtained by 15 students in a class test is $\frac{170}{290}$ and $\frac{165}{300}$ resp. What is the t-score if the SD of the heights is 21.05

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} = \frac{170 - 165}{\frac{21.05}{\sqrt{16}}} = \frac{5}{5.2625} = 0.950$$

Binomial Distribution

Properties:

1. The experiment consists of a sequence of n identical trials.
2. Two outcomes are possible on each trial. We can refer to one outcome as a success and the other as the failure.
3. The probability of a success is denoted by p and does not change from trial to trial. Consequently the probability of a failure is denoted by $(1-p)$ and does not change from trial to trial.
4. The trials are independent.

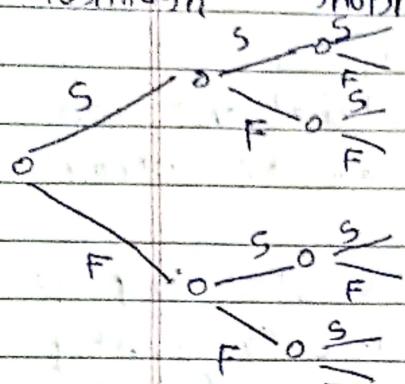
Example:

Suppose consider the experiment of tossing a coin 5 times and on each toss observing whether the coin lands with a head or a tail on its uppermost face. Suppose we have to count the no. of heads appearing over the 5 tosses. Does this experiment show the properties of binomial experiment?

- Ans:
1. The experiment consists of 5 identical trials i.e. each trial involves the tossing of one coin.
 2. Two outcomes are possible for a trial head or a tail. Head is can be success and tail can be failure.
 3. The probability of head and the probability of tail is same for each trial with p as 0.5 (50%)
 $1-p = 0.5$
 4. The trials are tossed for tosses are independent bcoz the outcome of any one trial is not affected by the outcome of other trials.

Binomial Probability Function.

Rushikesh Shubham Nairutya



Nairutya

SSS	$\frac{3}{3}$	$\binom{3}{3} = 3C_3 = 1$
SSF	$\frac{2}{3}$	
SFS	$\frac{2}{3}$	
SFF	$\frac{1}{3}$	$\binom{n}{\alpha} = \frac{n!}{\alpha!(n-\alpha)!}$
FSS	$\frac{2}{3}$	
FSF	$\frac{1}{3}$	
FFS	$\frac{1}{3}$	$\binom{3}{1} = 3$
FFF	$\frac{0}{3}$	$\frac{1}{2^1(3-2)} = \frac{1}{2}$

$$f(\alpha) = \binom{n}{\alpha} p^\alpha (1-p)^{(n-\alpha)}$$

α is no. of successes

p is the probability of success of one trial

$1-p$ is the probability of failure

n is the no. of trials

$f(\alpha)$ is the probability of α success in n trials

α values

$$f(0) = \binom{3}{0} p^0 (1-p)^{(3-0)}$$

$$\frac{3!}{0!(3-0)!} (0.3)^0 (0.7)^3$$

$$\frac{3!}{3!} \cdot 1 (0.7)^3$$

$$f(0) = 0.343$$

$$1 = \binom{3}{1} (p)^1 (1-p)^{3-1}$$

$$= 3! \cdot (0.3)^1 (0.7)^2$$

$$\frac{1!(3-1)!}{2!} \cdot (0.3) (0.7)^2$$

$$\frac{3 \times 2!}{3!} (0.3) (0.7)^2$$

$$2 = \binom{3}{2} (p)^2 (1-p)^{3-2}$$

$$= \frac{3!}{2!(3-2)!} \cdot (0.3)^2 (0.7)^1$$

$$\frac{3!}{2!1!} (0.3)^2 (0.7)^1$$

$$\frac{3}{3} (0.3)^2 (0.7)$$

$$0.189$$

$$3 = \binom{3}{3} (p)^3 (1-p)^{3-3}$$

$$= \frac{3!}{3!(3-3)!} (0.3)^3 (0.7)^0$$

$$\frac{1}{1} (0.3)^3$$

$$3 = 0.027$$

Adding $f(0) + f(1) + f(2) + f(3)$

$$0.343 + 0.441 + 0.189 + 0.027$$

$$= 1$$

Suppose for binomial experiment $n = 10$ $\pi = 4$
 Probability of success is 30%.

$$f(4) = \binom{10}{4} (0.3)^4 (0.7)^6 = 0.2001$$

Therefore the probability of making exactly 4 sales to 10 customers entering the store is 20%.

An agent sells life insurance to 5 equally aged, healthy people according to the recent data, the probability of a person living in this conditions for 30 years or more is $2/3$. Calculate the probability that after \uparrow success rate 30 years if

- i) all 5 people are still living
- ii) at least 3 people are still living
- iii) exactly 2 people are still living

$$p = 0.667$$

$$1-p = 1-0.667 = 0.333$$

$$n = 5$$

$$i) \pi = 5$$

$$f(5) = \binom{5}{5} (p)^5 (1-p)^{5-5}$$

$$\frac{5!}{(5-5)!} (0.667)^5 (0.333)^0$$

$$\frac{5!}{5!} (0.667)^5 (0.333)^0$$

$$\frac{1}{1} (0.667)^5 (1)$$

$$f(5) = 0.132$$

atmost ($F \leq 3$)

Page No.	
Date	

atleast

$$\text{ii) } F(x \geq 3) = F(3) + F(4) + F(5)$$

$$F(3) = \binom{5}{3} (0.667)^3 (0.333)^2$$

$$5! \quad (0.667)^3 (0.333)^2$$

$$\frac{5!}{(5-3)!} 3!$$

$$5! \quad (0.667)^3 (0.333)^2$$

$$\frac{5!}{2! 3!}$$

$$5 \times 4 \times 3! \quad (0.667)^3 (0.333)^2$$

$$\frac{5!}{2! 3!}$$

$$\frac{5!}{2!} \quad (0.667)^3 (0.333)^2$$

$$2$$

$$10 \quad (0.667)^3 (0.333)^2$$

$$10 \times 0.032$$

$$0.32$$

$$F(4) = \binom{5}{4} (0.667)^4 (0.333)^{5-4}$$

$$5! \quad (0.667)^4 (0.333)^1$$

$$\frac{5!}{(5-4)!} 4!$$

$$5! \quad (0.667)^4 (0.333)^1$$

$$\frac{5!}{4!}$$

$$5 \times 4! \quad (0.667)^4 (0.333)^1$$

$$\frac{5!}{4!}$$

$$5 \quad (0.667) (0.333)$$

$$0.32$$

$$F(5) = \binom{5}{5} (0.667)^5 (0.333)^{5-5}$$

Q In 800 families with 4 children each classify according to given criteria how many families would you expect to have

- 2 boys and 2 girls
- at least 1 boy
- no girl

Chi Square Test.

The chi square distribution is typically concerned with counts of subject or items falling into categories.

The chi square distribution is the probability distribution used in statistics for hypothesis testing and comparing observed and expected values. For example, you might want to test whether one variable (representing gender) is independent of another (boss promoted in job). The statistic that measures the extent to which result departs from the null expectation is the chi square statistic. The

Notes
IP-I

F distribution.

The F statistic measures the means which are greater than the expectation under normal random variation. It is the ratio of variability among the grouped means to the variability within each group.

TR
RAE

Poisson distribution:

The poisson distribution tells us the distribution of events per unit of time or space. When we sample many such units

Ex: It is useful when addressing queuing questions like (how much capacity do we need to be 95% sure of fully processing the internet traffic that arrives a server in any 5 sec period)

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

Exponential distribution:

Using the same parameters λ that we use in the Poisson distribution, we can also model the distribution of time between events i.e.

Ex1: Time between visits to a website

Ex2: Cars arriving at a toll plaza

Weibull distribution:

The weibull distribution is an extension of the exponential distribution in which the event rate is allowed to change as specified by a shaped parameter (β). If $\beta > 1$, the probability of an event increases over time. If $\beta < 1$, the probability decreases over time. Because the weibull distribution is used with time to failure analysis instead of event rate, the second parameter η (eta) (Greek letter) is expressed in terms of characteristics life instead of rate of events per interval.