

# ① Intro to BA : Role of Analytics for Data Driven Decision Making

→ Business Analytics means "means utilizing tools, techniques, and procedure to analyze past business data in order to gain insight and help in decision making & problem solving"

② Data is Everything, Analytics is necessary for survival.

Analytics can be used for process improvement (e.g. using it to reduce time to deliver order), problem solving (predicting customer cancellation & frauds), decision making (switch to a new market etc.).

② Introduction to concept of Big data Analytics.

# Big Data → refers to extremely large / complex sets of data that cannot be easily managed, processed or analyzed with traditional data processing tools.

→ The goal of big data analytics is to extract valuable insights, patterns and knowledge.

→ Big data is characterized by 3 Vs :-

① Volume (large amount), ② Velocity (high speed at which data is generated & processed) , ③ Variety (diff. type of data like structured / unstructured)

Sources of big data

• Social media data  
• Satellite data  
• Banking, financial data  
• Entertainment data  
• Internet data  
• Data from sensors, meters  
• Research data

④ Structured data :- Has clear structure and follows a regular sequence. Organized formatted in a way which is easily readable by both humans & machines. It generally follows tabular structure. e.g. Relational database entries :- Employee data in company database with columns Name, ID, Department, salary. ⑤ Spreadsheet :- Data in Excel & similar software like financial data with date, expense etc.

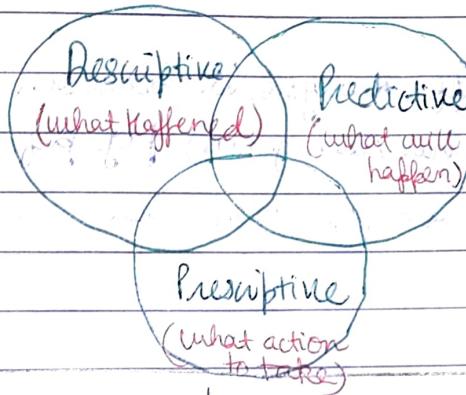
(Relational database is used to handle complex, large datasets while spreadsheets are simpler, used for basic data manipulation & smaller datasets and calculations)

② Unstructured data :- lacks formal structure. Doesn't conform to fixed schema and is often human generated.

Example: word docs, pdfs, video clips, images, audio recordings etc.

③ Types :- descriptive Analytics, Predictive Analytics

and Prescriptive Analytics



Descriptive is "Narrative function" - (covered in upcoming pages (one shot))

Descriptive	Predictive	Prescriptive
↓ Summarizes historical data to understand what happened in the past. [Tells you what happened in the past]	↑ use historical data & statistical algorithms to make prediction about future outcomes. [Take educated guess of what will happen in the future] Eg: Continuing with e-commerce one, it would forecast future sales or specific trends like sale of certain clothes ↑ during certain seasons.	↑ goes a step further and recommends actions to optimize or improve outcomes. gives advice / action to achieve desired result [Specific action to achieve best outcome] Eg: Predictive analytic told us sale would go up during this period I so prescriptive analytics might recommend raising inventory of those products, running targeted ads or discounts.
Q: You run a e-commerce business, it will involve looking at the past data(sales) and answer dues. like:- ① What were our best selling products last year? ② Which months had the highest sales?		

## Module 4 - Machine Learning

### Module 4 - Machine Learning

Both play crucial role in digital marketing and online commerce from these analytics, entities can make informed insights to optimize their online performance. Better management.

#### ④ WEB Analytics

→ It provides insights into performance of website ⇒ that how users are interacting with their online content. Tools like Google Analytics is used for this.

Key Aspects of this - Traffic analysis, user behaviour (time they spent), conversion tracking (sale/purchase/leads generate), page performance (load time), A/B testing (making two pages & compare with max images and seeing which perform better).

#### ⑤ Social Media Analytics

→ It involves collection and analysis of data from social media platforms. Its key aspects include - Audience insights, engagement metrics, sentiment analysis (+ve, -ve, neutral), influence impact, campaign performance etc.

→ Tools used are google analytics, facebook, insta, twitter, analytics, etc.

### ⑤ Online Machine Learning Algorithms

#### And Statistical Software Packages

#### ⑥ Machine learning Algorithms

They are set of rules/statistical model that enable computer to learn and make predictions without being explicitly programmed to do so. They learn from data and improve over time.

Eg. A very commonly used example of machine learning algo is "Spam email filter", which often uses a technique called Naive Bayes classification.

Example - Spam Email filter

Problem :- you receive lots of spam emails. You want to automatically filter out rather than check each manually.

Machine Based Algorithm :-

It learns patterns/probabilities of certain words being used in spam email like offer, free, discount, urgent and calculate probability of it being spam.

This is how machine learning can automate tasks.

#### ⑦ How Machine learning Operates :-

→ By training model on a dataset, allowing algorithm to identify patterns and relationships. The trained model can then be used to make decision when presented with new unseen data.

#### ⑧ Statistical Software / Software Packages

→ They are the type of computer program designed to perform statistical analysis on data. Kept, teachers, business analysts to make sense of data by providing various tools.

#### Benefits :-

① Helps in Data Analysis - Analyses trend/patterns in large, complex datasets.

② Visualization - Helps in visual representation of data.

③ Modelling - Helps in creation of statistical models.

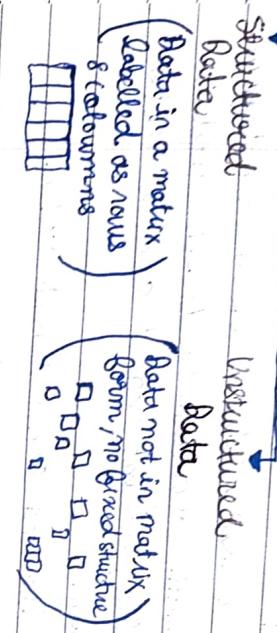
Popular Software - SPSS (statistical package for Social Sciences), R (R language), SAS (statistical analysis software), MATLAB (matrix laboratory)

# ○ Descriptive Statistics

①

Data

Based on Data Structure



Based on scale of measurement

Time

- Cross sectional data: Taking pictures at one moment in time (snapshot)
  - Eg. Survey 100 people on street today and ask their fav colour. It gives you picture of people's preference at that particular moment.
- Panel Data: (aka longitudinal) taking picture of the same group over a period, allowing to observe changes within that group.
  - Eg. Tracking Sales data made by your friend in month 20, 30, 40, 50, etc.

- Line series Data : Taking pic of same thing at regular time intervals like every 10 days to see how it evaluate fluctuate Eg. tracking stock price of a company on every 10 days.
- (Panel: over a period  
Time series i.e. regular intervals)

②

\* Types of Data Measurement Scales \*

Nominal

Ordinal

Interval

Ratio

If, with frequency;

- In it data refers to names which can be given a numerical value as frequent item (no necessarily means rank order).
- Eg. Single - 1 2nd order that is on 2nd order that is 1st is 2nd better than 3rd but the gap might be and 4th.
- Manus - 2 3rd better than 2nd but the gap might be and 2nd & 3rd are a lot.

- But here individual scale there is no true zero. Don't be 0 doesn't mean there is no zero. (it mean less worth).

③ Population and Sample

④ Population :- The entire grp you are interested in studying.

⑤ Sample :- Smaller grp, selected from population to represent & make inferences about the larger grp.

⇒ You want to find avg. height in college. You cannot ask 100 students ⇒ you select 10 students from each class and infer about form it about all.

★ Mean (Average)

⇒ It is the mathematical average value of data.

$$\text{Mean} = \bar{x} = \frac{n_1 + n_2 + \dots + n_n}{n} = \frac{\sum_{i=1}^n x_i}{n}$$

$$\sum_{i=1}^n (x_i - \bar{x}) = 0$$

★ "Median" is the ~~big~~ thing in descriptive stats which is not calculated based on the entire data.

### ★ Median

→ Divides data into two equal halves. Portion of observation below & above median will be 50%.

Individual Series

⇒ Arrange the data, find middle most value.

Discrete Series :-

Odd no. of observations,

$$\left(\frac{n+1}{2}\right)^{\text{th}} \text{ term} = \text{Median}$$

Even. no. of observations,

Continuous Series

$$\left(\frac{n}{2}\right)^{\text{th}} \text{ term} + \left(\frac{n}{2} + 1\right)^{\text{th}} \text{ term} = \text{Will give median value}$$

### ★ Mode

→ Mode is the frequency occurring value in dataset.

→ Mode is the only measure of central tendency which is valid for Nominal (qualitative data) since, the mean and median for nominal data is meaningless.

→ Eg. There is a data of marital status of customer namely (a) Married (b) Single (c) with children (d) without children. Mean, Median are meaningless here, but mode can tell what is the most occurring customer type.

Individual/Discrete Series

→ Series with highest frequency.

Continuous Series

$$\text{Mode} = l + \left( \frac{f_1 - f_0}{2f_1 - f_0 - f_2} \right) \times h$$

★ Modal class:- Class interval with highest frequency.

$l$  = lower limit of modal class

$f_0$  = frequency of class above modal class

$f_1$  = frequency of modal class

$f_2$  = frequency of class below modal class.

$h$  = class width

⇒ Modal class is the class which has the highest frequency.

④ Calculate CF.

\*  $l$  = lower limit of median class

\*  $f$  = frequency of median class

\*  $CF = CF$  of class above median class

\*  $h$  = class width

⇒ Median class is the class which has the highest frequency.

$$\# \quad \text{Mode} = 3 \text{ Median} - 2 \text{ Mean}$$

## (5.) Percentile, Quartile and Decile

$\Rightarrow$

Percentile, Quartile and Decile are used to tell the position of the observations (and its value) in the dataset.

and from it we can tell value.

④ Percentile score can be used in identifying 'per cent' of a student in class and in asset management.

$$\Rightarrow P_n =$$

gives value at which  $n\%$  of data is below that value.

$$P_{20} = \left[ \frac{2n}{100} (n+1) \right]^{\text{th}}$$

obs here,  $n = \text{no. of observation}$

## (6.) Measures of Variation

Time taken for burst of wire cut (mins)
2
22
3
24
5
24
4
26
45
21
41
75
78
76
66
67
77
78
78
79

$\Rightarrow$  Suppose this is a dataset.

(1) Company likes to know by what time 10% of the and 90% of the wire are cut.

$$\text{Ans: } n = 25 \text{ here, } P_{10} = \frac{10(25+1)}{100} = 0.6 \Rightarrow \text{th observation}$$

$\Downarrow$

$$\text{89, value at approx } (2.6)^{\text{th}} \text{ obs will be:} \\ \text{Set } = \text{value at } 2^{\text{nd}} + 0.6 \text{ (value at } 3^{\text{rd}} \text{ pos - value at } 2^{\text{nd}}) \\ = 3 + 0.6(5-3) \\ = 3 + 1.2 = 4.2 \text{ min Ans}$$

while,

$$P_{90} = \frac{90(25+1)}{100} = \frac{25 \times 9}{10} = \frac{225}{10} = 22.5 \text{ th obs}$$

$$\text{Value at } 23.4^{\text{th}} \text{ obs} = 23^{\text{rd}} \text{ obs} + 0.4^{\text{th}} \text{ obs} \\ = 78 + 0.4^{\text{th}} (\text{value at } 24^{\text{th}} \text{ obs} - 23^{\text{rd}} \text{ obs})$$

④

## Percentiles

## Quartiles

## Deciles



④ Special type of percentile which divide data in 4 equal parts.

④ Special type of percentile which divide data into 10 equal parts.

$P_2 = P_{25} =$  Contains first 25% of data

$P_{50} =$  Median

$P_{75} =$  Contains 75% of data

$D_1 = P_{10} =$  Contains 10% of data

$D_2 = P_{20} =$  Contains 20% of data

$D_3 = P_{30} =$  Contains 30% of data

This means the score spread across a range of 32 per cent

## V. Imp. FORMULA Trick

$$\text{Variance} = \sigma^2 =$$

② Inter Quartile Distance (IQR) [also k/a IQR (interQuartile)]

$\Rightarrow$  It is the measure of distance b/w  $Q_1$  and  $Q_3$ .

Eg. Dataset =  $[10, 15, 20, 25, 30, 35, 40]$

$$Q_2 = 25 \text{ (median)}$$

$\Rightarrow$  lower half of data =  $[10, 15, 20] = Q_1 = 15$

Upper half of data =  $[30, 35, 40] = Q_3 = 35$

$$IQR = IQR = 35 - 15 = 20$$

It means middle 50% of data has a range of 20 units.

★  $\Rightarrow$  IQR is used to identify outliers in data.

means, those values in dataset which differ significantly and thus prone to errors in conclusion.

$\Rightarrow$  Formula to identify outliers using IQR :-

- Lower Bound =  $Q_1 - 1.5 \times IQR$  (lower)
- Upper Bound =  $Q_3 + 1.5 \times IQR$  (upper)

$$\underline{\text{Eg.}} [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 90]$$

$$Q_1 = 3 \quad \Rightarrow IQR = 9 - 3 = 6 \quad \Rightarrow \text{outlier value} = 3 - 1.5 \times 6 = 3 - 9 = -6$$

$\rightarrow$  outlier value =  $Q_3 + 1.5 \times IQR$

$$= 9 + 1.5 \times 6 \\ = 9 + 9 = 18$$

Here, in our dataset there is a value which is higher than "18" (upper bound) which is  $90$ ,  $90$  is a outlier.

$$\text{Degree of freedom} = N - 1$$

③ Variance and Standard deviation

$$\text{④ Example } \sigma^2 = S^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{N-1}$$

Note,  $\bar{x}$  = sample mean  
 $\mu$  = population mean

(↑ in case of population variance)

$$\text{⑤ If we calculate variance, we divide by } n \text{ in case we are calculating variance for sample. This is k/a "Bessel's correction". It is denoted by } S^2 \text{.}$$

If we divide for calculating sample by "n" only then, it will give less value which is incorrect.  $\Rightarrow$  thus we also have downward bias. To come from it we divide by  $n-1$ .

⑥ Degree of freedom

Eg. You want to create a sample of 4 numbers with mean =

$$\text{⑦ 1st choice (first choice) } = 8 \\ \text{2nd choice (second choice) } = 12 \\ \text{3rd choice (third choice) } = 9 \quad \left\{ \text{sum} = 29 \right.$$

$\Rightarrow$  Now, from mean to be  $10$ , sum of 4 no. should be  $40$  then only  $\frac{40}{4} = 10$ . So 4th no. should be  $40 - 29 = 11$  or

$\Rightarrow$  so, we didn't have this freedom to choose this number. So,

## 7. Chebyshov's Theorem

⇒ Chebyshov Theorem | Inequality tells that how much your data lies b/w certain interval defined by mean & standard deviation.

⇒ Probability of finding a random selected value in an interval defined by  $\mu \pm k\sigma$  is  $1 - \frac{1}{k^2}$

$$P(\mu - k\sigma \leq x \leq \mu + k\sigma) \geq 1 - \frac{1}{k^2}$$

P=probability  
 $x = \text{value in dataset}$   
 $\sigma = \text{S.D.}$   
 $k = \text{no. of S.D. from mean}$

Q. Amount spent per month by a segment of credit card users of a bank has a mean value of 12000 and standard deviation of 2000. Calculate proportion of customers who are spending b/w 8000 and 16000.

Ans:

Given, mean =  $\mu = 12000$

S.D.  $\approx 2000$

Now,  $\mu \pm k\sigma \Rightarrow 12000 \pm 2(2000)$

$\Rightarrow 8000 \approx \text{from here } k=2$

so, Now,

$$P(\mu - 2\sigma \leq x \leq \mu + 2\sigma) = P(8000 \leq x \leq 16000) = 1 - \frac{1}{2^2} = 0.75$$

That is the proportion of customers spending b/w 8000 and 16000 is at least 0.75 or 75%.

## 8. Skewness and Kurtosis

### ★ Skewness

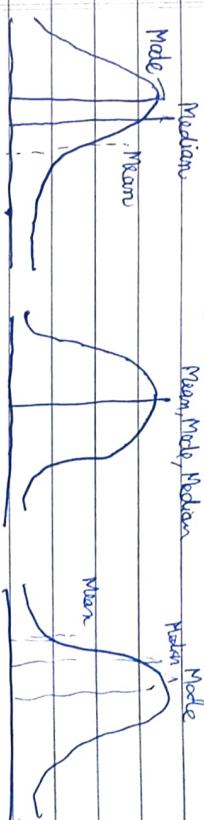
⇒ Skewness is the measure of asymmetry of probability distribution of a real valued random variable. In other words skewness tells us, the proportion of data b/w  $\mu$  and  $\mu + k\sigma$  is same as  $\mu$  and  $\mu + k\sigma$  ( $k = \text{some value constant}$ ). It tells us the extent to which distribution leans to the left or right of mean.

$$g_1 = \frac{\sum_{i=1}^n (x_i - \mu)^3}{n \sigma^3}$$

$x_i = \text{each value in dataset}$   
 $n = \text{no. of terms}$   
 $\sigma = \text{S.D.}$

$g_1 = +ve = +ve \text{ skewness} = \text{tail longer in right side}$   
 $g_1 = -ve = -ve \text{ skewness} = \text{tail longer in left side}$

$g_1 = 0 = \text{symmetrical distribution}$



(Mean > Median > Mode) (Mean - Median = Mode) (Mean < Median < Mode)

# Kurtosis

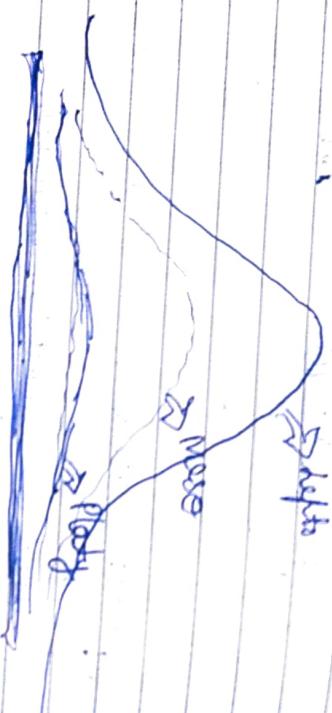
Kurtosis tells us the peakedness nature, as skewness told where the tail is.

$$(K) \text{ Kurtosis} = \frac{\sum_{i=1}^n (x_i - \bar{x})^4}{\sigma^4}$$

If,  $K < 3 \Rightarrow$  platykurtic

$K = 3 \Rightarrow$  mesokurtic

$K > 3 \Rightarrow$  leptokurtic



$\sigma^4$

$\Rightarrow$  If you are not given data in frequency form, make it.

Frequency

32.

1, 3, 22, 32, 5, 63, 49, 50, 20, 30, 46, 50, 27, 29, 30,

With draw a frequency:

Intervall	Frequency
0-10	3
10-20	4
20-30	5
30-40	5
40-50	3

11

10

20

30

40

50

Intervall

If we take mid values

5

4

3

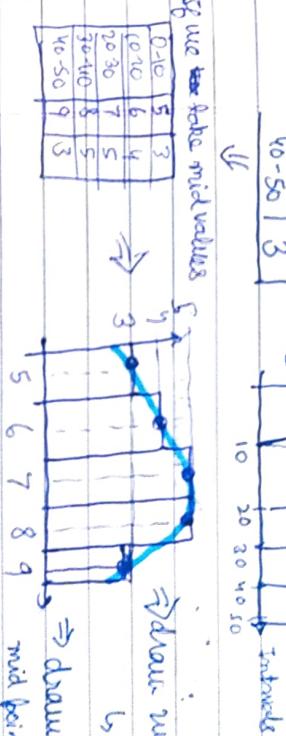
2

1

0

draw with free hand

or klo Frequency curve



⇒ draw st. lines through

mid points of intervals

(b) Kla "Frequency polyto make it closed since polygons are closed"

[make sure to extend freq. polyto make it closed since polygons are closed]

# Doline Curve

⇒ curve made by segment cumulative frequency (CF)

⇒ It can be less than or more than 100%

⇒ When they both meet, the intersection point is median. Although we can calculate median from each one separately above

Ex	0-10	3
10-20	27	
20-30	15	
30-40	3	
40-50	5	
	53	

Ex	more than 0	C.F.
more than 0	53	
more than 10	48	
more than 20	30	
more than 30	15	
more than 40	3	
more than 50	0	

Ex	more than 0	C.F.
more than 0	53	
more than 10	48	
more than 20	30	
more than 30	15	
more than 40	3	
more than 50	0	

Ex	less than	C.F.
less than 10	3	
less than 20	30	
less than 30	45	
less than 40	30	
less than 50	3	
	83	

Ex	less than	C.F.
less than 10	3	
less than 20	30	
less than 30	45	
less than 40	30	
less than 50	3	
	83	



#### ④ 5 Point Summary & Box-Plot (or Box Whisker Plot)

Concise description of dataset which gives 5 things:-

- ① Minimum Value
- ② Q<sub>1</sub>
- ③ Median
- ④ Q<sub>3</sub>
- ⑤ Maximum Value

⇒ graphical representation of 5 point summary is Box plot.

Ex Data = [55, 66, 71, 75, 80, 85, 90, 95, 99]

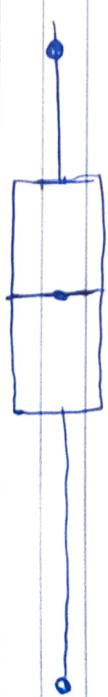
$$\text{① Min} = 55$$

$$\text{② } Q_1 = P_{25} = \frac{25(5+1)}{100} = \frac{1}{2}(Q_5)^{\text{th}} = 2^{\text{th}} + 0.5(3^{\text{th}} - 2^{\text{th}}) = (6 + 0.5)(71 - 66) = 66 + 2.5 = 68.5$$

$$Q_2 = 80$$

$$\text{③ } Q_3 = P_{75} = \frac{75(5+1)}{100} = \frac{3}{4}(Q_9)^{\text{th}} = 7^{\text{th}} + 0.5^{\text{th}} = 90 + 0.5(95 - 90) = 92.5$$

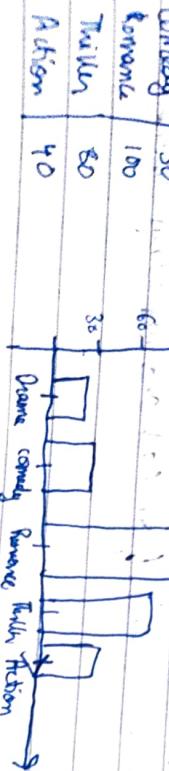
$$\text{④ Max} = 99$$



② Bar Graphs

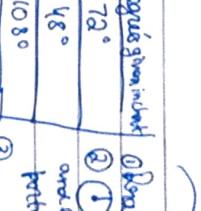
↳ used in case of Nominal data where names are associated with numerical value.

Eg Category People



#### ③ Pie Chart

Ex	subject	frequency	Degrees subtended
Social	6	72°	① draw a circle on paper. Put 'O' on it.
P.E.	4	48°	② draw n line segments. Put 'O' on it and start working angles using the sector
Maths	9	108°	and start working angles using the sector portion of it.



④ Box Plot

Box Range from Q<sub>1</sub> to Q<sub>3</sub>. Dihes from max & min and median is marked inside box. If case of outliers. They are displayed as point. Pen