

"what they are"

\* Data Science → Mining large amounts of data to identify patterns

↳ + PL, Stats, ML, algorithms

\* Big Data → ✓ Huge Data Volume

+ data capture, storage sharing  
processing -

Data Analytics → process and perform statistical analysis of data

✓ How data can be used to draw conclusions, solve problems

① → "what they do"

Data Scientist → predict future based on past patterns

→ capture / examine data from multiple disconnected sources

→ Develop New Analytical Methods,  
ML Models

② Big Data Prof: Analyse by Bottlenecks  
large scale data processing  
systems creation

✓ Architect scalable distributed  
systems

- \* Data Analyst - Acquire, process, summarize data
  - package data for insights
  - Design / create data reports using various reporting tools

→ where used?

\* Data Science → Search Engines, Fin. services, E-commerce.

Big Data → Fin., Comm., Retail

Data Analytics - Healthcare, Travel, IT Industry

→ what skills?

\* DS → Prog skills - Python, Stats, Math, ML, Hadoop, SQL, Shiny, R, visualization

Big Data → Java, Scala } PC, NoSQL - MongoDB, Cassandra  
 → Apache Hadoop framework, distributed systems

→ Data Analytics → Prog skills - R, Stats, Math, visualization skills

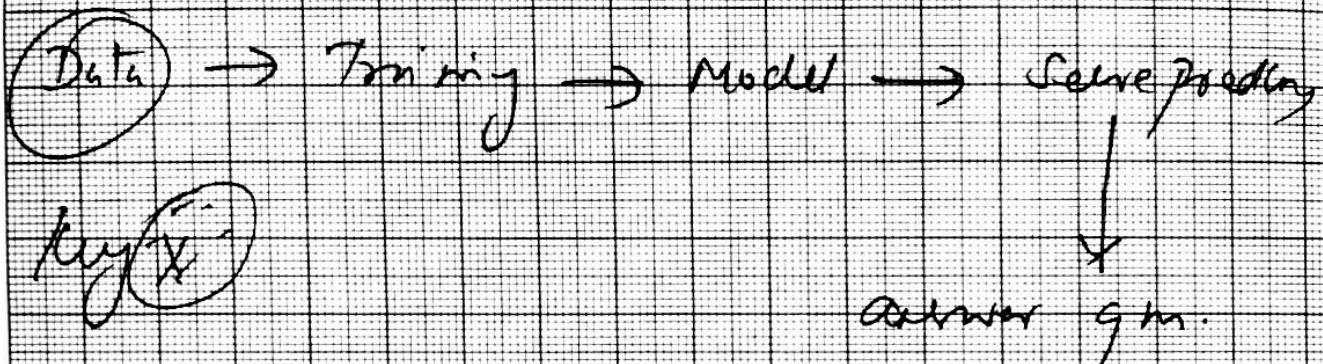
\$123K | \$88 | \$61



✓ Tagging people and objects

✓ Google Search  
 driving, ⇒ Recommendation systems, Self driving  
 Use a website for a company  
 ML is mandatory for products.

⑧ ML = using data / to answer question  
"Training" / "predictive / Inferential"



7 steps of ML :-  
wine or beer?  
Colours, % of features

Colour | & | wine or beer -

Data preparation

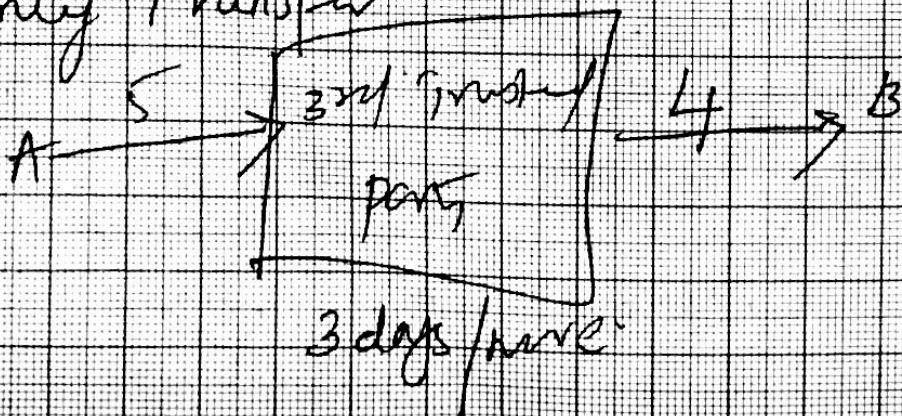
Training / Evaluation

- ✓ Internet of Things (IoT)
  - ✓ Cloud Computing
  - ✓ Big Data
  - ✓ Blockchain  $\leftarrow$  new ~~X~~
  - ✓ Artificial Intelligence
- ~~Machine Learning~~

Blockchain :=

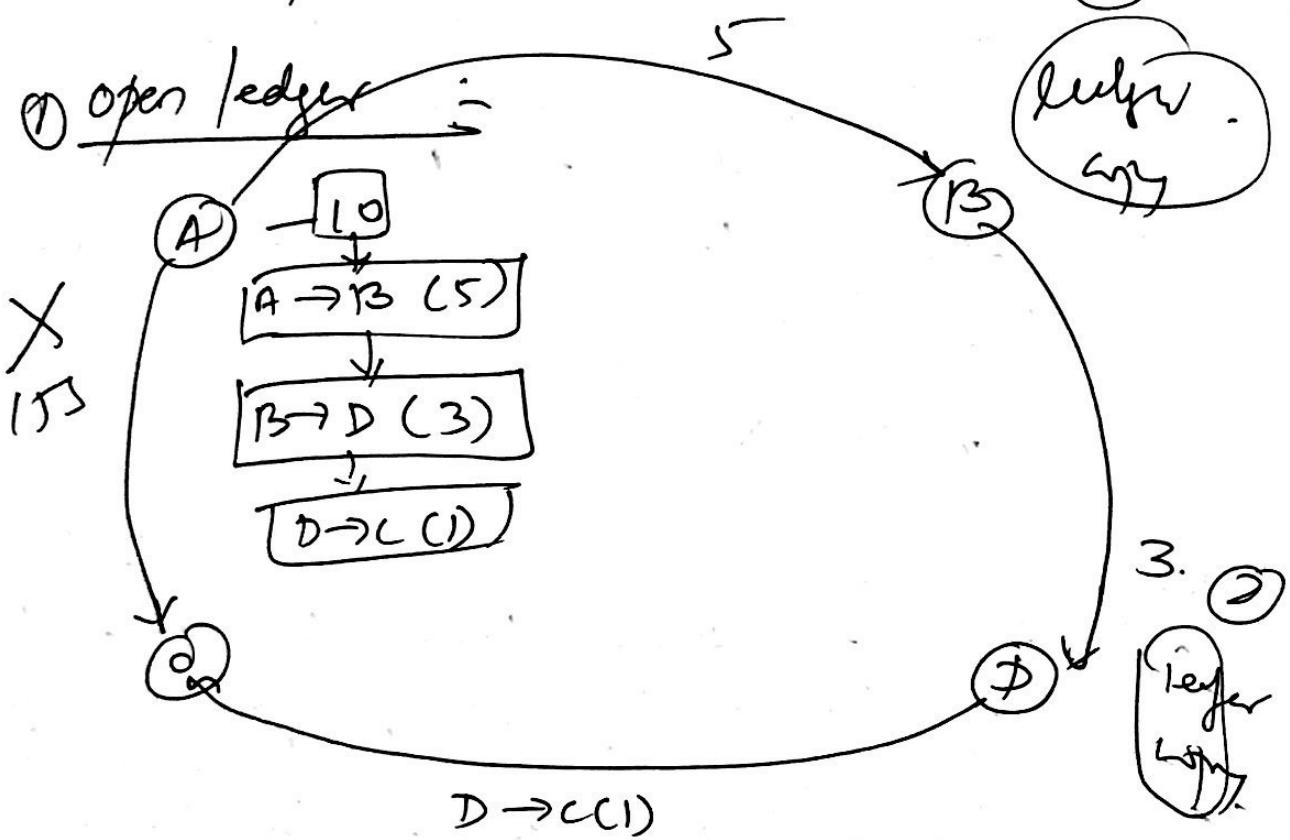
Not coins  $\neq$  Blockchain  
 $\downarrow$   
 Big Money enables movement of digital assets

✓ Money Transfer



- \* No third party
- \* Immortal
- \* Cheaper

Transfer → Blockchain



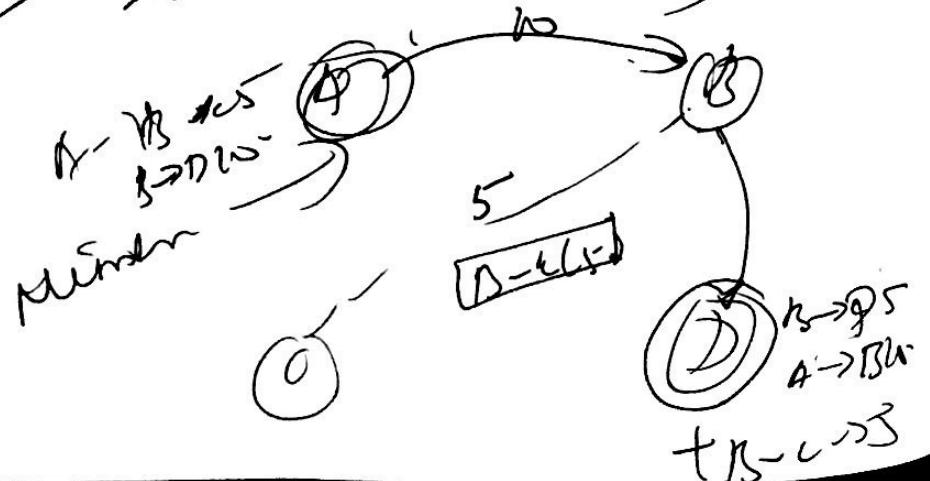
✓ Chain of X is open to everyone

## ② Distributed Ledger

- ✓ multiple ledger copies
- ✓ All parties are same Verity

### of ledger

- ✓ Validate
- ✓ Key



- ✓ Descriptive statistics, Exploratory Data Analysis
- Data description
  - Data visualization techniques
  - Graph forms, Description of a Variable, Relationship between two or variables graphically
  - Data summarization. — ||Descriptive||

Descriptive → Central Tendency Measures  
Mean, Median, Mode.

— Measures of Dispersion,

+ Prob Distn for engineering data  
— More than Mean / Median /  
20% data <  $x_i$ ; 30% data <  $y_i$  . . .

∴ Not only idea of dispersion, centrality —  
entire characterization of the data set

- ✓ Inferential statistics — data is a sample  
from a larger population

Ht of IITDM students — finite data - but  
may not be available completely  
— Sampling of 50 students

→ Descriptive alone may not suffice

Use of Amra in Inferential Statistics.

→ Something about the population

e.g. avg ht of IITDN students  $\leq 150$ , = 140

Sample observation → population prediction  
is Inferential statistics.

" Sample Observation → population Inference  
finite / infinite population

### ✓ Regression Analysis

↳ Inferencing, Descriptive, optimizations

↳ Overlap with m/c learning

↳ SLR → Dep (o/p) variables relationship with  
indep (i/p) Variables

↳ line fitting for data

↳ prediction

(Indep) | rainfall → crop yield (o/p var)

speed → mileage

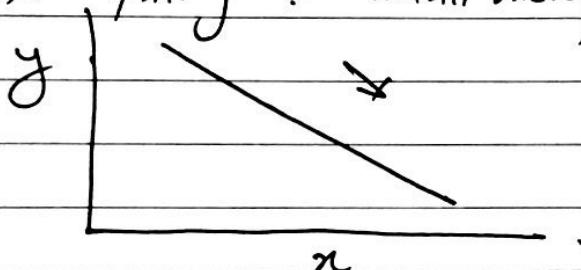
exercise → weight loss

examples

Reg. Analysis  $\rightarrow$  Relationship b/w dep var (y) with 1/more indepen. variables (x)

x axis - Indep var - Rainfall  
y "  $\rightarrow$  crop yield - Dep.

Line fitting :- Relationship b/w two variables.



" Given Indep variable can I predict Dep Var?"  
 $\downarrow$  Regression Answer.

ML (Essence of prediction)? Accuracy of given,

" Historic Data"  $\rightarrow$  ML

(2) Major part Next - Machine Learning

Data Mining / Pattern Recognition / Statistical Learning related topic, Not same.

ML  $\rightarrow$  focus on algorithm that converts data to knowledge.

(L2 Contd)

Regression for Continuous Var. others.

Classification  $\rightarrow$  o/p variable is discrete

Categorical — Nominal

" ADT of Supervised / Unsupervised learning



Clustering



Association Rule  
Mining

clustering  $\rightarrow$  task of grouping set of objects

into clusters / groups based on

similarity defined across common set  
of attributes / features.

Variety of Distance (metric).

k-means, Hierarchical, etc - .

ARM another USL :-

— relationships between feature across  
a lot of objects.

Machine Learning → Broadly classified as supervised / unsupervised.

↪ SL → Create a function / relationship from training data : one explicit output variable - atleast. labelled data / training data

"Algorithmic mapping between I/o/p Variable and

✓ Unsupervised → Task of creating patterns from data which have no explicit measure, patterns from unlabelled data (No o/p Variables)

Sup. learning problems ↴ Classification  
Regression problems.

SL → establish relationship between I/p and O/p.

Classification — o/p variable is Discrete Category and not a Nominal Categorical Variable  
(No explicit Orderly of classes)

e.g. Male / female opposed to Amount of crop yield prediction  
↪ Continuous Variable.



Grouping of customers based on  
age, gender, nationality.

- Design of experiments - Data gen.
- Active learning
- Reinforcement learning

Active learning → (Some data) → Start w/  
of  
look @ it to understand type  
of data required — DUE.

+ Deep learning NN, etc.



## Lecture 03

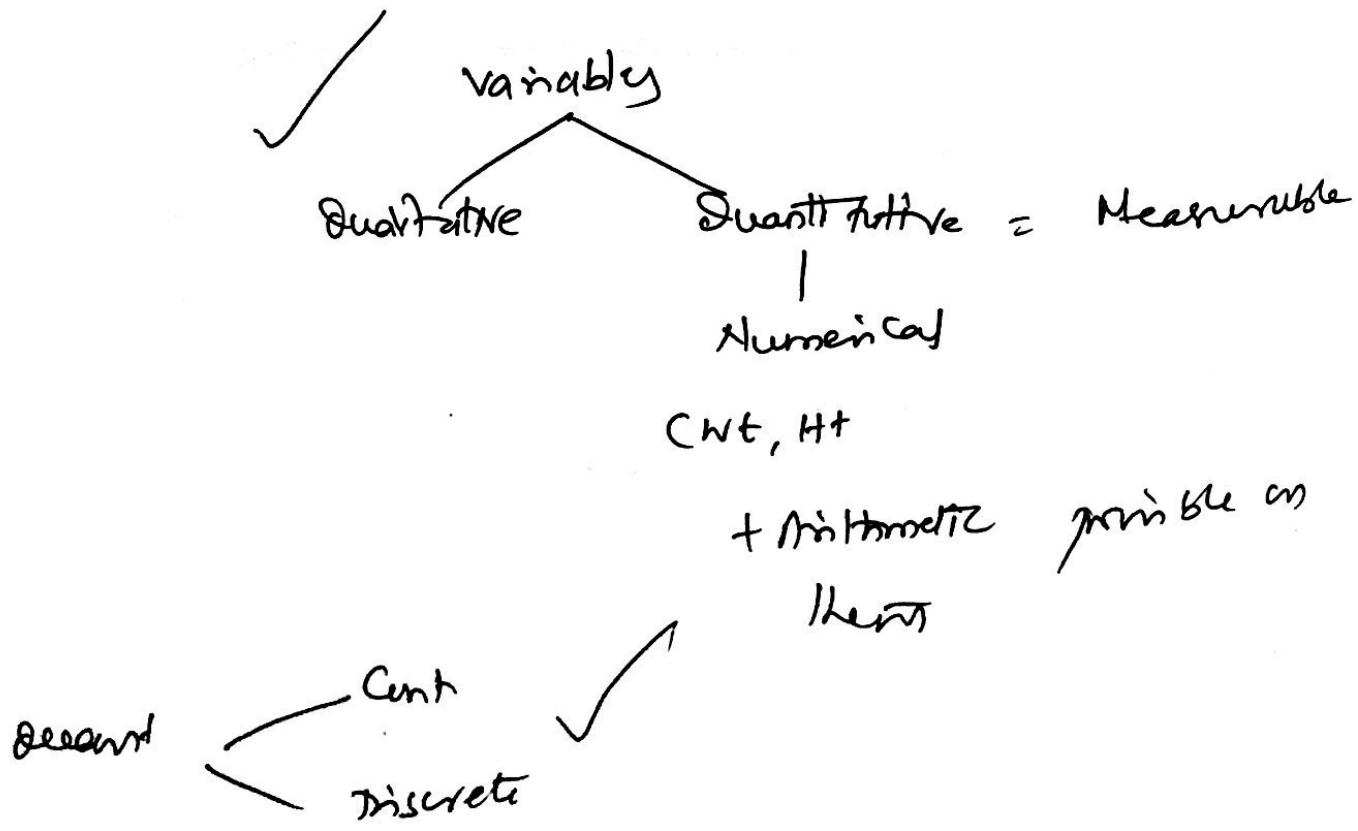
~~Descriptive statistics~~: Science of learning from data

- Data — Numbers / Text / Symbols — represent some information
- 'values' of Qualit. / Quant. Variables.
- Num / Quant. ( cont / Discrete)  
Catego / Qualit. ( Always discrete)  
Nominal , Ordinal .

Des. stats — varies according to data types

Data — No / Text / Symbols  
↳ frequent .

Data = Values .



$CV \rightarrow$  within interval any value is possible  
 $\rightarrow$  one square within a certain interval

1000 students;  $Ht$  (student/s)

(lowest - highest).

e.g.: (120 - 140 cm)

135 possible, 135.5 also possible.

$\rightarrow$  No measuring scale for a part accuracy  
"Nothing that prohibits 135.658 for height"

$CV \rightarrow$  any value between certain intervals

No. of uniq. people entering MRTDY/day :-

Multiple entries not  $x$

No of uniq. ppl per day  $\rightarrow$  variable of interest  
for a year  $\rightarrow$  365 data points / one per day

0 possible) lower bound; up bound - e.g. 10000

only 0, 1, 2, ...

1.5 not possible.

$\therefore$  Discrete Variable

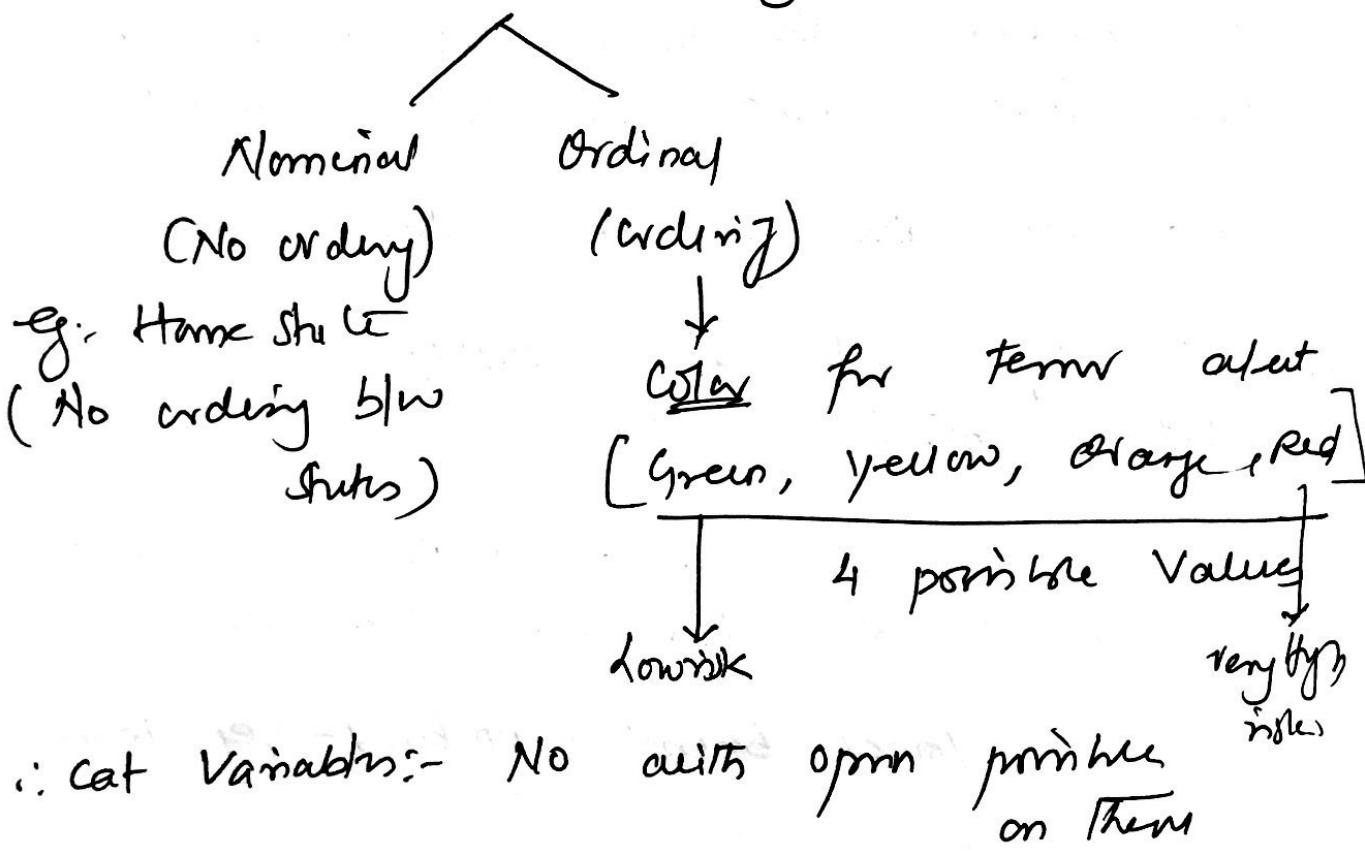
~~Discrete in sense that only integer values are possible~~

Qualitative:- represents characteristics that can be grouped / characterized

e.g.: Gender (male/ female)

Home town → Aniti set of names  
= diff states of India

Cat. Variable (Inherently) discrete



Greater / Lesser | Comparison possible with  
ordinal variables.

Descriptive Statistics:-

Quantitatively describe data

Means — visual → Graphs / tables

"Crunch the data"

✓ 1/2/3/4 types / numbers that quickly summarise the data.

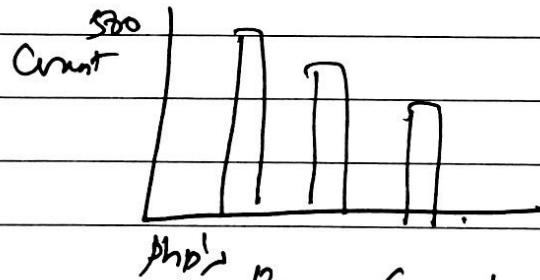
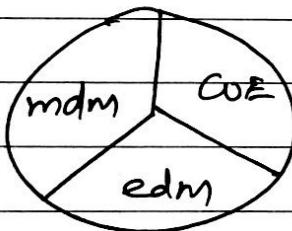
✓ data about data / metadata → Descriptive stats.

Descriptive → Cannot make conclusions beyond the data one has

"No generalisation / predictive / influence possible"

lot of data → Concise Repn. required for simpler interpretation — graphs, nos, etc.,

Des. stats can be for multiple variables / also.



Pie of Branch st. Dens.

Bar Graph

✓ describe each variable / also show inter-relationship between variables.

Graph repn of Single Variables (absc)

Highest education (cat. variable) 4 values — (School / Bach / Masters / Doctorate)

Highest education vs frequency of occurrence

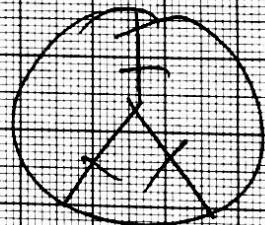
↳ Quick summary of max, least, average, etc.

another view ↳

↳ BSc < Masters < PhD (years)

⇒ Ordinal ?

Categorical variables — another way is  
the chart.



Coincidentes

✓ Not suited for Ordinal / Nominal  
Categorical variables

✓ More departments — the last better

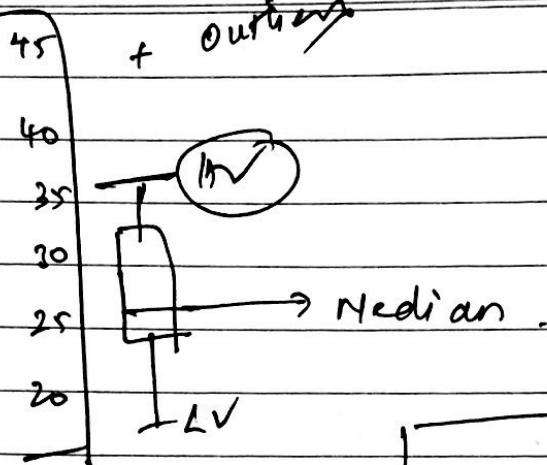
Quant: — Boxplot + Num. data

(range, variance, ...)

↳ Central tendency Captured  
(median)

This is  
feet  
(inches based)  
Content

Repeated Next page  
from slide



range of (25 - 33)

Histogram → Ricket for num. data

n axis - diff values



Column - Rows

6 points.

→ empirical construction

multiple variables → scatter plots, Box plots, Contingency

✓ (2 Quant)

(1 cat, 1 Quant)

p45

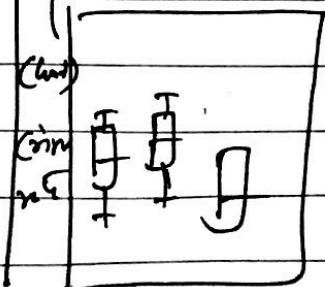
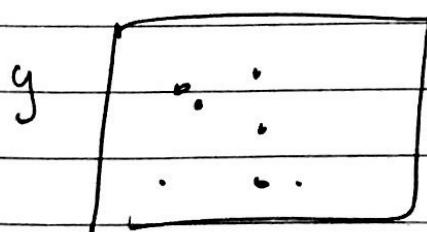
Both Numeric

(Quant)

(2 cat)

✓ Relationship b/w x, y

↳ Categorical var.



MBox  
plot

County  
(cat)

from this page :- Next Lecture;

- ✓ pie charts — Not suited for ordinal vars.
- ✓ for Nominal Variables
- ✓ left out depts / zones / . . . etc is not a good representation.

" pie gen shud represent all Values"

2 ways for rep Quant Variables:-

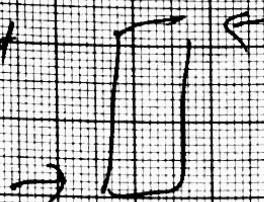
Box plot, Histogram

\* might be interested in average, Variance, Summary statistics

Box plot  $\rightarrow$  Captures Central Tendency

"line within the box" — Median of dataset

+ two boxes of dataset



" data variability

hence, upper quartile

25%

75% percentile.

✓ Histograms - pick of forms for Numerical data

? How many data points within a range  
is answered by Histogram.

✓ Columns → Rows

Histograms are first step towards  
→ "Empirical Construction"

Graphical Repn for Multiple Variables:-

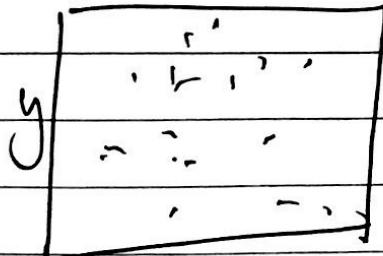
Scatter | Box plots | Contingency tables.

Sp:- Rep 2 Quantitative Variables

Bp:- One Cat v/s One Quant. Variable

CT:- 2 Cat variables with freq of occurrence as  
1 term.

Sp:- good to capture relationship b/w variables



$x \uparrow \rightarrow y \uparrow$ .

\* Bad @ Understanding each  
individual Variable.

Extended Box plots:- Country v/s Crime Rate  
(Cat) (Quant)

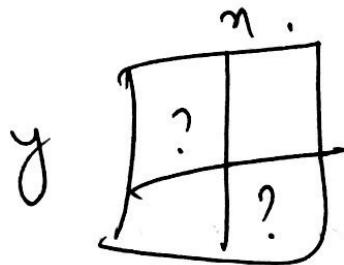
✓ Multiple Box plots on the same graph  
"Country make understanding of crime Rate"

⇒ Contingency

Contingency table :- 2 cat variab.

$$y < \frac{MBA - Y}{MBA - N}$$

$n <$  worked as n<sup>r</sup> no.



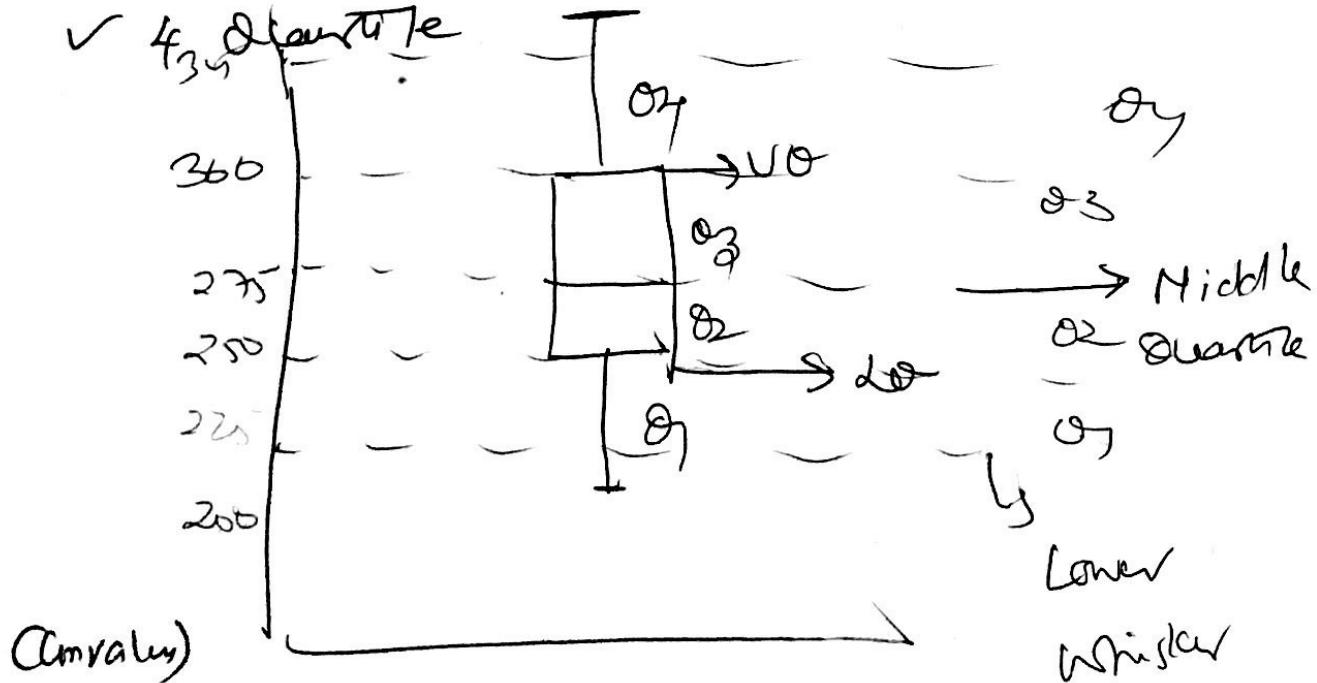
Box plot :- distributional characteristic of group of scores, levels of scores

✓ scores are sorted

✓ for equal sized groups are made from ordered scores.

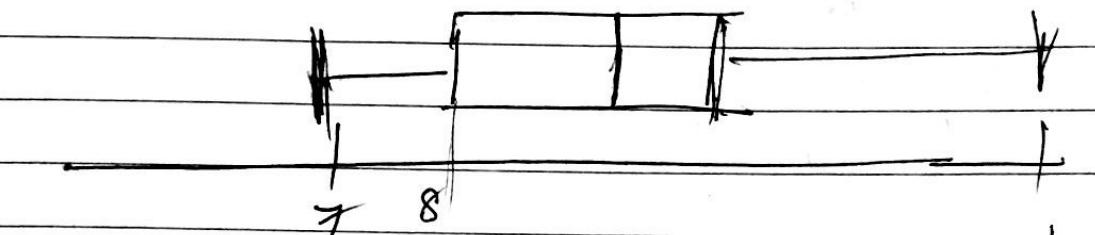
⇒ 25% scores in each group

✓ 4 Quartile



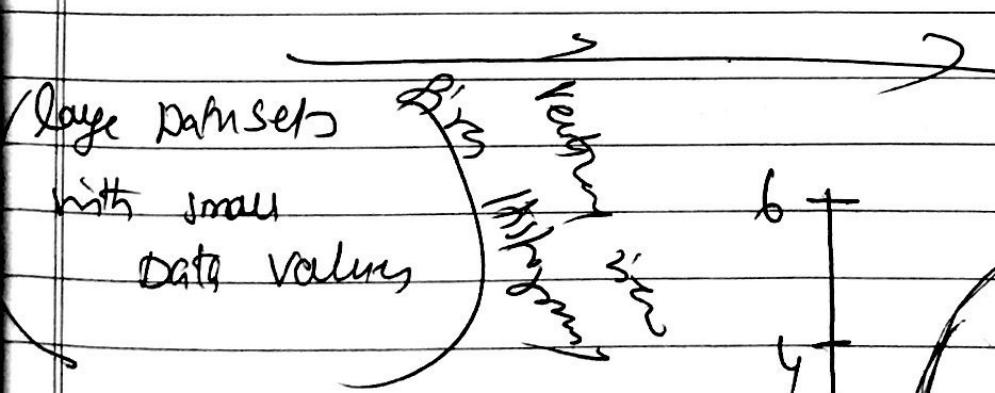
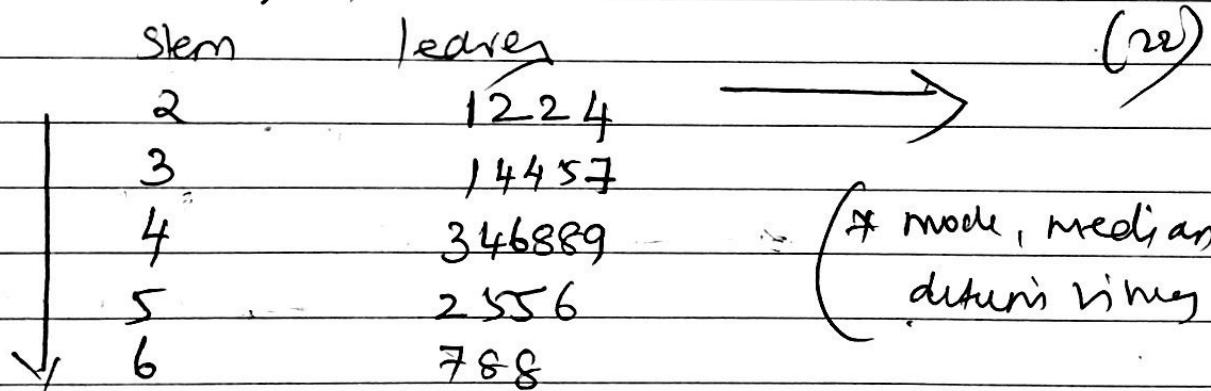
Hist → | Bar chart → Ordinal/Nominal also

7 8 8 8 9 9 | 9 9 10 10 13

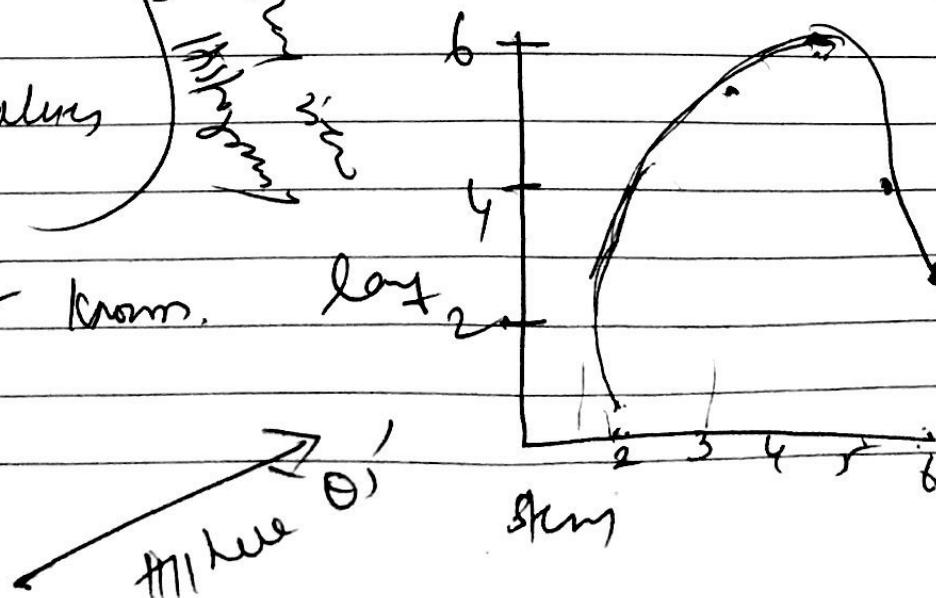


→ Stem + Leaf Plot:

21, 22, 23, 24, 31, 34, 34, 35, 37, 43, | 44, 46, 46, 47, 49, 53  
55, 55, 56, 67, 68, 68



Range / spread - known.



⇒ *from here on* / Lecture 04.

Descriptive Statistics — Summary, Central Tendency

✓ Summarizing data thru numbers.

→ Measures of Central tendency

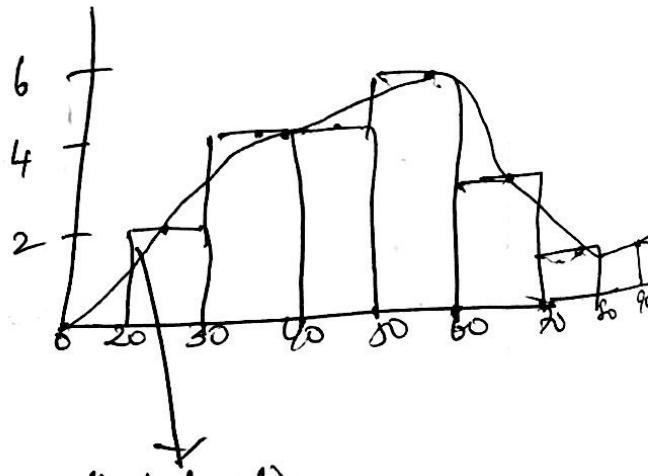
→ Dispersion

→ Skew and Kurtosis

✓ Histogram — rich representation; no of values within any range.

✓ Shows if frequency counts

20-30	→	2		3
4		4		1
4		5		0
5				1



✓ Simpler than distribution?

✓ Central tendency

@ Centre → min → max — in b/w

? pt in histogram → convergence to 50% of area

? ip central

✓ Balancing act? → x axis is balance  
and bays are weights. ? forces for  
balancing.

✓ Measure of Central tendency - ? Central Value.

Measures of dispersion — ? data dispersed  
around central value. far / close from centre

→ Skew / Kurtosis

↳ shape of the dispersion distribution.  
"dishes leans one side / other"

Kurtosis → ? fat are tails of distribution.

Mean = Balance in Seesaw (Central tendency).

Median — Another central tendency measure.

↓ sorted order — middle / value.

Mode → most freq. occurring

"common range"

? to use Mean / Median / mode.



outlier — outside majority /

| ↴ reason — errr data — Bad outlier

↳ imp part (slng) — Good outlier

Mean is impacted by outliers more.

Median ~~isnt~~ for outliers

Bad outliers  $\rightarrow$  Median is preferred

8, ~~80~~

Median remains same

✓ Mean and nearly day salaries  
 $\rightarrow$  outliers.

99.1 days - (or 1 mup) day.

1.1. "  $\rightarrow$  10 cmes.

Stock market trading  $\rightarrow$  normal bump

1000 days

Not contributing to sum  $\rightarrow$  Bad outlier

Instead of 87  $\rightarrow$  800

3, 4, 3, 1, 2, 3, 1, 1, 3, 6, 7, 4, 8, 00

Old Mean = 4.583  $\rightarrow$

Median :- 1, 2, 3, 3, 3, (4, 4), 5, 6, 7, 8, 9, 9

8th median unchanged

✓ Good outlier  $\rightarrow$  part (shy) to tell

lose 1 as every day on 99% days

108 (loss) 1. "outlier"

10,000 days :-

Median = -1, -1, -1, ..., 10,000.

Median = -1

Mean =  $\frac{-1 + -1 + \dots + 10000}{1000} \Rightarrow$  +ve no.

Mean is a better measure here

Mode :- Most popular / frequent value  
for distributions that are symmetric

↳ can be used for Nominal variables

✓ useful in multimodal distributions



many peaks to distribution

↗ Cont. prob distn. with two or more modes

0985

## Khan Academy

Box plot for odd

Q1      6, 15, 19, [21, 24, (38), 47, 55,] 58, [55, 78]

↑  
Q1

n=11 points  
↓

Q3 (upper quartile)

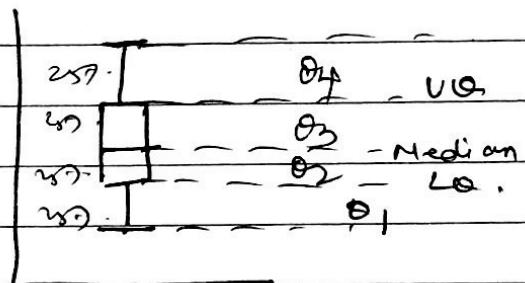
Q2 (middle quartile)

Generic Box:

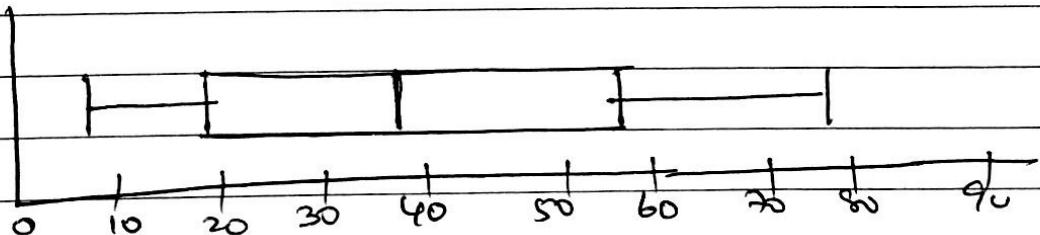
$$IQR = Q_3 - Q_1$$

$$\text{Outlier} < Q_1 - 1.5 IQR$$

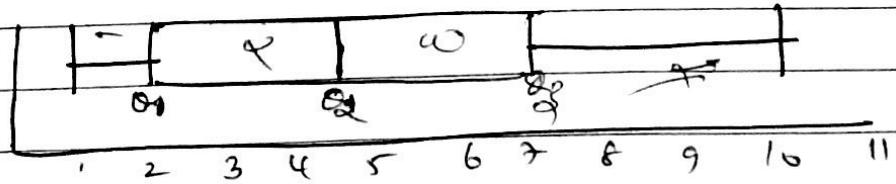
$$> Q_3 + 1.5 IQR$$



Q1



even: [1, 2, 2, 2], 3, 3, (4), (5), 5, 6, (7), 8, 8, 10      (14 pts)  
Median = 4.5 (Q2)



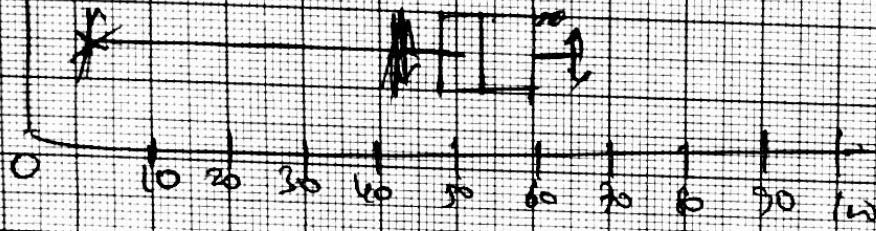
whisker for outer

$$\textcircled{2} \quad 2, 43, [49] 50, 51, [51] 53, 54, [60] 62, 63$$

$\Theta_1$        $\Theta_2$        $\Theta_3$

(11)

(11)



$$\Theta_1 - 1.5 \text{ IQR} = 49 - 1.5 \times 11$$

$$49 - 16.5$$

(32.5)

15  
12  
5  
2

$\alpha^2$  is an outlier

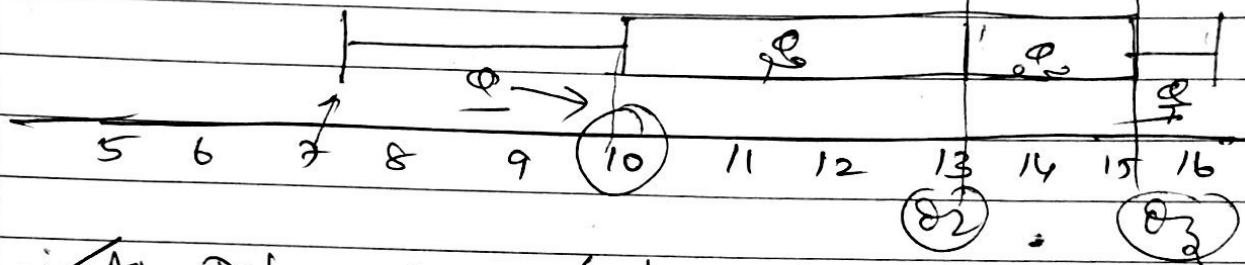
Lower whisker =

$$\Theta_3 + 1.5 \text{ IQR} = 60 + 16.5 = (76.5)$$

Upper whisker is fix

Interpretation:

(4)



(i) All Data are  $<$  than 17

(ii)  $\geq 10$  year old

e.g.  $\exists [10 \dots 13 \dots 15 \dots 16]$  (i) (7 data)

$6/7 \in \geq 10 \Rightarrow > 75\%$

e.g.  $\frac{7}{9} (11, 12, 14, 15, 15-16)$  (ii)

$$6/8 = 3/4 \geq 75\%$$

(iii) Only 1 person is  $\geq 17$  or 18 years old.

"can't say" | multiple  $\exists$  /  $\forall$   
can be there as well.

(iv) exactly 50% are  $\geq 13$ .

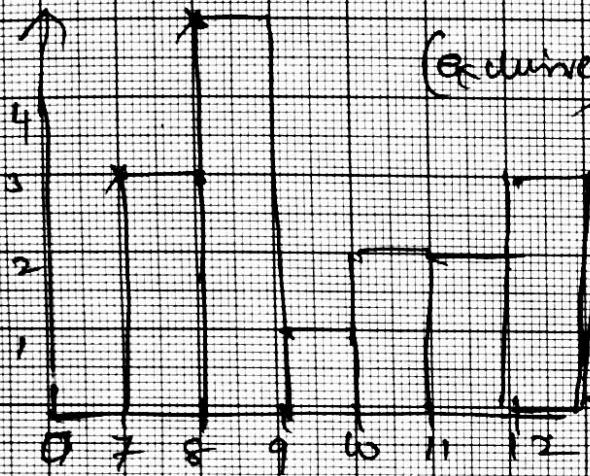
"can't say"  $\rightarrow$  True w/ (ii) but  
not (i).

Histogram:-

7, 7, 7, 8, 8, 8, 8, 9,

10, 10, 11, 11, 12, 12

$$12 - 7 = 5 / 5 = 1$$



Khan.

8 of students.

5, 7, 5, 9, 3, 7, 6, 9, 9, 9, 10,

12, 12, 7

11

Q) freq. table.

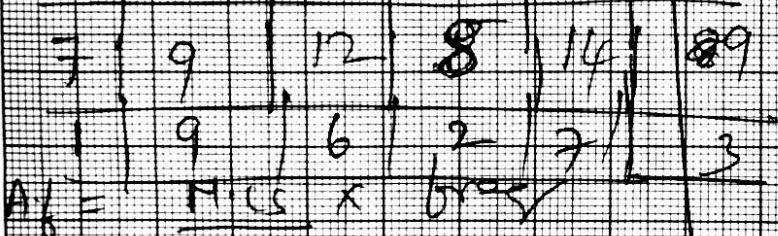
Age	5	6	7	8	9	10	11	12
#	2	1	3	4	1	2	0	2

4

3

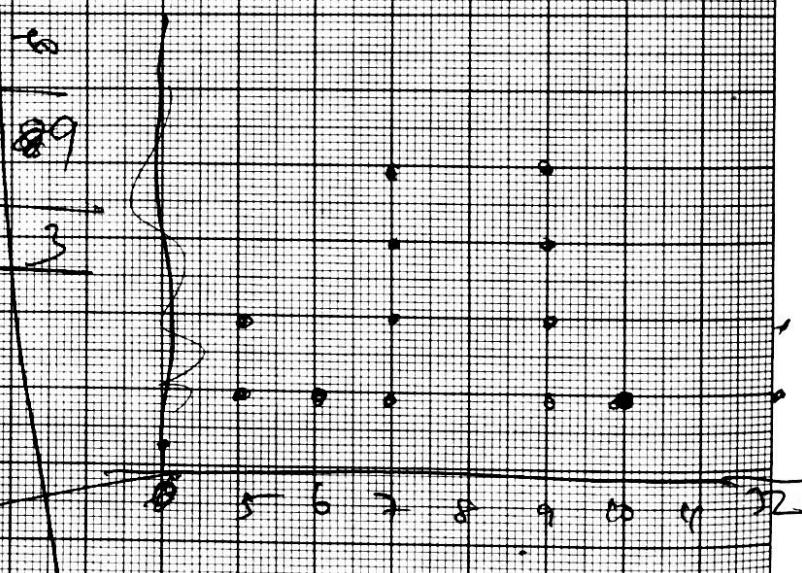
① Dot plot.

2) 0 → 5 | 5 → 10 | 10 → 15 | 15 → 20 | 20 → 25 | 25 → 30 | 30 → 35 | 35 → 40



A/C

$$= \frac{5}{5} \times 2 = 2 \rightarrow 10$$



how many  
> 9 (3)  
(i, ii, iii)

✓ freq of x larger shock

✓ range of age?

- max - min

4 = 17

look @ data

5/50

App 1  
 1, 3, 27, 33, 5, 63, 26, 25, 28, 16, 4, 40  
 29, 19, 22, 51, 58, 9, 42, 6, 4

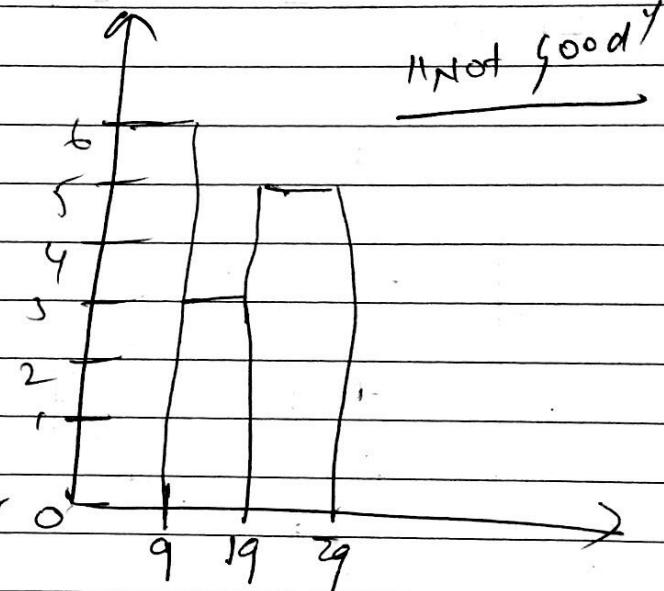
(Inclusive)  
 Bucket

"Bucketing" Matrix

Count (#)

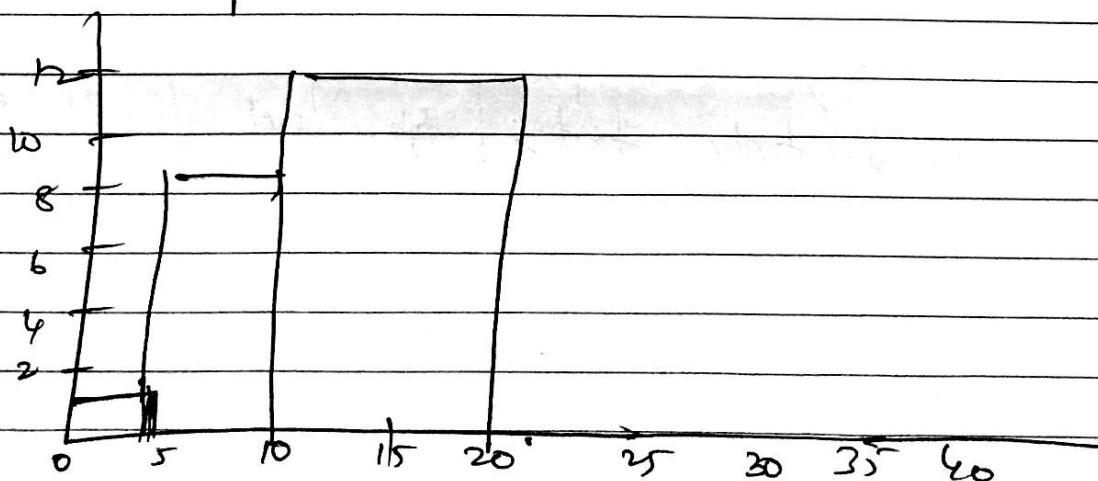
Visualizing

0 - 9	6
10 - 19	3
20 - 29	5
30 - 39	1
40 - 49	2
50 - 59	2
60 - 69	1



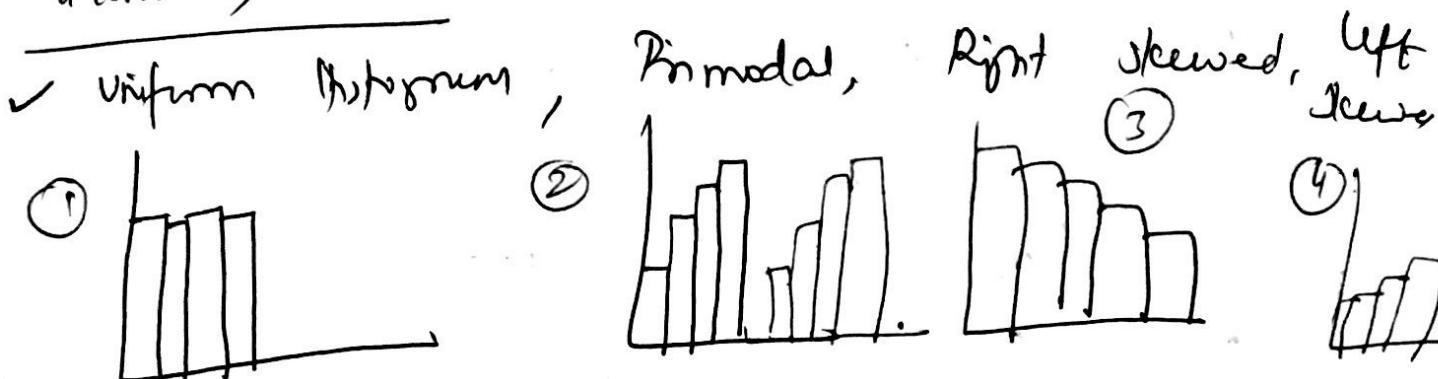
"Not good"

⑦

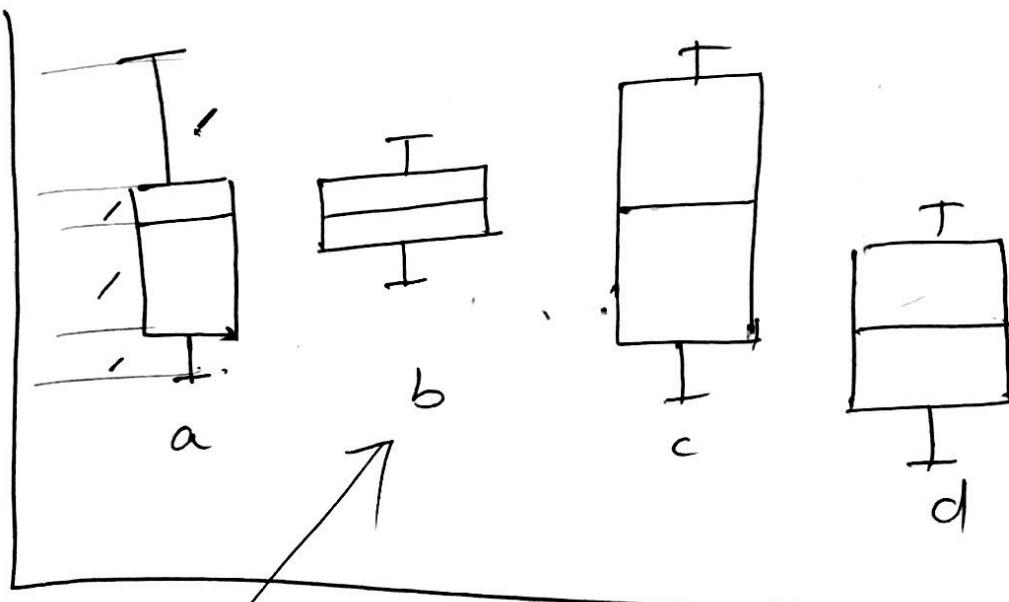


Make it larger cyl.

✓ Histograms vs Bar chart (Categorical Var on X axis)  
(Numerical)



→ Box plot Interpretations



✓ (b) → Comparatively short; High level of agreement.  
- Spread out factor is small

(a,c) → Comp. full; Varied views

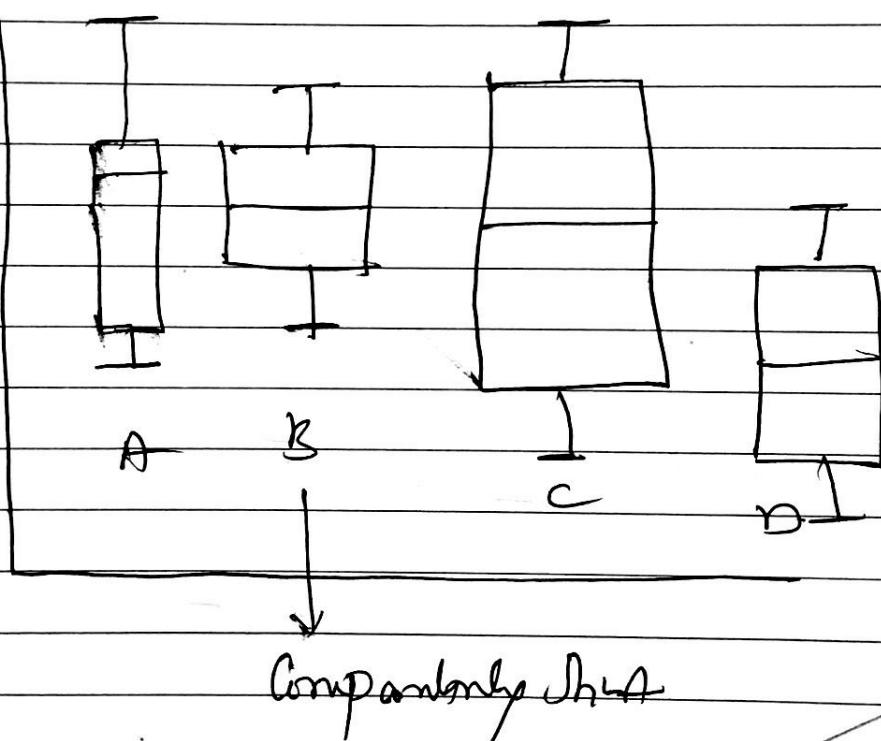
(a,d) → one box plot is higher/lower than other

4 sectors are varied  $\rightarrow$  (a)

$\nearrow$  longer upper - winter  $\rightarrow$  maximum variation  
in that quarter

shorter winter  $\rightarrow$  similarity  $\uparrow$

(\*) Same Median, Different Distributions (a, b, c).  
"Normal significant"



120.

80.