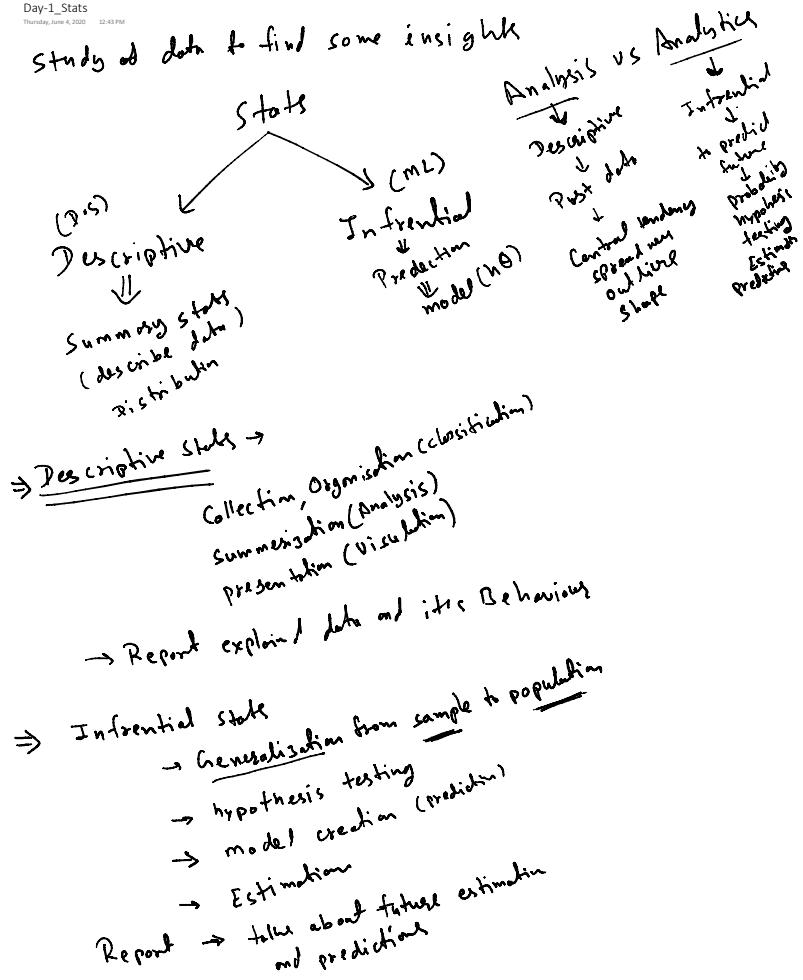
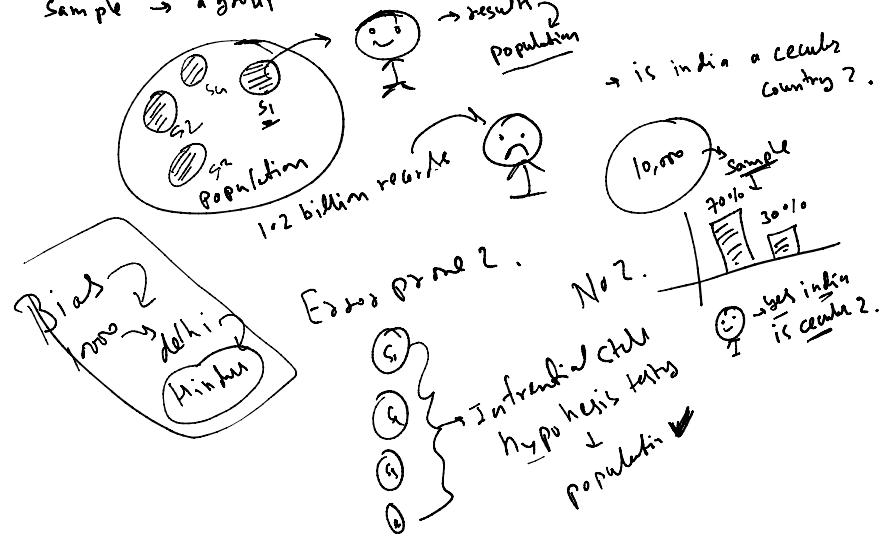


Study of data to find some insights



Population → Entire Data Set

Sample → a group or example from data



① Data type

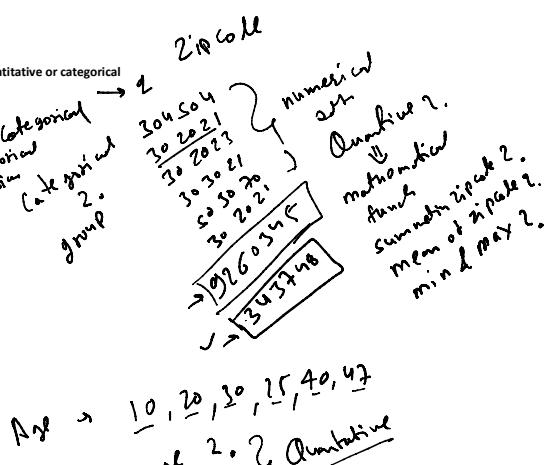
i) Quantitative Data (Numerical)
ii) Qualitative Data (Categorical)

i) Quantitative Data

ii) Qualitative Data (Categorical)

- Quiz:
for each variable below, identify each as either quantitative or categorical
- Zip Code → Categorical
 - Age → Quantitative
 - Income → Quantitative
 - Marital Status (Single, Married, Divorced etc.) → Categorical
 - Height → Quantitative
 - Letter Grades (A+, A, A-, B+, B, B-) → Categorical
 - Travel Distance to Work → Quantitative
 - Ratings on a Survey (Poor, OK, Great) → Categorical
 - Temperature → Quantitative
 - Average Speed → Quantitative

Age group categorical
child - 0 - 10
tween - 10 - 20
middle - 20 - 30
adult - 30 - 40
versus no adult



Age → 10, 20, 30, 40, 47
avg age 2. Quantitative
min age 2. Quantitative
max age 2.

at ↓
at ↓
order max and 2.
→ order

we can not order them

- O Nominal
Letter Grades (A,B+,B,B-)
Types of Fruit(Apple,Banana,etc.)
Ratings on a Survey(Poor,OK,Great)
Types of Dog Breeds(German Shepherd,Collie,etc.)
Genres of Movies(Horror,Comedy,etc.)
Gender
Nationality
Education(HS,Associates,Bachelors,Masters,PhD,etc.)

Continuous vs Discrete

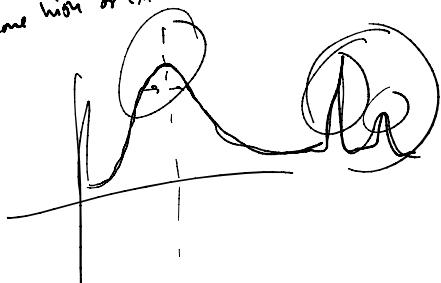
- Continuous
Travel Distance from Home to Work → 12m, meters, centi, milli ...
Number of Pages in Book → count
Amount of Rain in Year → cm, millimeter, mm
Time to Run a Mile → min, sec, sec
Number of Movies Watched in a Week → countable
Amount of Water Consumed in a Day → ml, l
Number of Phones per Household → countable

{160, 100, 110, 105, 160, 165, 140, 137}

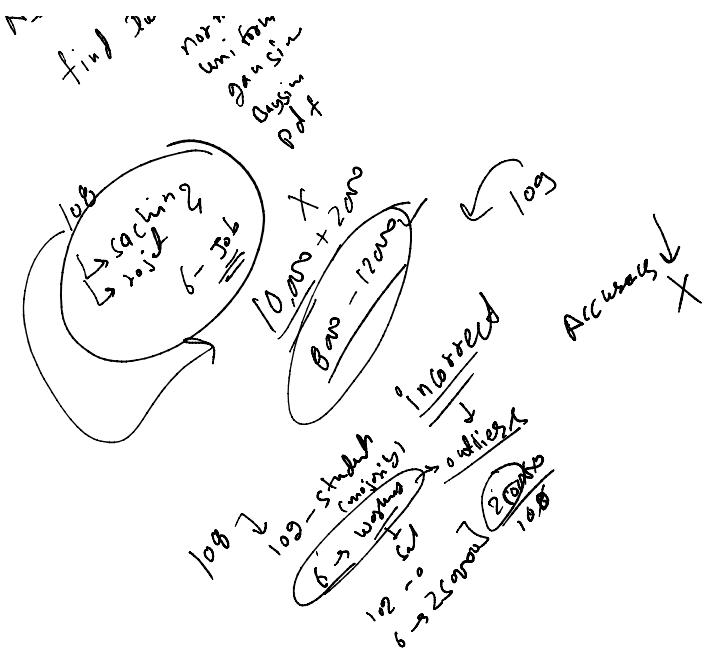
140

Quantitative Data Analysis

- Measure of Center → Single value that describes the center point of your data set. It's the most descriptive property of data.
- Measure of Spread → Represents variability in data. It generalizes your central tendency.
- Shape of Data → Distribution of data, properties of data.
- Outlier → extreme high or extreme low values



Assignment
find out distributions?
normal
uniform
gaussian
exponential



① Numpy vector & matrices

↓
vector & matrices
(array)

Data can be stored in →

- ① Scalers
- ② Vectors
- ③ matrix
- ④ Tensors

① Scalar is a single discrete observation
only magnitude

5
10
35.5

Each discrete observation is a scalar type

② row vector or column vector
1D array

Vector → collection of scalars

$$v_1 = \begin{bmatrix} 10 \\ 20 \\ 30 \end{bmatrix} \quad \text{column vector}$$

... for

$$v_1 = \begin{bmatrix} 10 \\ 20 \\ 30 \end{bmatrix} \text{ column vector}$$

$$v_2 = \begin{bmatrix} 50 & 100 & 200 \end{bmatrix} \text{ row vector}$$

③ Matrix → collection of vectors

$$\begin{matrix} p_1 \\ p_2 \\ p_3 \end{matrix} \rightarrow \begin{bmatrix} \text{qnt} \\ 3 \\ 5 \\ 6 \end{bmatrix} \quad \begin{matrix} \text{price} \\ 200 \\ 150 \\ 300 \end{matrix}$$

$$\text{Shape} \rightarrow \text{row} \times \text{coll} \Rightarrow (3, 2)$$

$$\text{Size} = \text{row} \times \text{coll} \Rightarrow 6$$

row → records, tuple
cols → features, attributes

Dimension of matrix - 2D

$$\begin{matrix} n \\ \downarrow \\ \text{y} \end{matrix} \quad 3 \times 2$$

Square matrix

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \quad 3 \times 3$$

Diagonal diagonal elements [1 2 3 3]
off-diagonal elements [3 2 3 3]

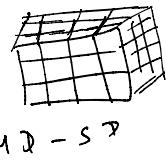
$$I = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \text{eye, Identity matrix}$$

$$\begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \quad 3 \times 2 \text{ zero matrix}$$

$$\begin{bmatrix} 1 & 1 & 1 \end{bmatrix}_{2 \times 3} \text{ ans matrix}$$

\rightarrow Tensors \rightarrow high dimensional arrays



$4^2 - 5^2$

\rightarrow Operations \times on a matrix

① Transpose

$$A = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix}_{3 \times 2} \quad A^T = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{2 \times 3}$$

② Symmetric matrices

if $A = A^T$ then A is symmetric

$$A = \begin{bmatrix} 3 & 4 \\ 4 & 6 \end{bmatrix}_{2 \times 2}$$

$$A^T = \begin{bmatrix} 3 & 4 \\ 4 & 6 \end{bmatrix}_{2 \times 2}$$

$$\boxed{A = A^T}$$

③ Addition

element wise add

Cond: shape of both matrix should be same

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \quad B = \begin{bmatrix} 7 & 8 & 9 \\ 10 & 11 & 12 \end{bmatrix}$$

... 12 7

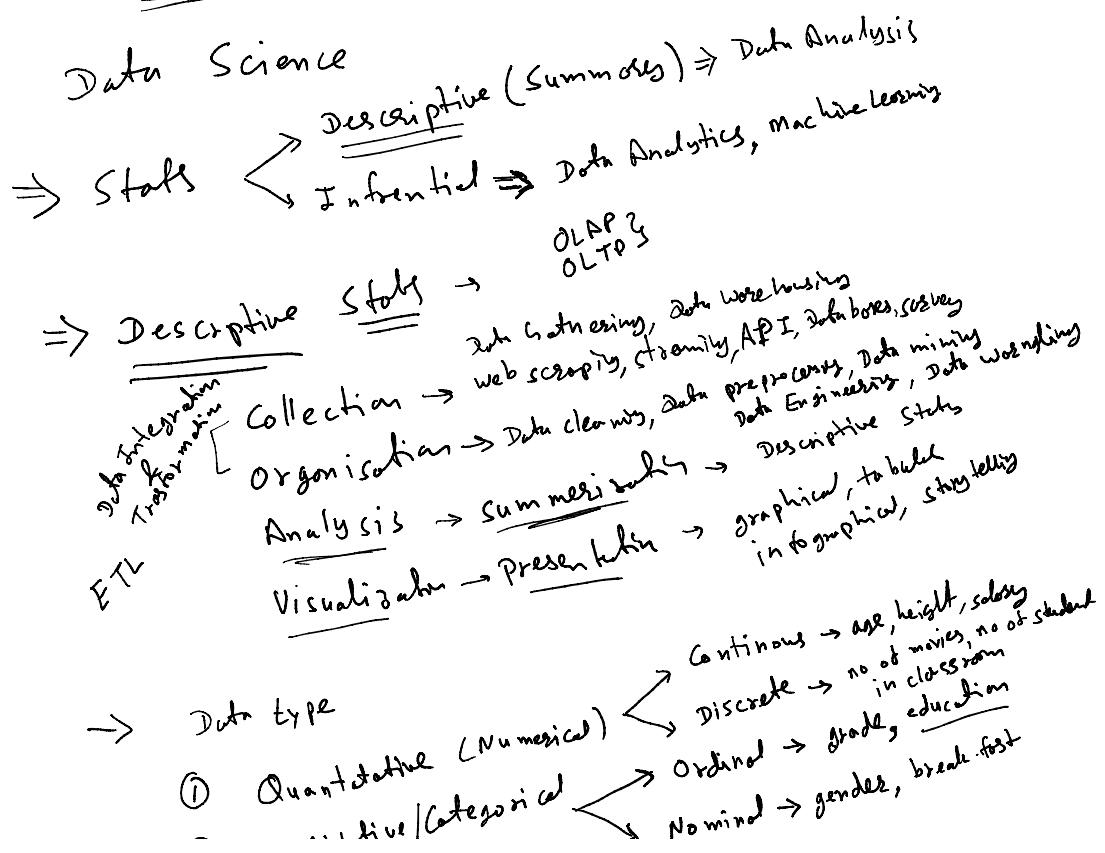
$$A = \begin{bmatrix} 4 & 5 & 0 \\ 1 & 10 & 12 \\ 14 & 16 & 18 \end{bmatrix}$$

$$B = \begin{bmatrix} B_{11} & B_{12} & B_{13} \\ B_{21} & B_{22} & B_{23} \end{bmatrix} \Rightarrow \begin{bmatrix} B_{11} & B_{12} & B_{13} \\ B_{21} & B_{22} & B_{23} \end{bmatrix}$$

$A + B$, commutative
 $A + B = B + A$



Review



- ① Quantitative (Numerical) → Ordinal → σ
 ② Qualitative/Categorical → Nominal → gender, break toes

→ Data Analysis

→ Qualitative (Numerical)

Sums of
Desribe S

- ① Central tendency → Best value that can represent the whole sum → mean, median, mode
- ② Measure of Center → continuous variable
- ③ Measure of spreadness → tells about variability of data → range, variance, standard deviation, quantiles, IQR, 5-number summary, deciles, percentiles
- ④ Shape of data → relation of data with standard normal distribution
- ⑤ Outliers → visualization of data distribution, skewness, correlation, dispersion

extreme values which does not follow characteristics of data

→ Central tendency

→ Let's take salary of 50 employes

$$x = 26, 71, 53, 96, 23, 93, 55, 87, 37, 38, 65, 75, 27, 46, 22, 23, 76, 87, 51, 32, 44, 77, 20, 97, 20, 21, 29, 63, 52, 31, 32, 35, 64, 45, 24, 84, 62, 71, 34, 60, 67, 21, 49, 32, 51, 20, 21, 77, 54, 43$$

So → ① Mean

$$\begin{matrix} n_1 = 38 \\ n_2 = 25 \\ \vdots \\ d_1 = 27 \end{matrix}$$

$$\bar{x} = \frac{n_1 + n_2 + n_3 + \dots + n_n}{N}$$

$X = [x_1, x_2, x_3, x_4, \dots, x_N]$
 $n = \text{no of observations, size}$
 $\bar{x} = \text{Sample Data, Random Variable}$

$$\bar{x} = \frac{\sum_{i=1}^N x_i}{N}$$

$$d = (\bar{x} - n_1) + (\bar{x} - n_2) + \dots + (\bar{x} - n_n)$$

$$\therefore (n_1 + n_2 + \dots + n_n)$$

$\sum_{i=0}^{N-1} d_i = \bar{x} - \bar{x} = 0$
 if \bar{x} is center then $d = 0$.

data distributions
 normal
 binomial
 gaussian

$$d = (\bar{x} - n_1) + (\bar{x} - n_2) + \dots + (\bar{x} - n_N)$$

$$0 = n\bar{x} - (n_1 + n_2 + \dots + n_N)$$

$$0 = n\bar{x} - \sum_{i=1}^N n_i$$

$$\bar{n} = \frac{\sum_{i=1}^N n_i}{N}$$

$$\bar{x} = \frac{\sum_{i=1}^N n_i}{N}$$

$$d = 0$$

from \bar{x}

$$A = 1, 2, 3, 4, 5$$

$$\bar{x}_B = 2 \times \begin{cases} d = 1+0-1-2-3 \\ d = -5 \end{cases}$$

$$\bar{x} = x_A + d \Rightarrow |2-5| \Rightarrow 3$$

$$\bar{x} = \frac{1+9+3+4+5}{5}$$

$$\sum d = 0$$

$-19 - 18 - 17 - 16 - 15 + 85$
 $-85 + 85$
0

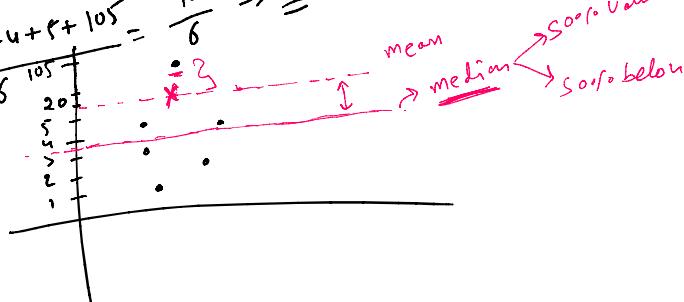
mean can be shifted by outlier

$$x = 1, 2, 3, 4, 5, 105$$

$$\bar{x} = \frac{1+2+3+4+5+105}{6} = \frac{120}{6} \Rightarrow 20$$

$$\bar{x} = 20$$

average of values



→ median ⇒ it always bisect data into two half

average at position

median will be $X[i]$, where $i = \frac{N+1}{2}$

① Odd observation ⇒ median will be $X[i]$, where $i = \frac{N+1}{2}$

$$x = 5, 2, 4, 3, 1$$

mean = 3 → uniform Normal

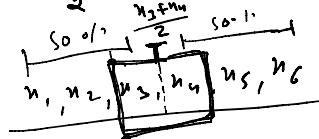
Step ① sort data (asc)
 $x = 1, 2, 3, 4, 5$

$$\text{median} = \frac{(N+1)}{2} \text{ value}$$

$N=5$ median = 3

$\frac{5+1}{2} \rightarrow 3$ rd value of x will be center

② Even observation



median will be average of i^{th} and $i+1^{\text{th}}$ of sorted X data set

$$\text{where } i = \frac{N}{2}$$

$$\text{mean } \bar{x} = \frac{\sum x}{n}$$

$$X = 1, 2, 3, 4, 5, 105$$

median

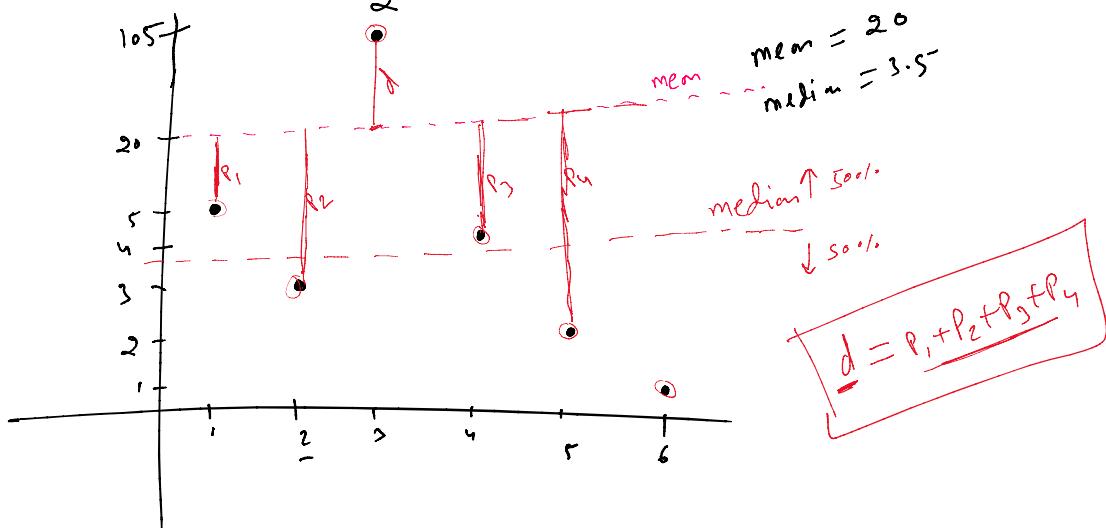
$$i = \frac{6}{2} = 3$$

$$i = 3,$$

$$\underline{\text{mean}} = 20$$

$$\underline{\text{median}} = 3.5$$

$$x = [1, 3, 105, 4, 2, 1] \quad \left[\frac{n_i + n_{i+1}}{2} \right] \Rightarrow \left[\frac{3+4}{2} \right] \Rightarrow 3.5$$



\Rightarrow Discrete Data points

mode \rightarrow highest frequency

Case 1. single mode

$$1, 3, 1, 3, 2, 1, 3, 2, 1, 1, 5, 1, 8, 3$$

1	6
2	2
3	4
5	1
8	1

mode = 6

Case 2. multiple modes

$$3, 5, 2, 1, 3, 5, 2, 3, 5, 1, 4, 5, 3$$

X	f
1	2
2	2

mode = 3, 5

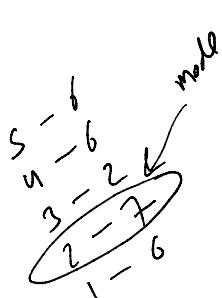
1	2
2	2
3	4
4	1
5	4

mode = 3, 5

Case 3 :

No mode

2, 3, 1, 2, 4, 3, 1, 4



? No mode

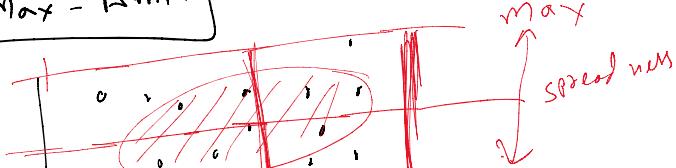
② Measure of Spread

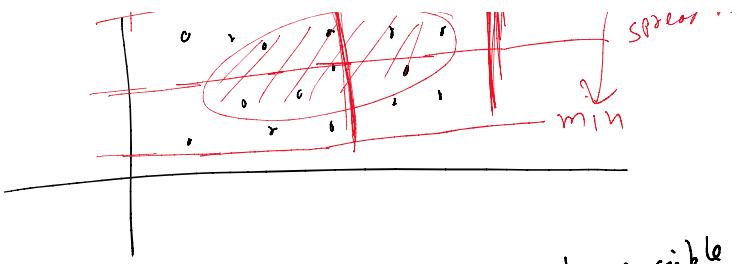
- ① Range
- ② Variance
- ③ Standard deviation
- ④ Quantiles
- ⑤ Inter quartile range
- ⑥ 5-point summary
- ⑦ Deciles
- ⑧ Percentiles

Distribution of Data Around Center

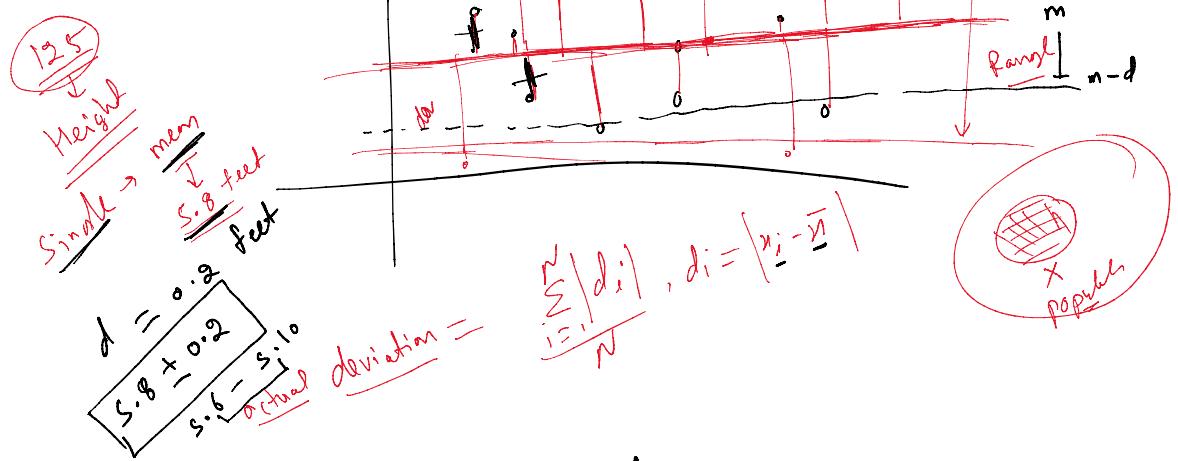
- ① Range \rightarrow Variance of entire data set

$$R = A_{\text{Max}} - A_{\text{Min}}$$





② Variance → Minimum Variability possible



Variance → squared deviation

$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N}$$

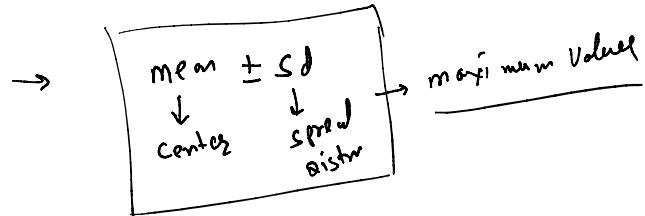


③ Standard deviation →

$$\sigma = \sqrt{V\sigma^2 (\sigma^2)}$$

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

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



\rightarrow

mean
5.6

sd
.4

\rightarrow

mean \pm sd
 $5.6 \pm .4$

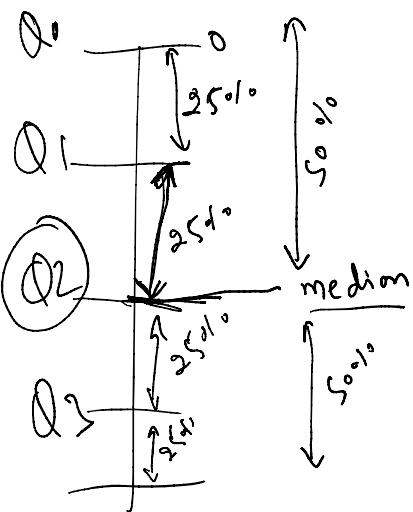
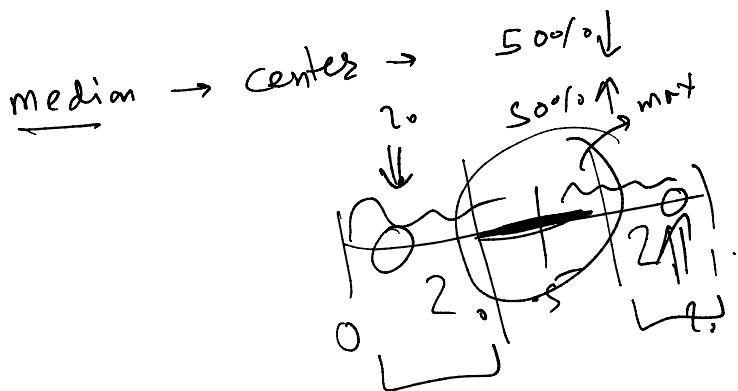
4.2

.2
.6
.10

6.5



④ Quantiles



$d_{1:n} = 4, 1, 3, 8, 6, 2, 5, 7, 9, 6, 10$

mean = 5.5

$d_{2:n} = 1, 7, 3, 4, 5, 6, 2, 8, 9, 10$

$Q_2 = (\text{median}) \rightarrow \boxed{\text{5th and 6th}} \text{ Avg}$

$\frac{5+6}{2} \Rightarrow 5.5$

median = 5.5

$$\overline{2}$$

median = 5.5

$Q_1 \rightarrow$ median of $0 - Q_2$ values

$Q_3 \rightarrow$ median of $Q_2 - Q_u$ values

$Q_1 = 1, 2, 3, 4, \rightarrow 2.5$

$, 6, 7, 8, 9, 10 \rightarrow 8$

$Q_3 =$

1 - 10

$\rightarrow \underline{\underline{5.5}} \rightarrow \text{center}, \underline{\underline{sd}} = 2.$

$$\frac{5.5 - 2.5}{3.0}$$

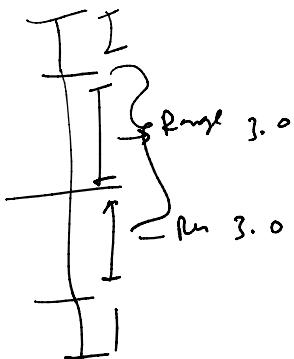
$\rightarrow 1ct 25\% \rightarrow \underline{\underline{2.5}}$

$$\underline{\underline{8.5 - 5.5}} \\ 3.0$$

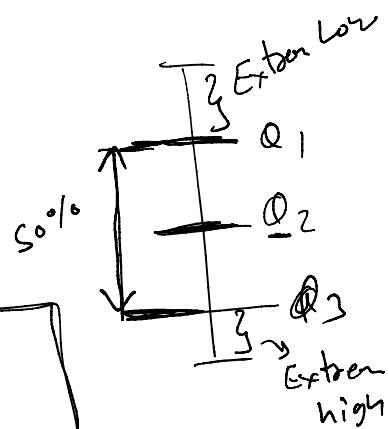
$\rightarrow 50\% \rightarrow \underline{\underline{5.5}}$

$$5.0$$

$\rightarrow 75\% \rightarrow 8$



$\rightarrow Q_1 - Q_3$



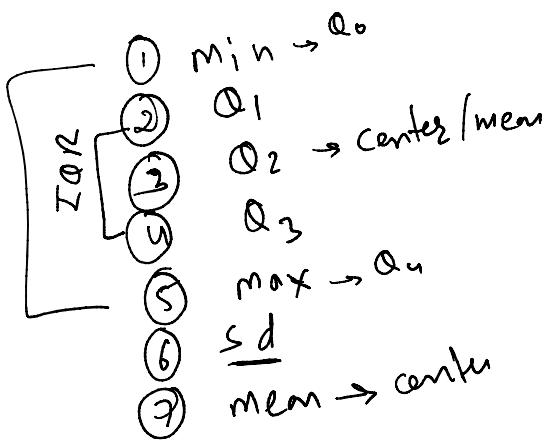
(S) IQR \Rightarrow $Q_3 - Q_1$

(S) 5-point summary

$\therefore \min \rightarrow Q_1$

$\min, Q_1, Q_2, Q_3, \max$
C minst summary

Exploratory Data Analysis
Range



$[min, Q_1, Q_2, Q_3, \dots]$
5 point summary

median - mean
significal differ

Outliers or Skewness

Data Distributions

Shape

Quartiles \rightarrow

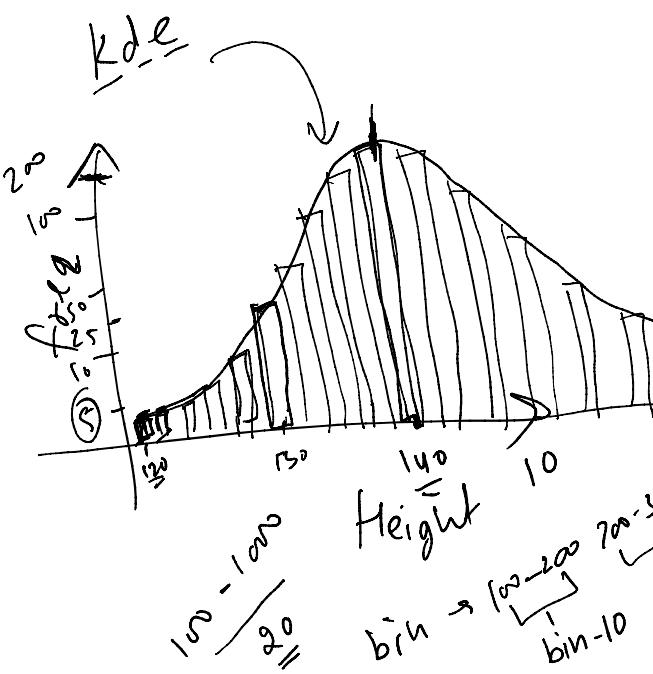
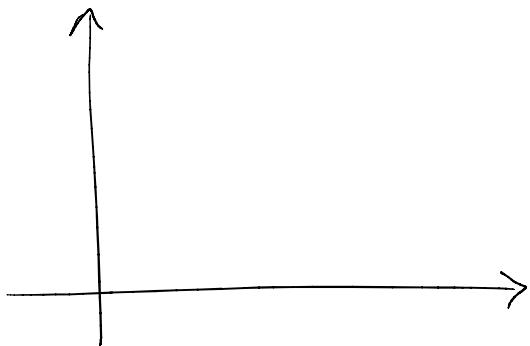
population mean =

$$\frac{\sum_{i=1}^N x_i}{N}$$

sample mean =

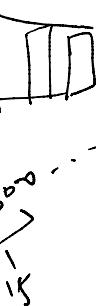
$$\frac{\sum_{i=1}^{N-1} x_i}{N-1}$$

(i) Histogram \rightarrow Data Distribution (numerical)



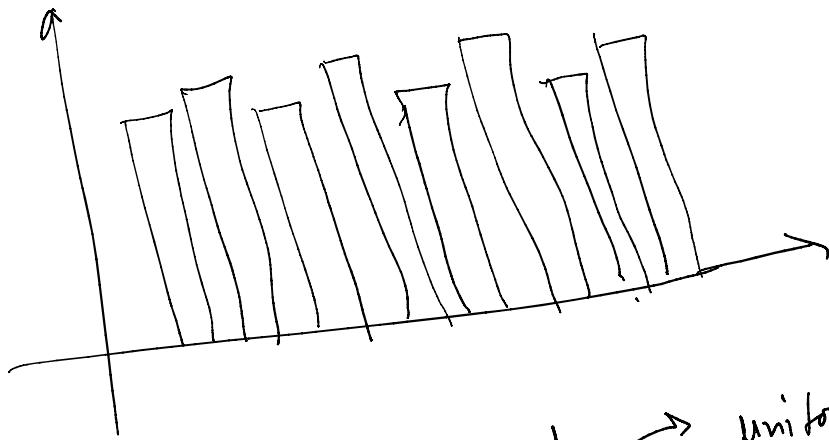
Probability \times distribution

... \rightarrow distribution



① uniform distribution

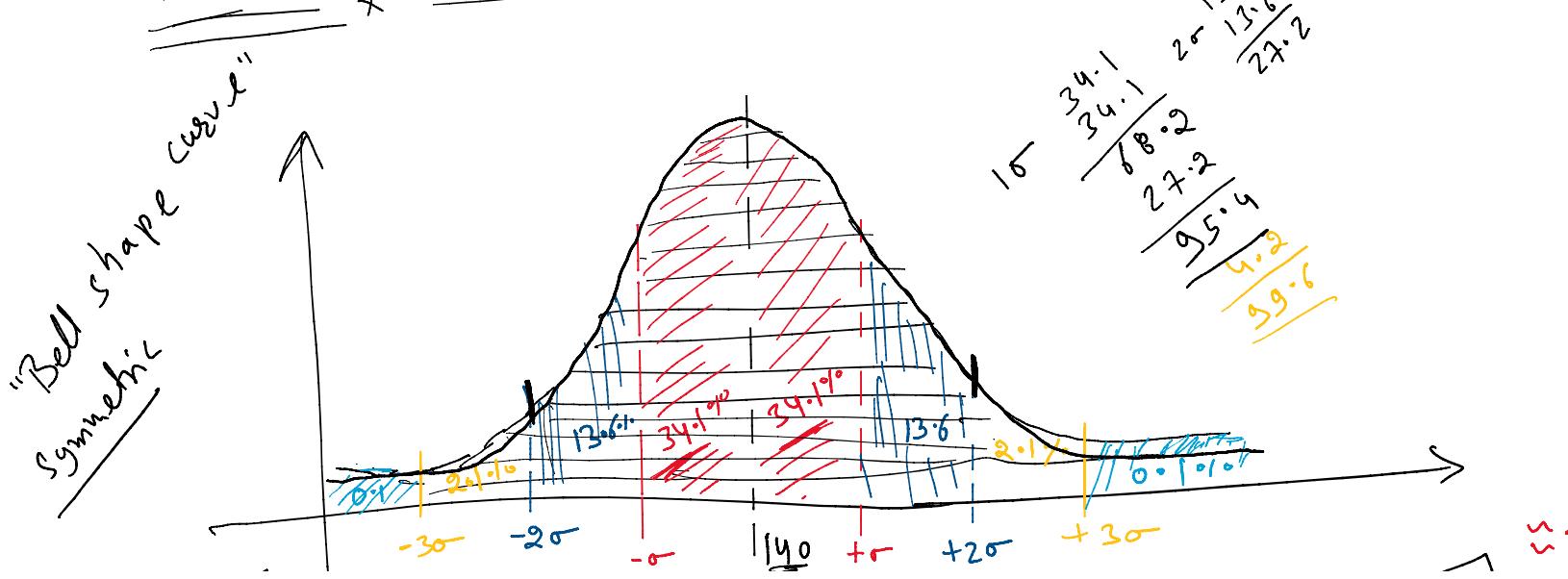
where probability of occurrence at any point in data set is equal.



1 - 50

$$f(n) = \frac{1}{b-a} \quad \text{random} \rightarrow \text{uniform}$$

Normal Distribution (Gaussian Distributed)

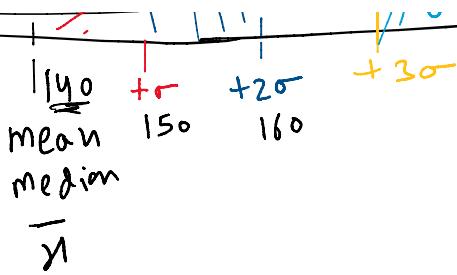


70%

$$\bar{x} = \mu$$

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

mean / median

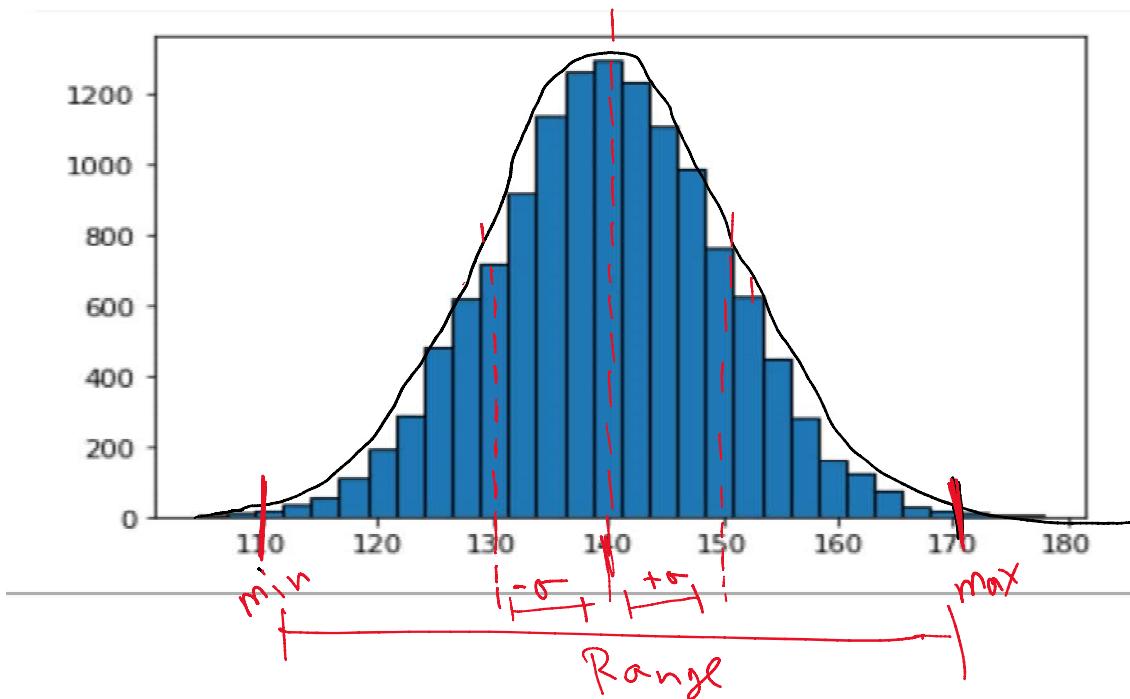


$$\bar{x} \pm 2\sigma \rightarrow 120, 160$$

$$Sd = \sigma$$

$$\bar{x} \pm \sigma$$

$$140 \pm 10$$



$$f(h)$$

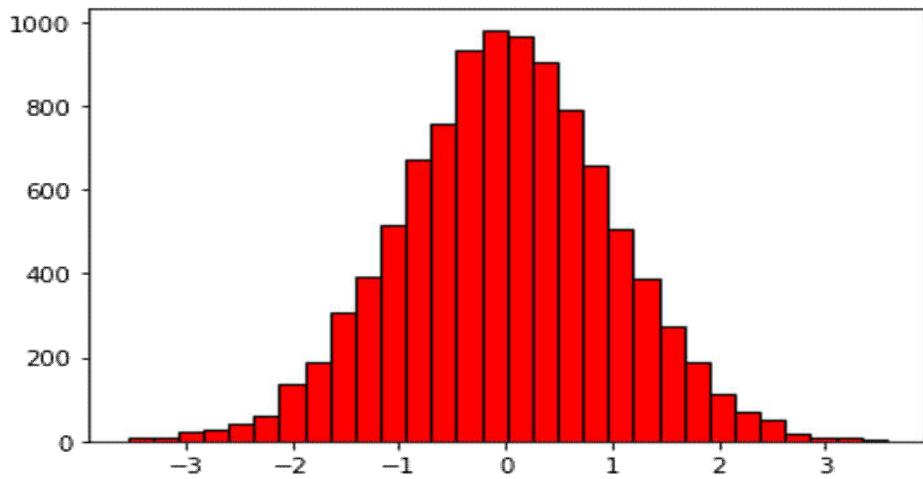
$f_{0.5}, 130-150$

→ Standard Normal Distribution

70%

0 - 15°
70%

)



mean = 0
std = 1

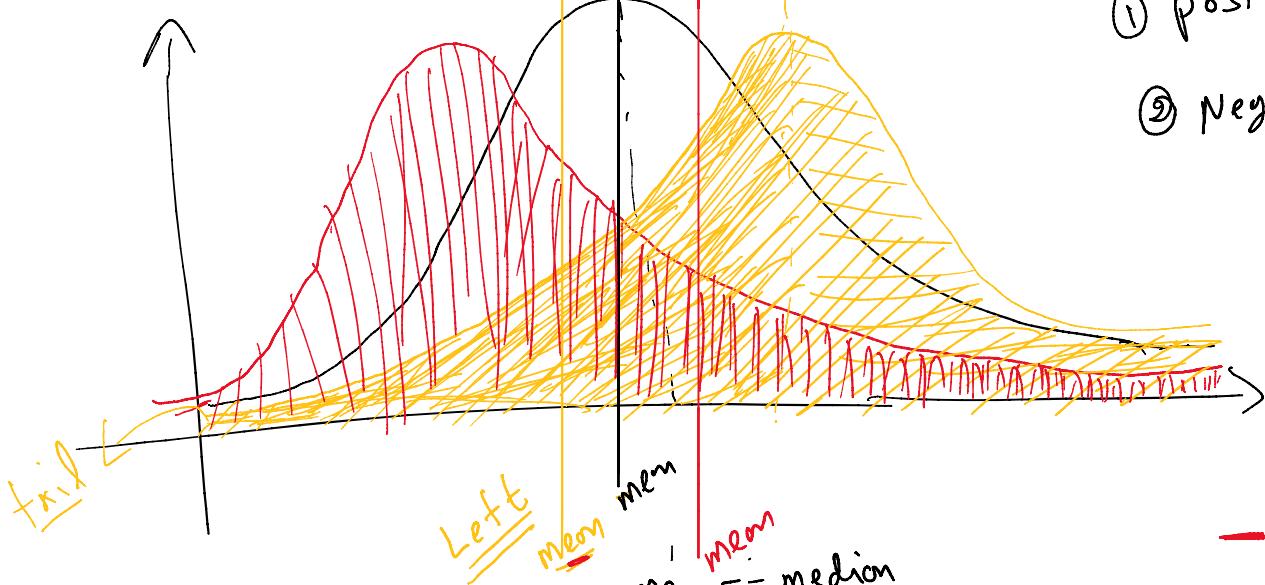
⇒ study of shape

① Symmetric or Asymmetric

Skewness

① positive skewness

② Negative skewness



① Positive skewness

mean > median

② Negative skewness (Left shift)

mean < median

③ Normal or symmetric

mean == median

- pos /
- neg /
- pos /
- neg /

g
em

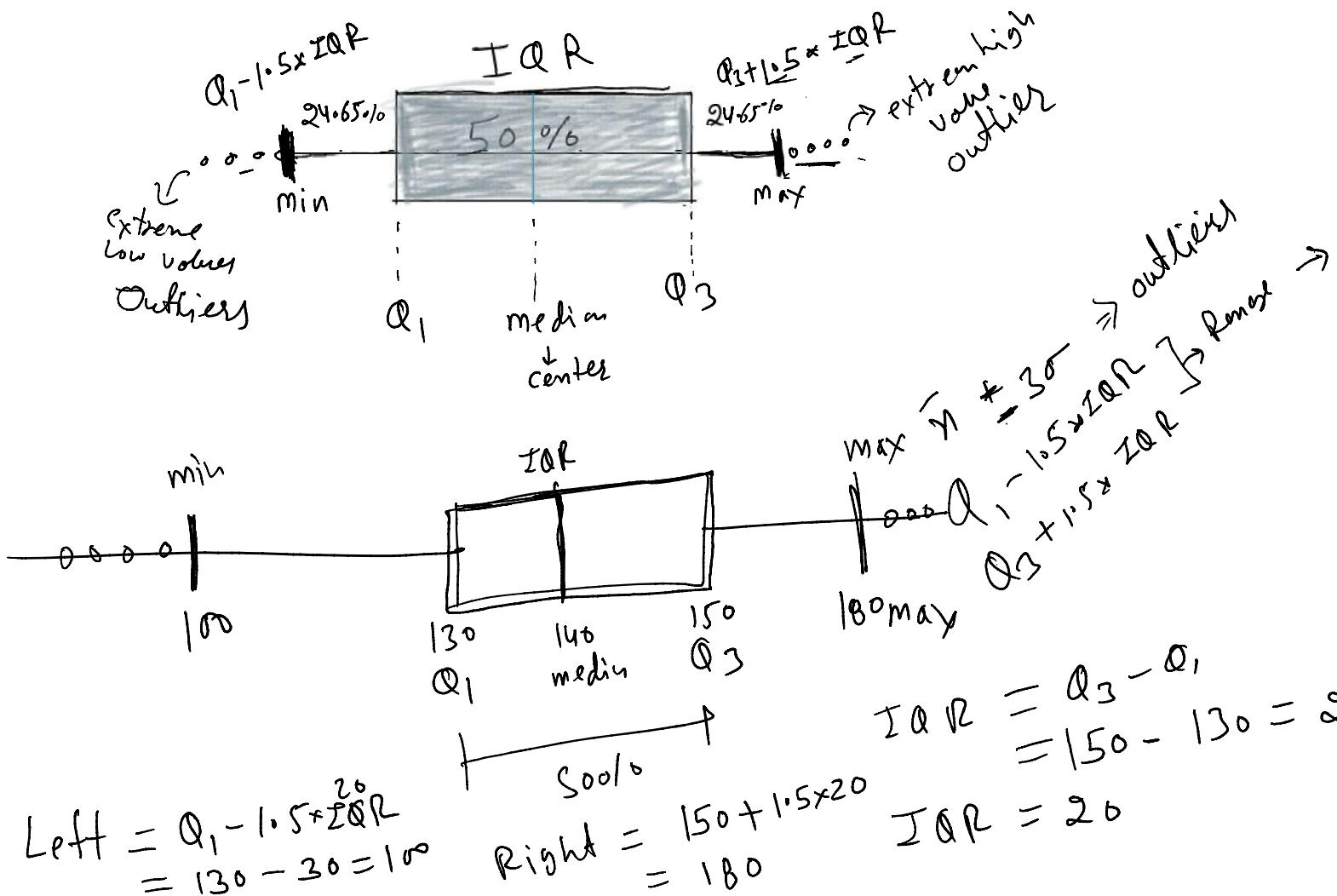
/fight
Left

right
left

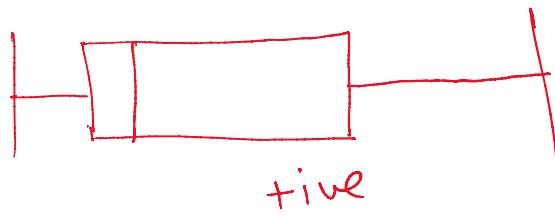
$\Rightarrow \underline{\text{Quartiles}} \rightarrow$

$$\begin{aligned}Q_1 &= 25\% \\Q_2 &= 50\% \\Q_3 &= 75\% \\Q_4 &= 100\%\end{aligned}$$

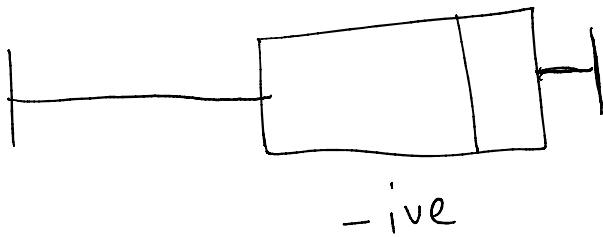
Box-plot



①



②



→ Summary stats

min

Q_1

Q_2

Q_3

max

std

mean

box plot

histogram

Box

