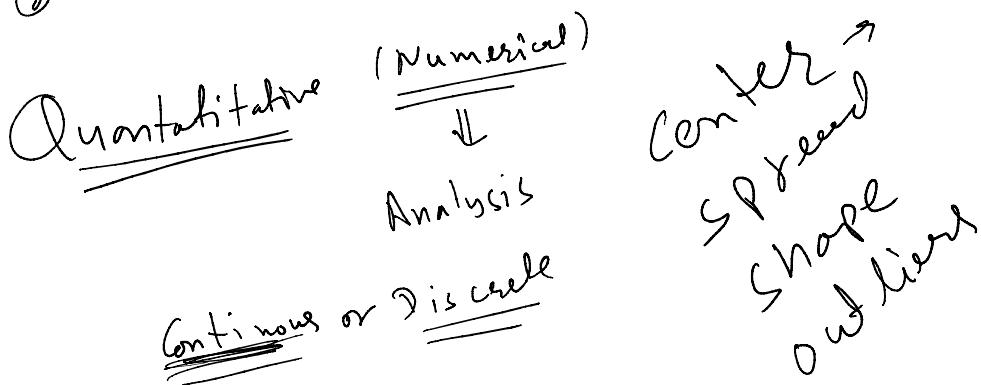


⇒ Data type

- ① Quantitative (Numerical)
- ② Qualitative (Categorical)



→ Analysis → summary

30 employ's salary

31, 86, 28, 77, 88, 34, 60, 32, 66, 52, 83, 64, 90, 37, 41, 44, 37, 69, 70, 34, 87, 45, 37, 80, 29, 83, 41, 37, 34, 45

① Central tendency or measure of center → mean, median & mode

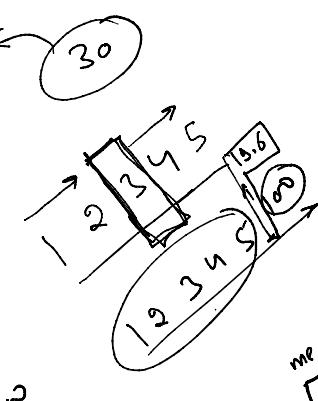
Continuous

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

30

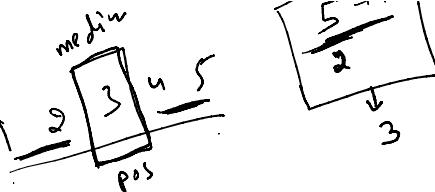
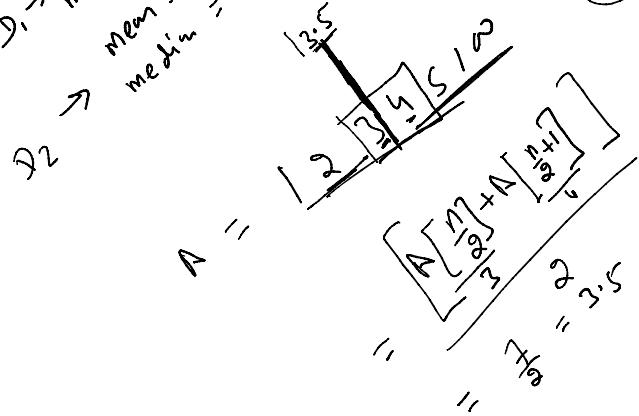
$$\text{mean} = \frac{115}{6} \approx 19.6$$

$$\begin{aligned} \text{mean} &= 3.5 \\ \text{median} &= 3.5 \\ \text{mode} &= 2.0 \\ \text{min} &= 1.0 \end{aligned}$$



$$\begin{aligned} \text{median} &= \frac{5+1}{2} \\ &= 3 \end{aligned}$$

\rightarrow mean
mean = 19
median = 3.5



\rightarrow Discrete



Q. I want to represent 1 billion

mean \rightarrow bias
median \checkmark

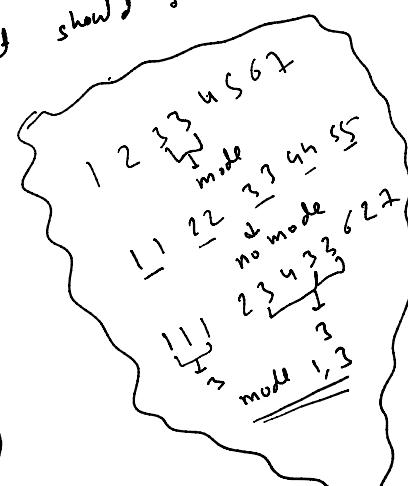
center of Indian people salary?

1. \rightarrow reason

outlier
Billionaires & outliers

Q. I want to summarize high ad mean in a single digit what should I use

mean \checkmark
median \checkmark



$$\text{mean} \quad \bar{x} = \frac{\sum x_i}{N}$$

$$\bar{x} = \frac{n_1 + n_2 + n_3 + n_4 + \dots + n_k}{N}$$

Median \rightarrow ① sort your values

(use) ② if even observation (N)

$$\text{cal } i_1 = \frac{N}{2}$$

$$i_2 = \frac{N}{2} + 1, i_1^+$$

$$\text{then median} = \frac{A_{i_1} + A_{i_1^+}}{2}$$

{ 100 }

$$A = \{1, 2, 3, 4, 5, 100\}$$

$N = 6$

$i_1 = \frac{N}{2} - 3$

$i_2 = i_1 + 1 = 4$

$\text{median} = \frac{A_{i_1} + A_{i_2}}{2}$

$= \frac{3 + 4}{2}$

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

then $i_2 = \frac{N}{2}$
 $\text{median} = \frac{A_i + A_{i+1}}{2}$
 A is array of size

if N is odd, here N is no. of observations

Ques ②

$$\text{eg. } 1, 2, 3, 4, 5$$

$N = 5$

$i = \frac{N+1}{2} = 3$

$A_i = A_3 \Rightarrow 3$

$\boxed{\text{median} = 3}$

$$i = \frac{N+1}{2}$$

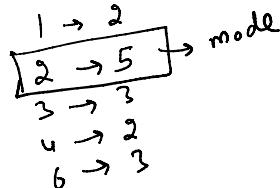
$$\text{median} = A_i \quad \text{where } A_i \text{ is } i^{\text{th}} \text{ element in data}$$

① Mode \rightarrow (Discrete)

single mode

$$\text{eg. } 1, 2, 3, 6, 2, 4, 3, 6, 2, 2, 4, 2, 6, 3, 1$$

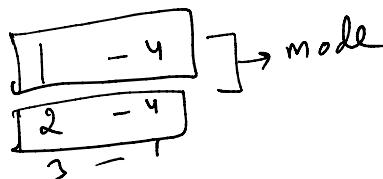
Ques ①



multiple mode

$$1, 2, 5, 2, 6, 3, 2, 1, 6, 4, 1, 2, 1$$

Ques ②



∴ mode

$$1, 2, 5, 2, 6, 3, 2, 1, 6, 4, 1$$

No mode
Ques ②

$6 - \alpha$
3, 2, 5, 6, 3, 6, 4, 2, 1, 9, 5, 4, 1

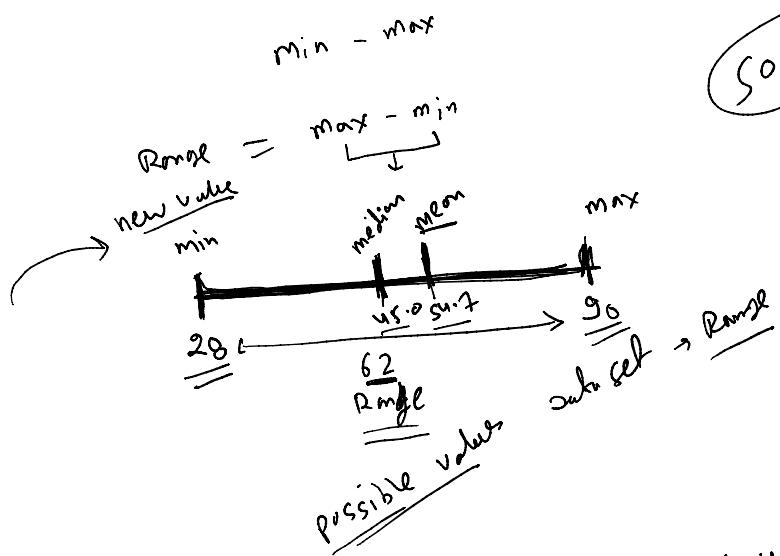
3 - 2
2 - 2
5 - 2
6 → 2
9 → 2
1 → 2

All freq same so
No mode

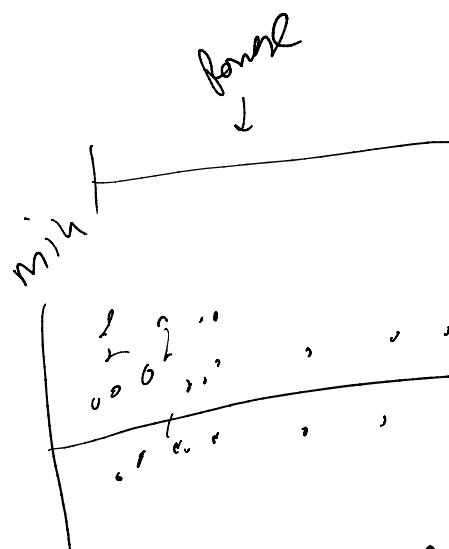
③ Measure of Spread

- ① Range
- ② Variance
- ③ Standard deviation
- ④ Inter Quartile Range
- ⑤ Five point Summary
- ⑥ deciles
- ⑦ percentiles

① Range → maximum variability



→ Symmetrize
max values are around $\frac{45.0}{2}$
min 28 and max 30



Range

+ many

,

→ Summarize

max value are or
min 28 and max 30

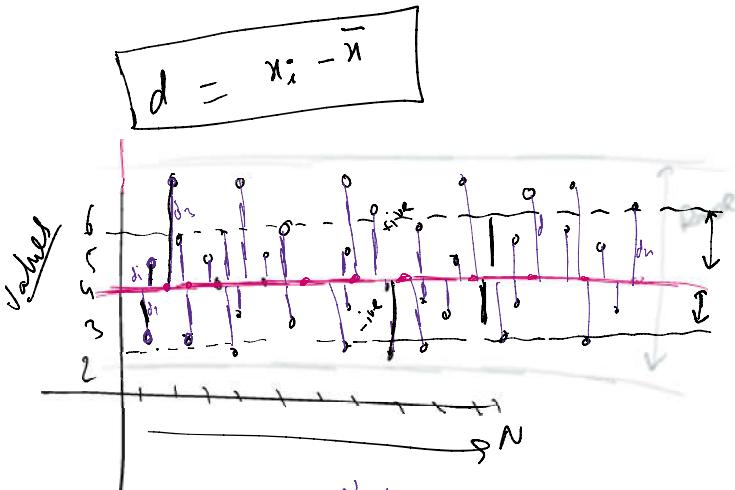
→ Variance

↳ it calculates maximum variance at data around center

X	$ x - \bar{x} $
5	0
3	2
6	1
2	3
9	4
$\bar{x} = 5$	$\sum d = 8$
	$\frac{8}{N}$

$$M \pm 2$$

$$d = x_i - \bar{x}$$



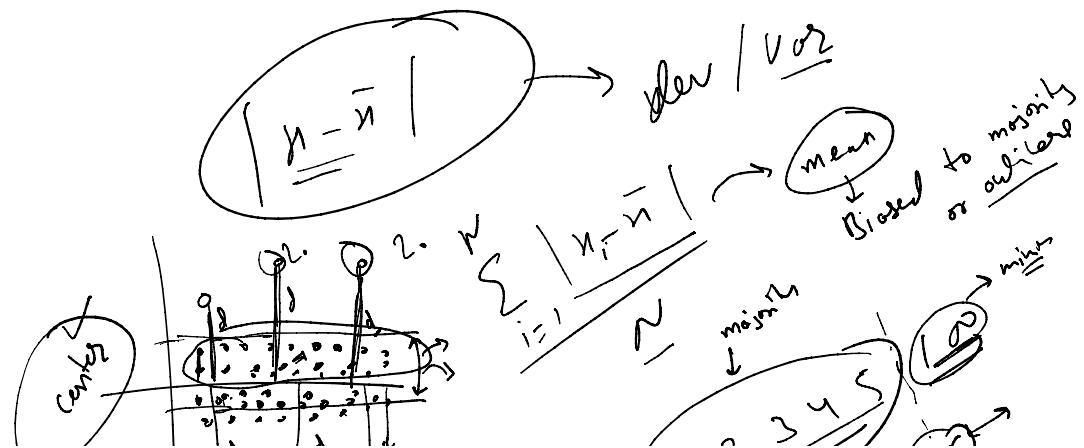
for deviation = $\sum_{i=1}^n d_i$

$$\frac{\text{abs } d}{\text{no } d} = \frac{\sum_{i=1}^n |x_i - \bar{x}|}{N}$$

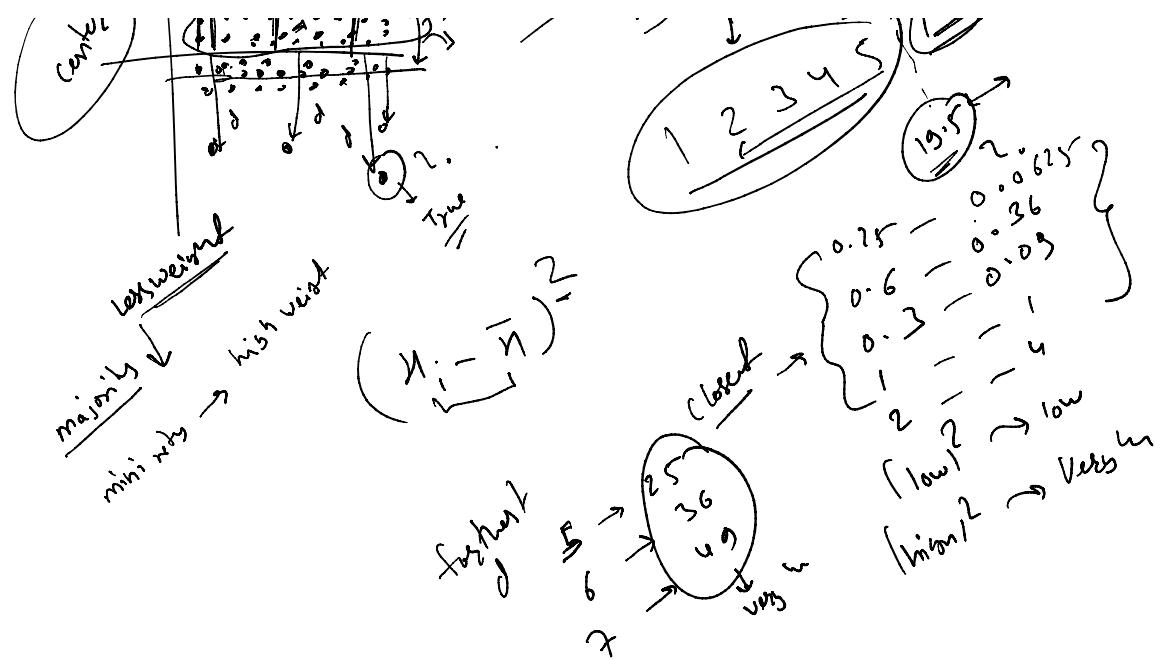
$$d^2 \rightarrow \sigma^2$$

→ Variance → maximum variability

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







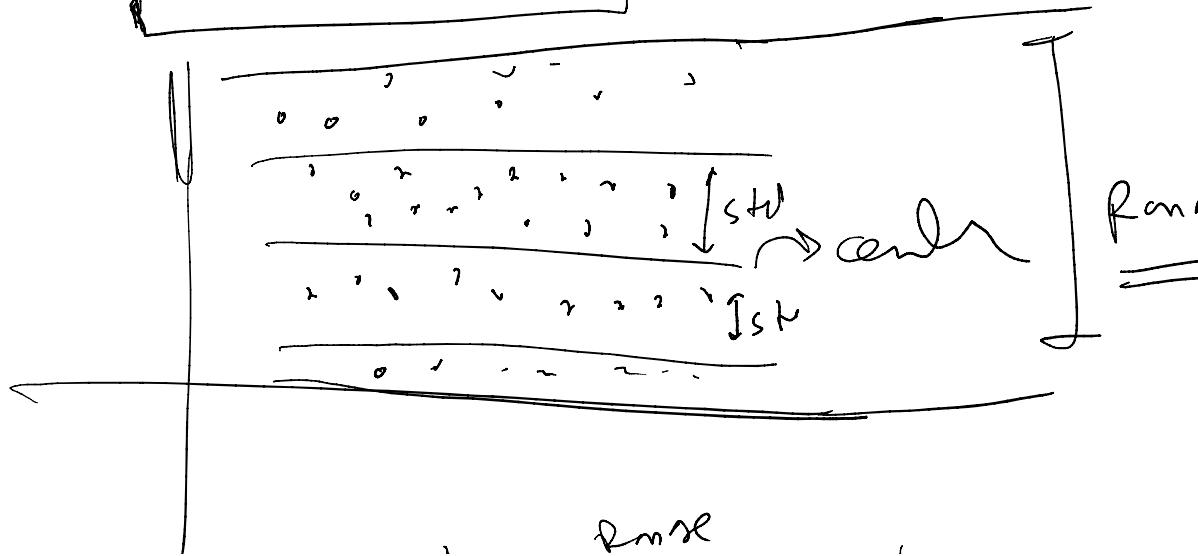
Variance

 $\sigma^2 \rightarrow \text{squared value}$
 $s^2 \rightarrow 25$
 $10 \rightarrow 100$

Standard Deviation

$$\sigma = \sqrt{\sigma^2}$$

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

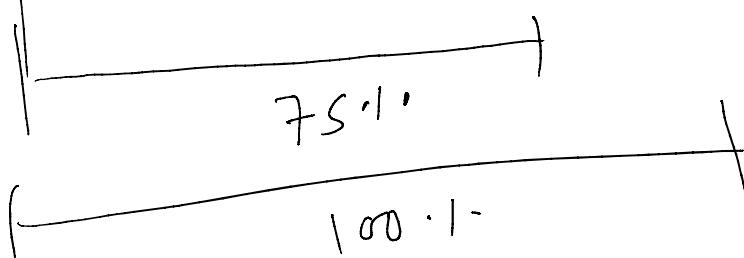
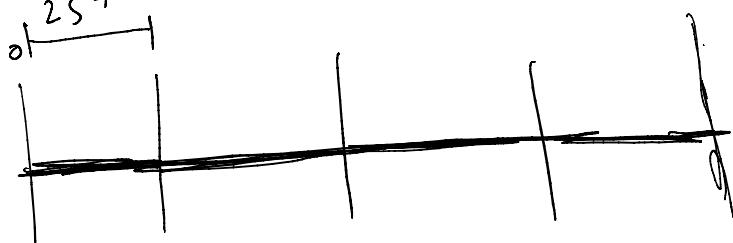
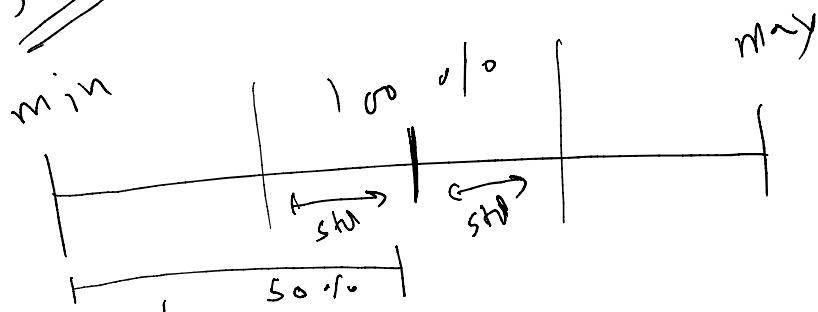


χ
=

Range
Std

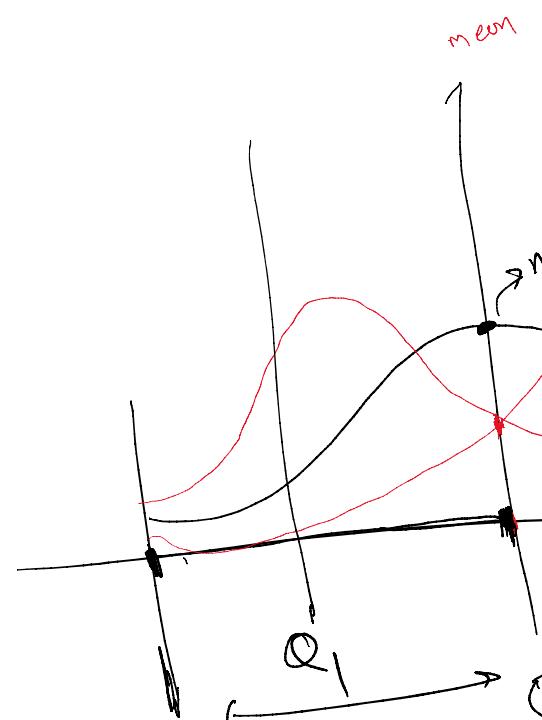
→ Quartiles

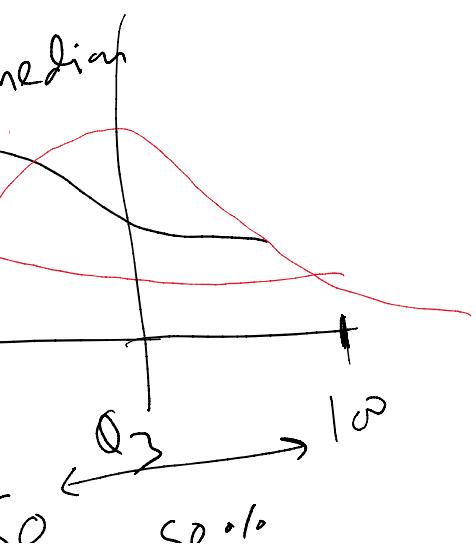
Standard



Calculation

$$Q_2 = \text{median}$$



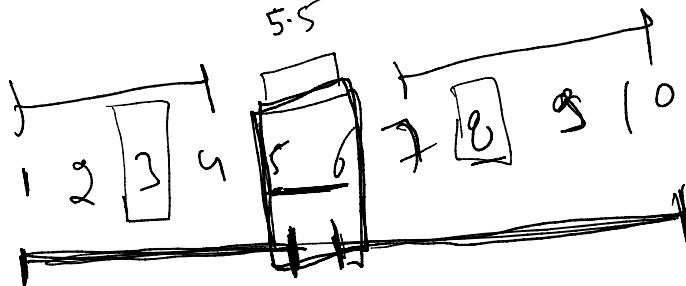


$\chi_2 -$

$Q_1 = \text{median of } \frac{\min - Q_2}{\max}$

$Q_3 = \text{median of } \frac{Q_2 - \max}{\min}$

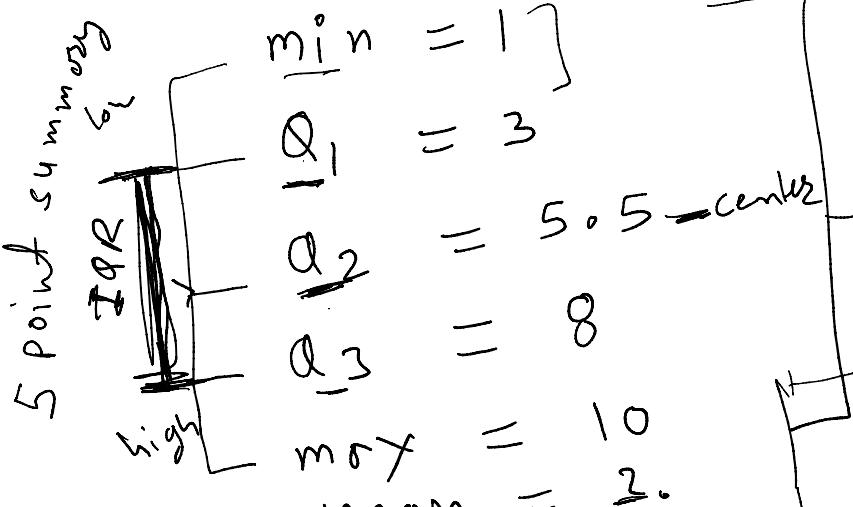
Example →



$$Q_2 = 5.5$$

$$Q_1 = 3$$

$$Q_3 = 8$$

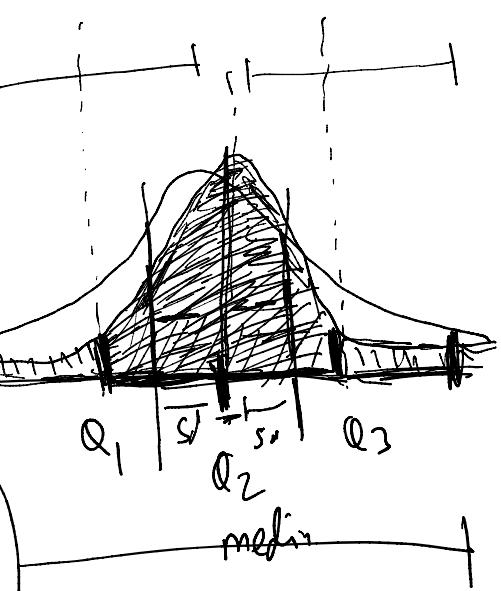


⇒ IQR → Inter Quartile Range

$$\boxed{IQR = Q_3 - Q_1}$$

~~50~~ ← →
50 ± 10
~~Q_2~~
Q_2 Q_3
- 50 - ~~75~~

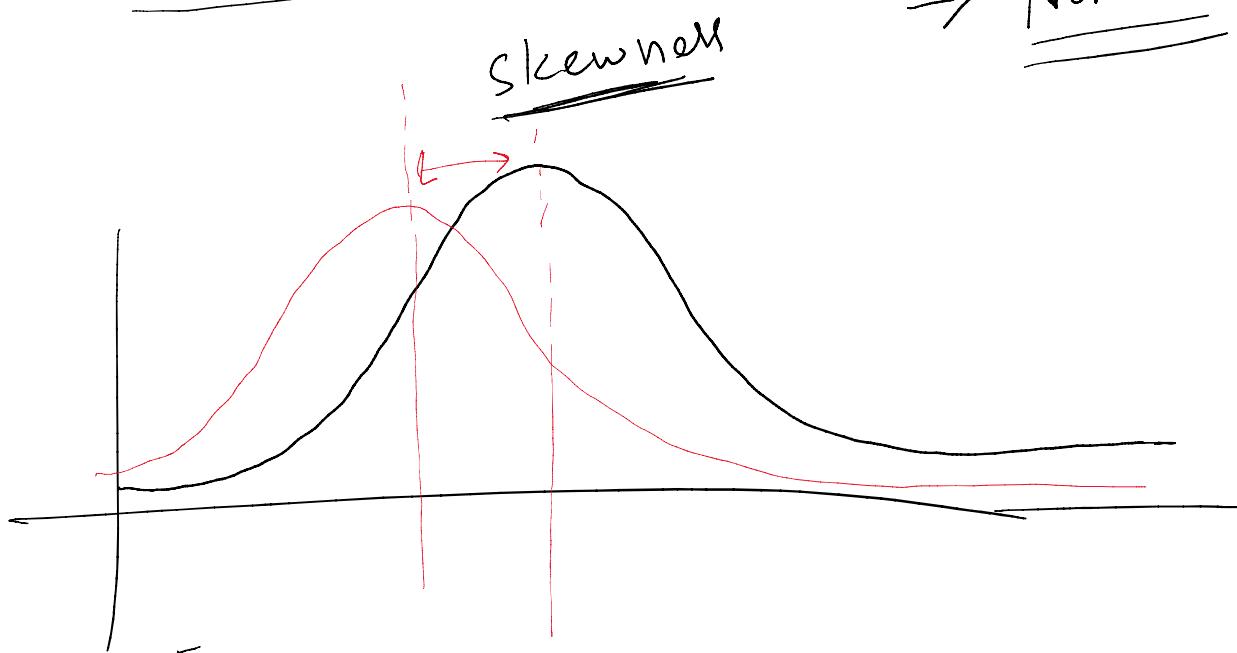
Range



\Rightarrow Decile

\Rightarrow Percentile

Shape of Data \rightarrow



\Rightarrow Normal Distributed

\rightarrow Sym

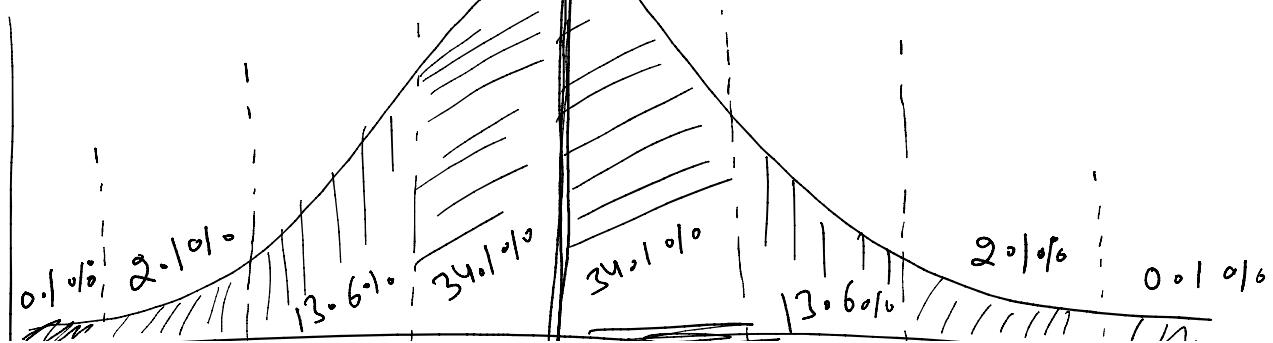
\rightarrow Asym

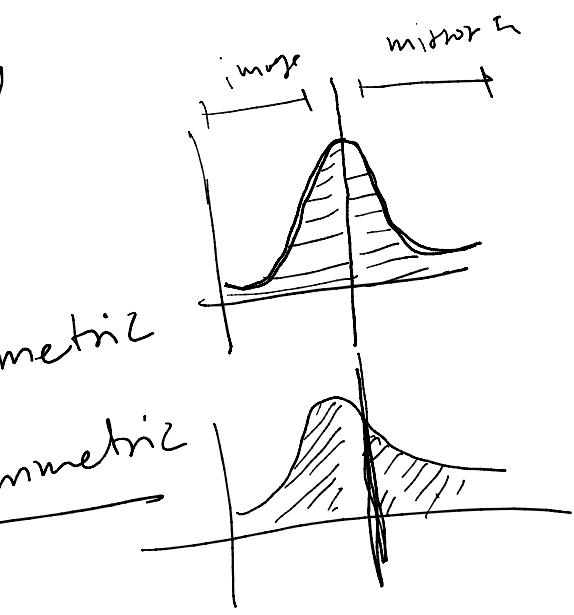
① Normal Distribution - I

0 σ +1

(i) mean

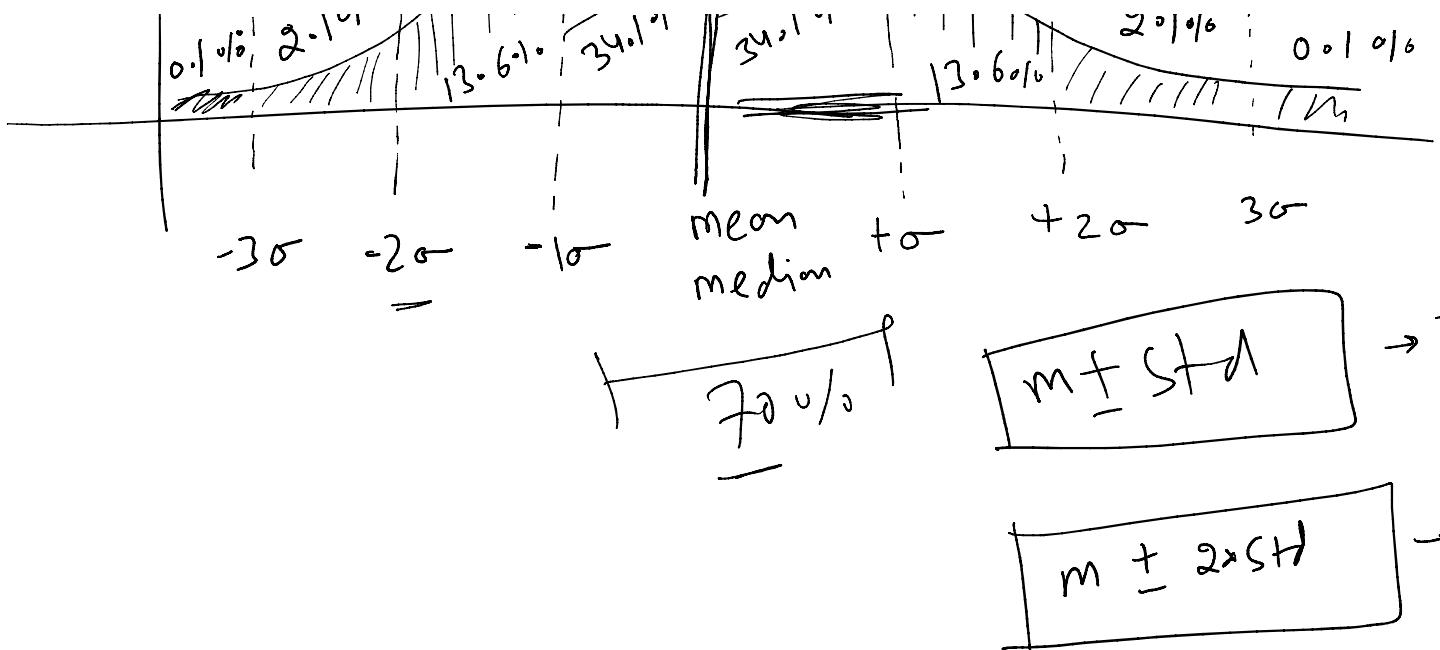
(ii) σ



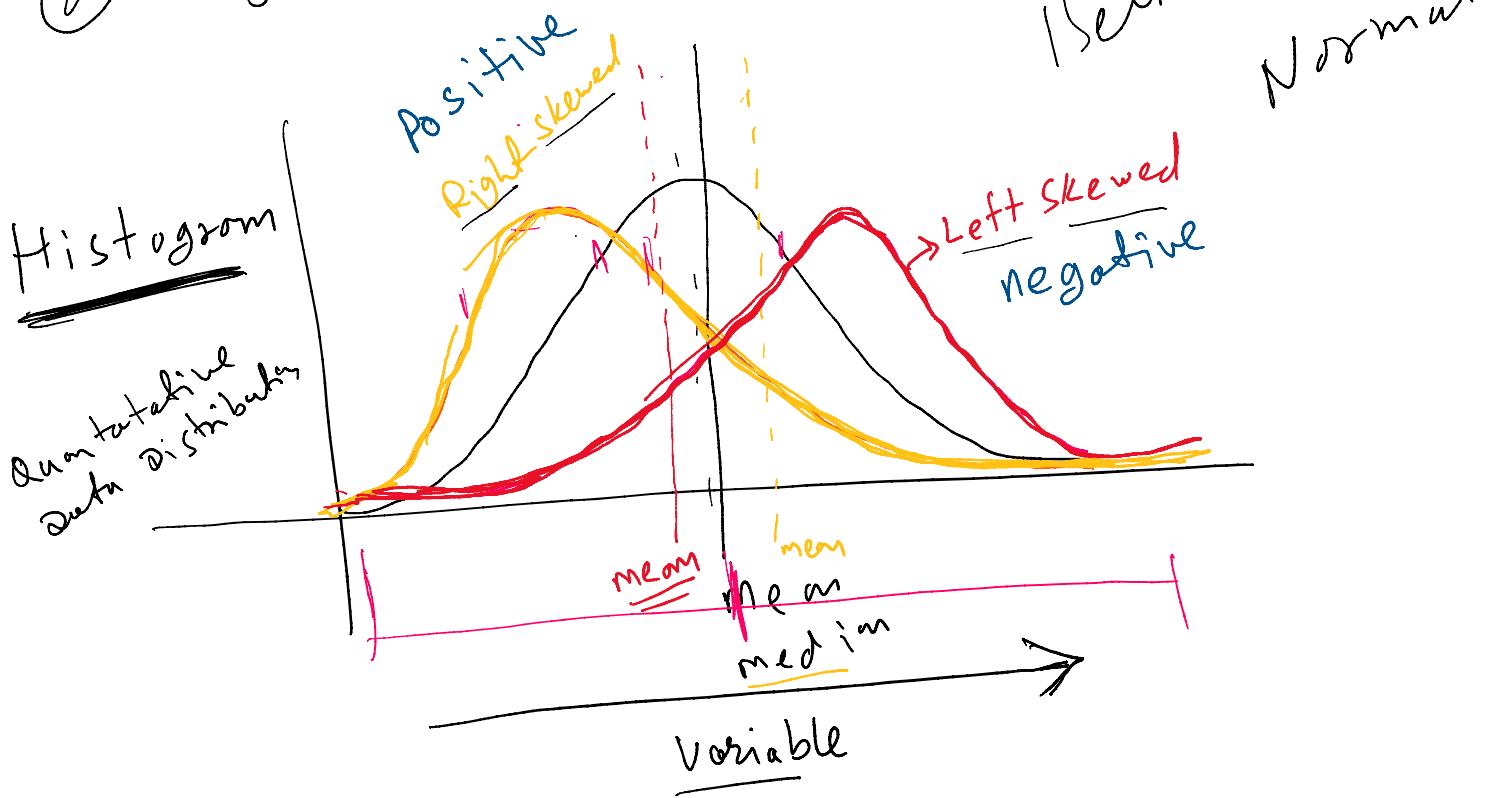


~ median
std

$$mean = 0.7$$



- ① left skewness $\rightarrow \underline{\text{mean}} < \underline{\text{median}}$
- ② Right skewness $\rightarrow \underline{\text{mean}} > \underline{\text{median}}$



$$\begin{array}{l} \text{mean} = 0 \\ \text{std} = 1 \end{array}$$

Standard
Normal
distribution

$$0 \pm 1$$

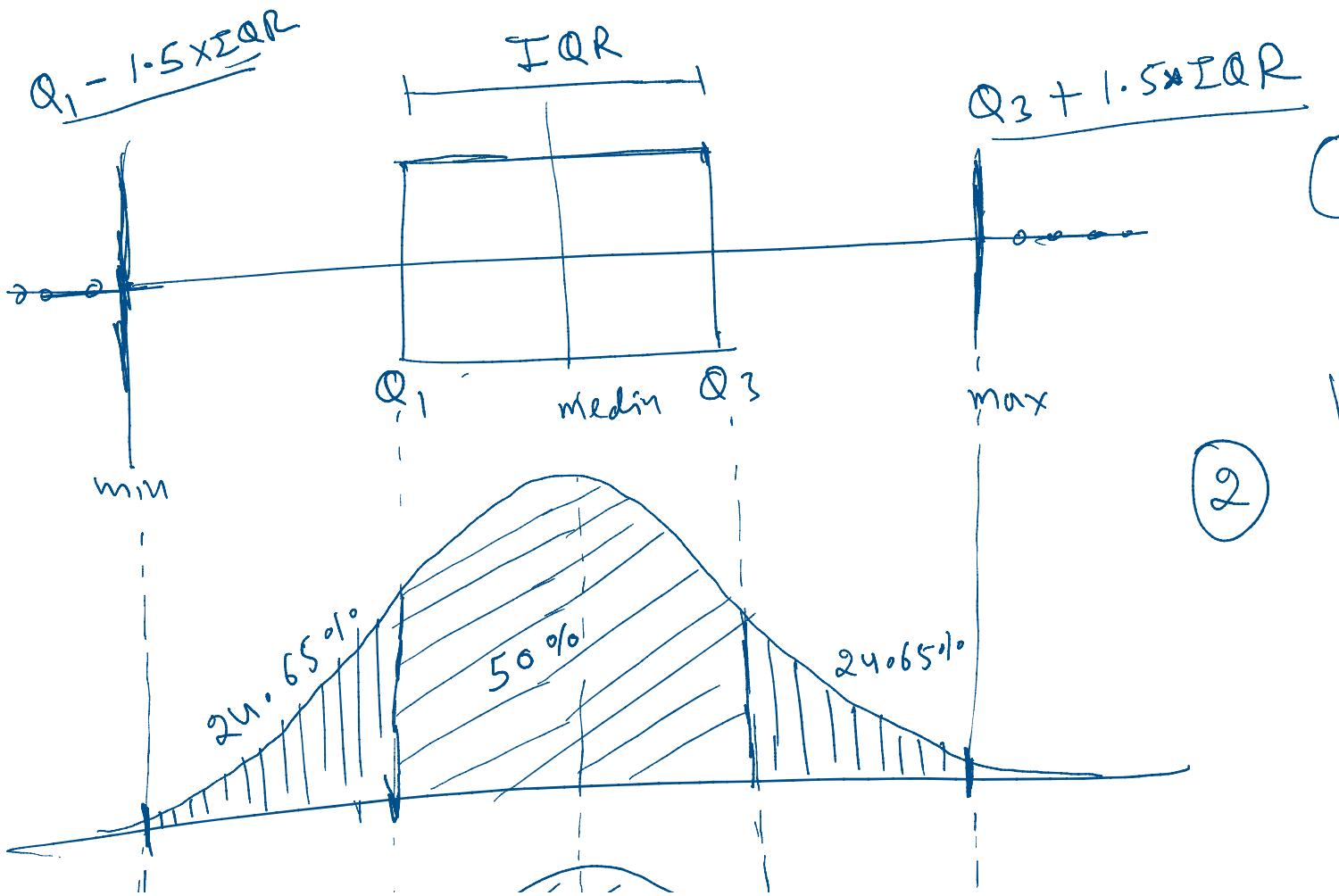
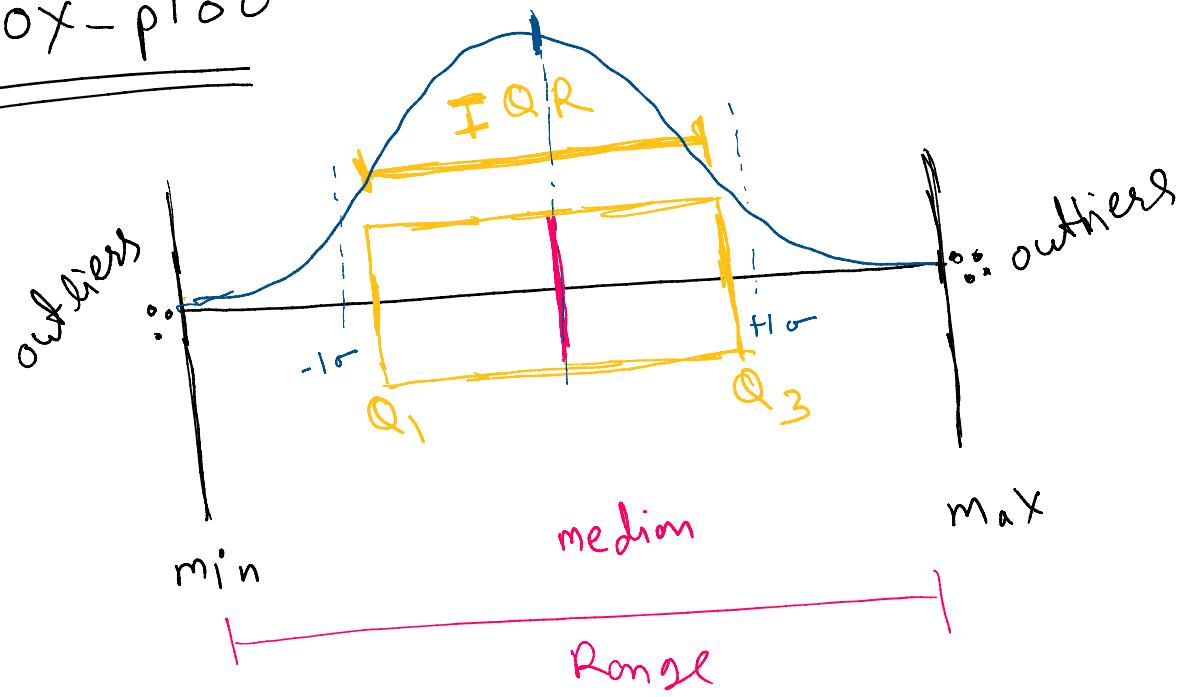
$$0 \pm 2$$

$$0 \pm 3$$

99.99

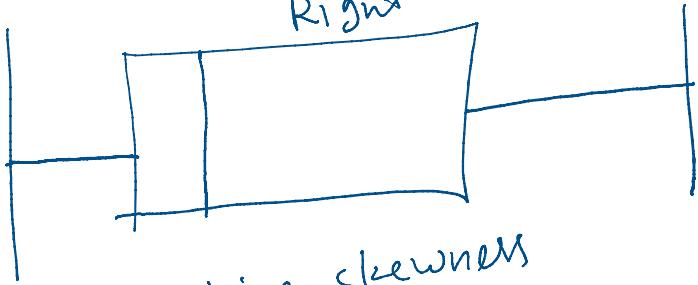
normal
distri

→ Box-plot



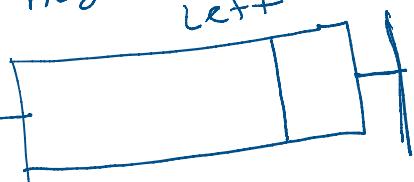
positive skewness

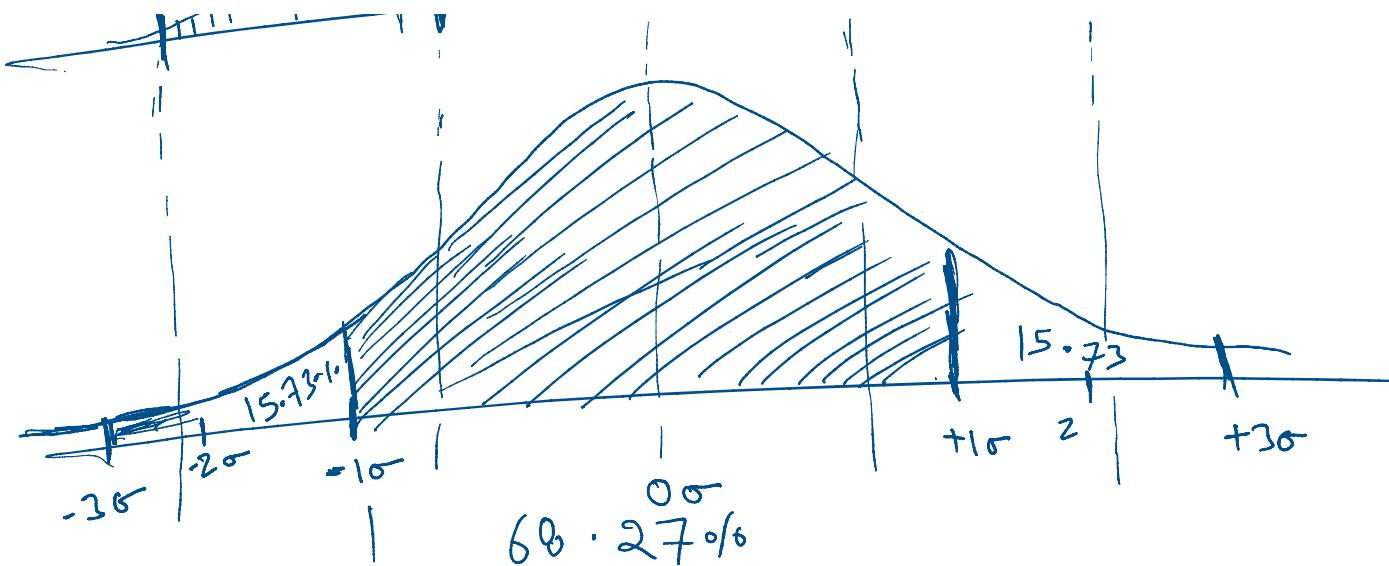
Right



negative skewness

left



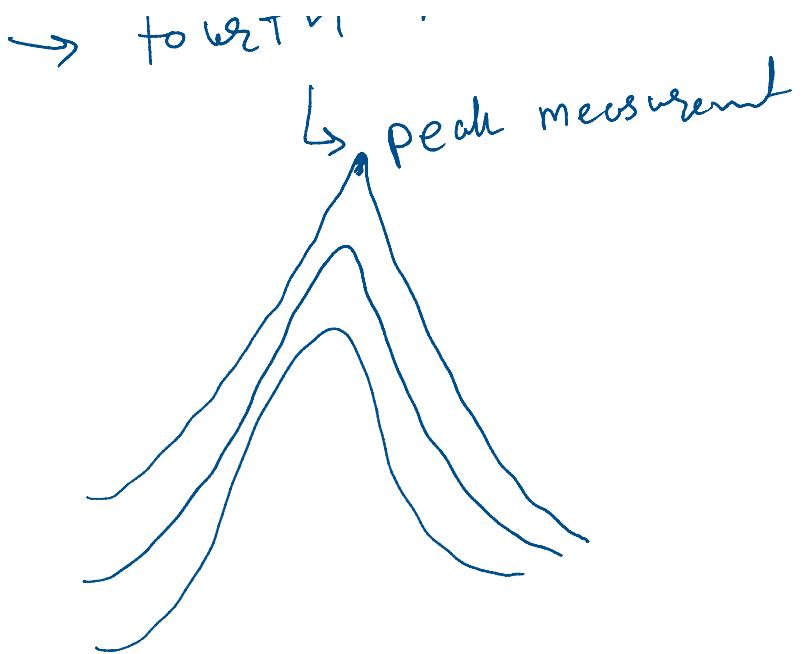


$\Rightarrow \underline{\text{Moments}}$

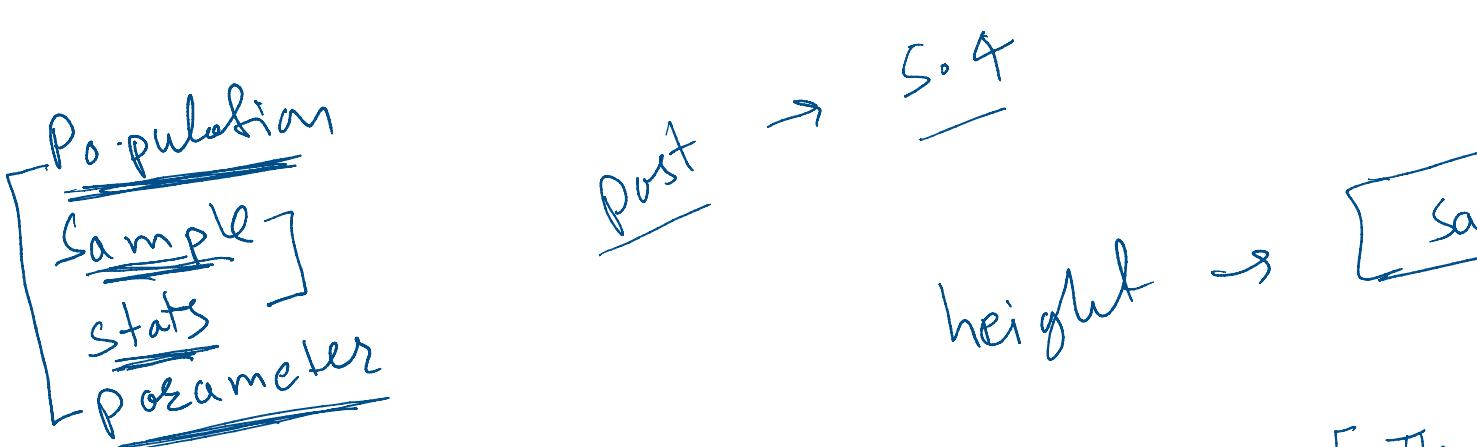
$$\mu_n = \int_{-\infty}^{\infty} (x - \bar{x})^n f(x) dx$$

- first moment \Rightarrow mean (\bar{x})
 - second moment \Rightarrow variance (σ^2)
 - third moment \Rightarrow skewness (γ)
 - fourth moment \Rightarrow kurtosis
1. moment





- ① Descriptive Stats → Distribution of data
- ② Inferential stats → predictive Analysis



① Null hypothesis (h_0)

② Alternative hypothesis (h_A)

(Sample)

(Population)

sample

[There is no difference at all]
[There should be difference in data]

↳ ② Alternative hypothesis ✓

→ hypothesis testing ✓ X

→ P-value

if $P > 0.05$

we

population

⇒ Tests

1 sample proportion test

- ①
- ②
- ③

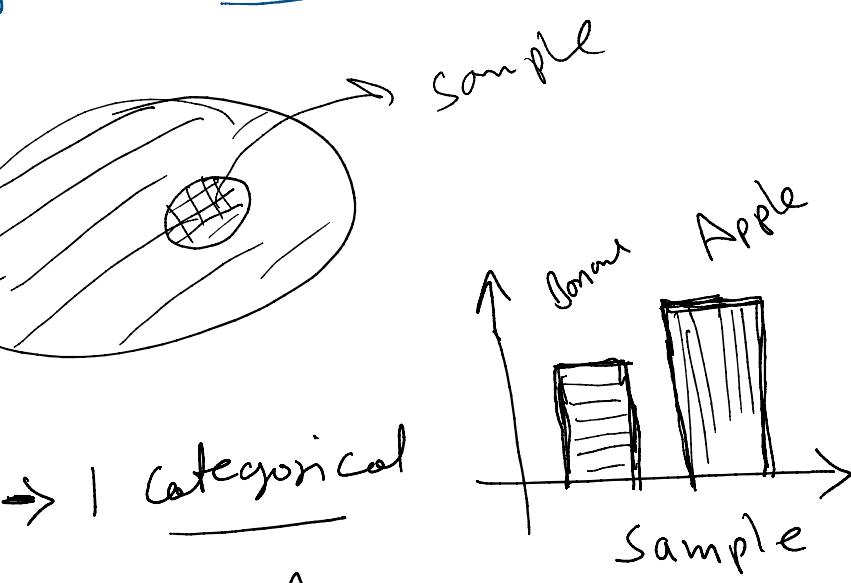
Chi-square test → 2 categories

t-test → single numerical

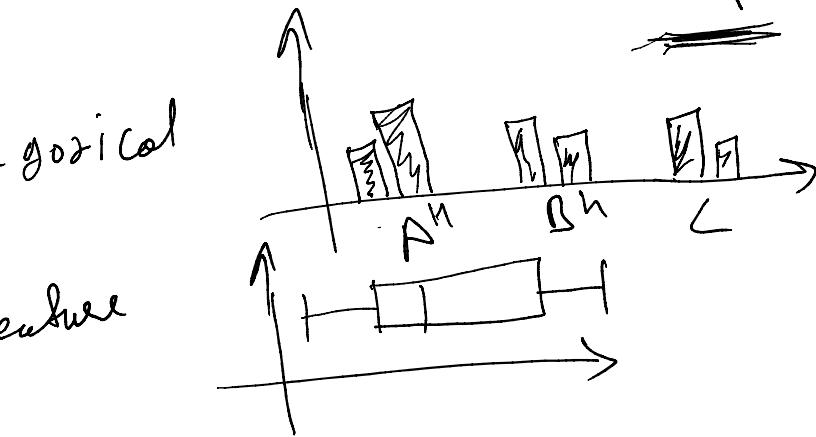
· 1st ANOVA → one numerical

$$P(h_0) + P(h_A) = 1$$

reject null hypothesis

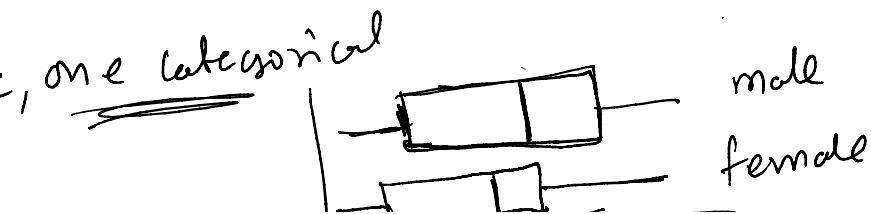


→ 1 categorical



categorical

variable



(C)

(S)

(S)

t-test, ANOVA → statistic

Correlation → 2 numerical

Scat

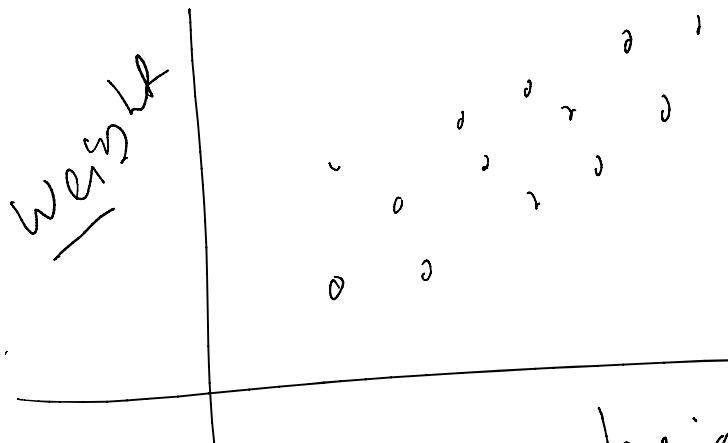
Distribution
Box Plot
Scatter
Bar
Pie

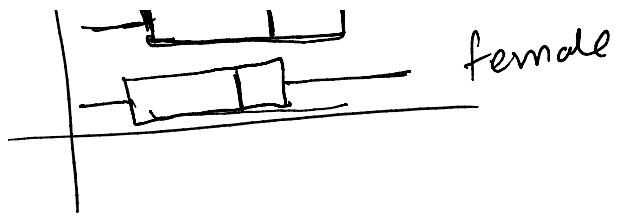
numerical
categorical

Correlation Coeff → +/+
P-value

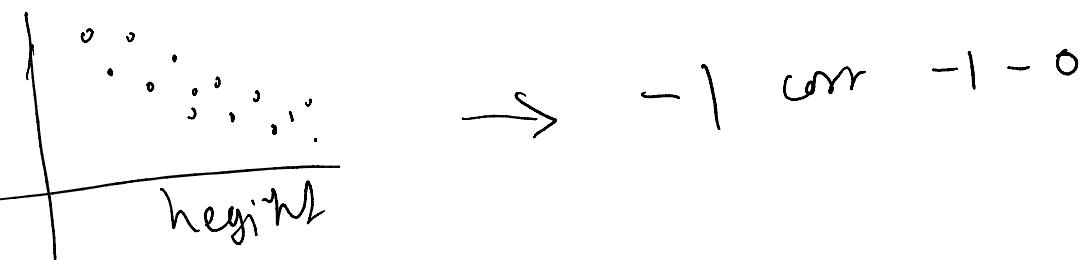
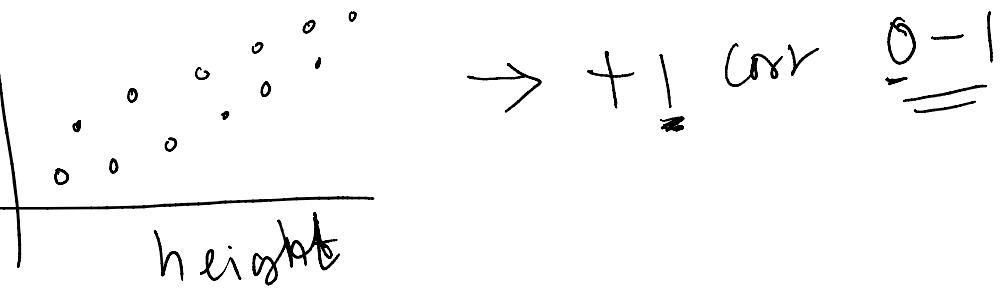
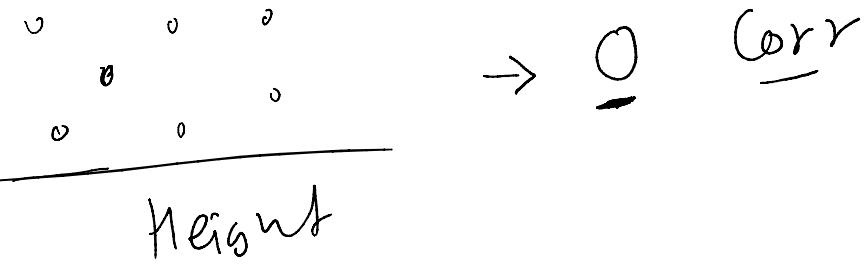
Positive
negative

Weighted





-er plot



Co-var

height

→ $X = \text{height}$

→ $y = \text{weight}$

Co-var →

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

$$\sigma_y^2 = \frac{\sum_{i=1}^N (y_i - \bar{y})^2}{N}$$

Co-var $\sigma_{xy}^2 = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{N}$

Co-var $r = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^N (x_i - \bar{x})^2 \sum_{i=1}^N (y_i - \bar{y})^2}}$

$$\frac{1}{n} \sum_{i=1}^n$$

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

N

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

$$[-\infty \text{ to } +\infty]$$

$$\begin{aligned} \text{Population} &= N \\ \text{Sample} &= N-1 \end{aligned}$$

$$\begin{aligned} \text{Std Dev} &> \text{SDev} \\ \therefore (\text{d}) &\rightarrow \text{SDev} \end{aligned}$$

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

$$Co-var = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

neg
pos

$$-1 \leq r \leq 1$$

$$r = \frac{\text{Cov}(x, y)}{s_x s_y}$$

$$\text{Cov}(x, y) = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

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

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

$$RMS(d) \rightarrow S^m$$

$$\frac{-d}{Std} \rightarrow -1, 1$$

$$(n - \bar{n})(y - \bar{y})$$
$$\frac{N}{\bar{n})^2}$$

$$\sum_{i=1}^n$$

$$S_g = \sqrt{\sum_{i=1}^n}$$

→ goal → S_{fut5}
→ matrix & vectors

$$\frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2$$