

Agenda

* Mean deviation

① Variance

② Std dev.

* Measures of Symmetry

* Skewness

* Kurtosis

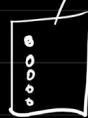
* Covariance and Correlation

* Probability

{

infantil

Measures of spread \rightarrow range, IQR, Percentile, Boxplot.



How ??

Compare Kusum's marks
as compared avg
of the class.

Avg marks = 2:1

* Mean deviation

each of the data point



$$\sum_{i=1}^n |x_i - M|$$

n = no. of dp

On an avg how much away
each date point is away
from mean value ??

$$\frac{1+2+1+2+0}{5} = \frac{6}{5} = \boxed{\underline{\underline{1.2}}}$$

* Variance → The average of the squared differences from the mean

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

population mean.

$$\begin{aligned}\text{Var}_{\text{Sample}} &= s^2 \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}\end{aligned}$$

$N \rightarrow n$

$\bar{x} \rightarrow \bar{x}$

$n \rightarrow n-1$ (denominator)

Variance

→ Calculate mean

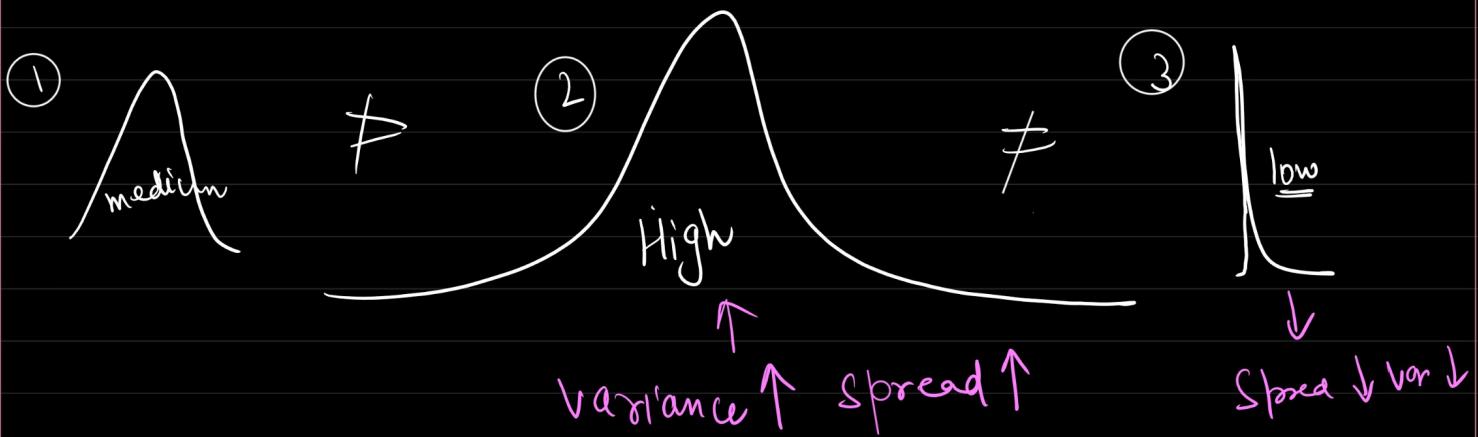
→ for each no in date, subtract the mean and no.

→ Square the difference

→ Calculate the avg square of difference

$$\text{date} = \{1, 2, 3, 4, 4\}$$

<u>(Part of some sample)</u>	x	\bar{x}	$x - \bar{x}$	$(x - \bar{x})^2$	$\underline{\text{Var}_{\text{sample}}} = \frac{6.82}{n-1} = \underline{\underline{s^2}}$
1	2.83	2.83	-1.83	3.34	
2	2.83	2.83	-0.83	0.68	
3	2.83	2.83	0.17	0.03	
3	2.83	2.83	0.17	0.03	
4	2.83	2.83	1.17	1.37	
4	2.83	2.83	1.17	1.37	
<hr/>				$\underline{\underline{6.82}}$	$\Rightarrow \underline{\underline{1.37}}$
$\underline{\underline{\sum x_i}} \Rightarrow \frac{17}{6} \Rightarrow 2.83$					



* Standard deviation

Standard devⁿ is a measure of how spread out numbers are.

Square root of Variance.

$$\text{Std} = \sqrt{\text{Var}} = \sqrt{1.37} = 1.17$$

Standard devⁿ of population $\Rightarrow \sigma = \sqrt{\text{Var}_p}$

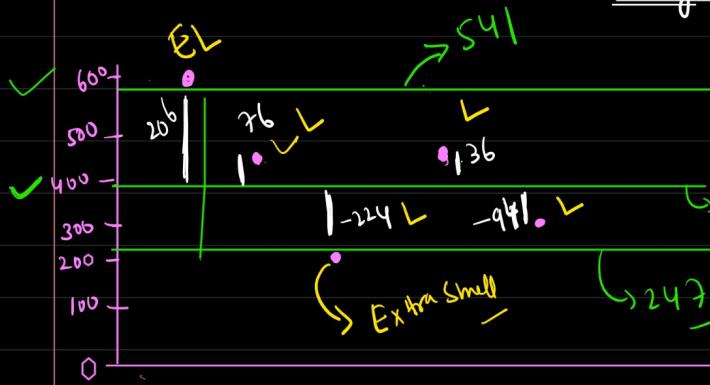
" " " " " $\Rightarrow s = \sqrt{\text{Var}_s}$

Use Cases

This is the class

✓ Height

→ Soumya 600cm, Rahul 470cm, Deep 170cm, Farhan 430cm, Suman 300cm



$$\bar{M}_{\text{ean}} = \frac{600 + 470 + 170 + 430 + 300}{5}$$

$$394$$

$$\sigma^2 = 206^2 + 76^2 + (-24)^2 + (36)^2 + (-94)^2$$

$$\sigma^2 = 21704$$

Variance $\rightarrow 21704$

Standard devⁿ is

a standard way of knowing how normal

Standard devⁿ

$$\sigma = \sqrt{21704}$$

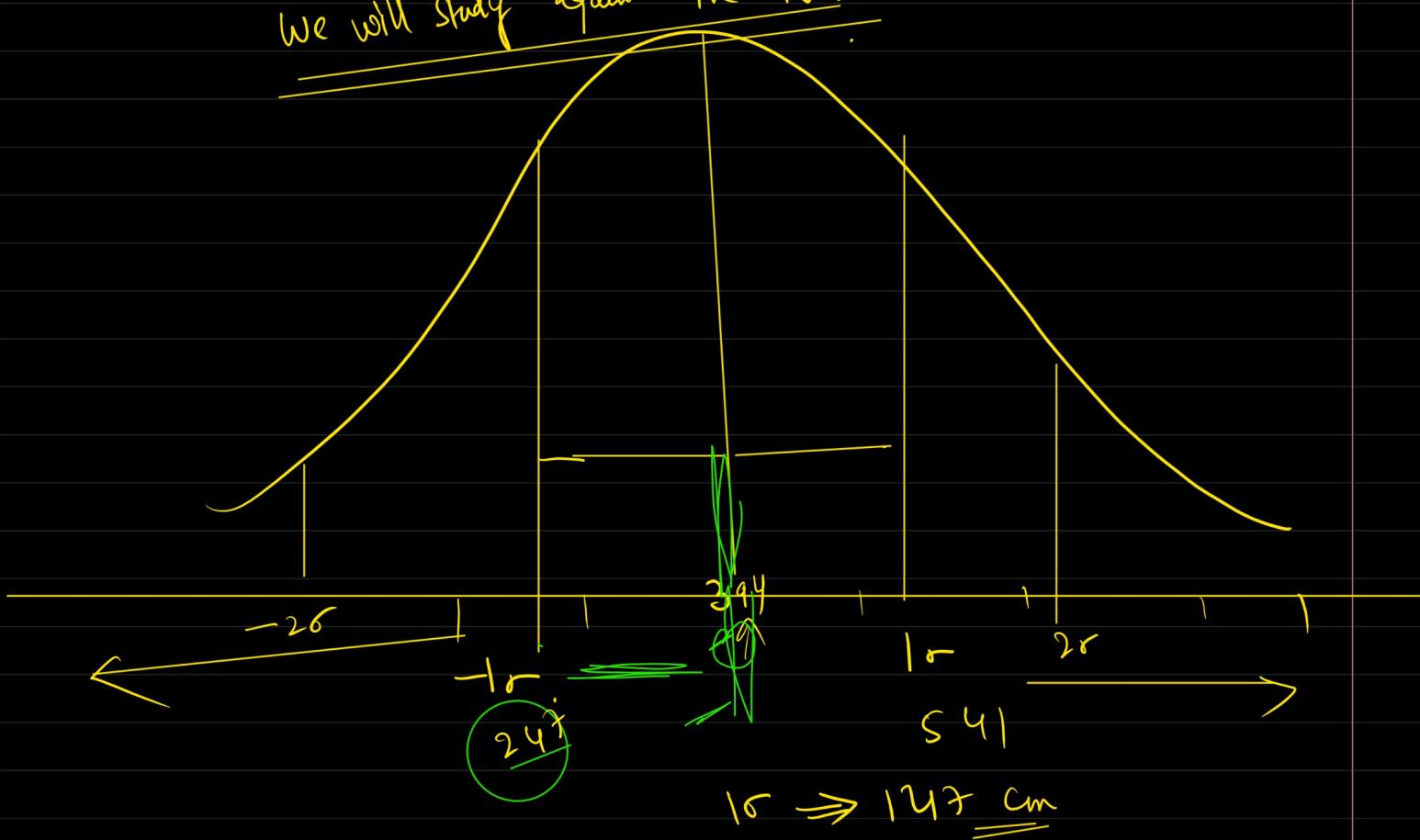
$$= 147$$

large, extra large is something.

$$394 + 147 = 541$$

$$394 - 147 = 247$$

We will study again in N.D.



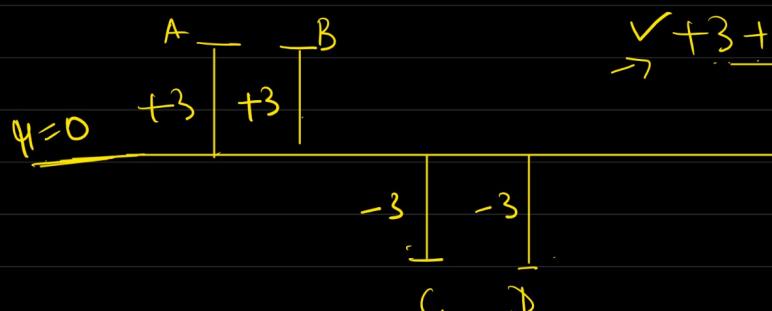
What was the mean deviation

$$\text{Var}_{\text{p}} = \sigma^2 = \sum_{i=1}^n \frac{(x - \mu)^2}{N}$$

To calculate standard dev.

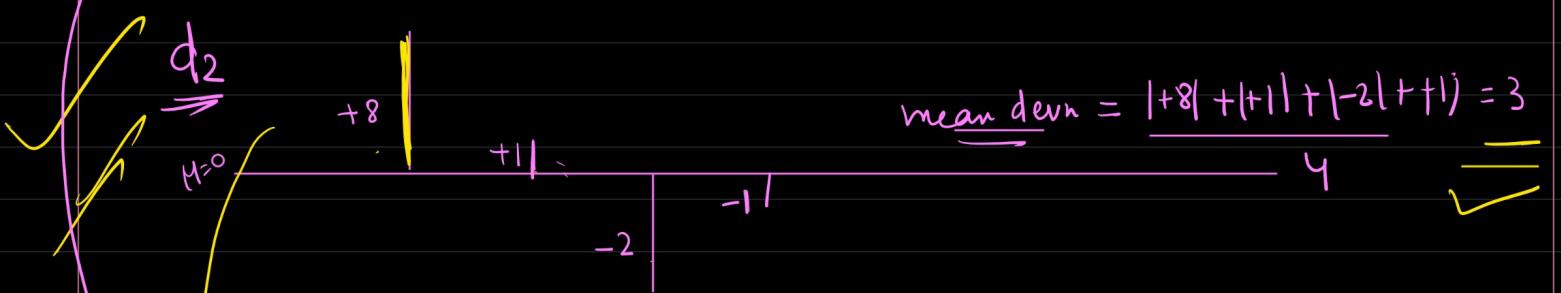
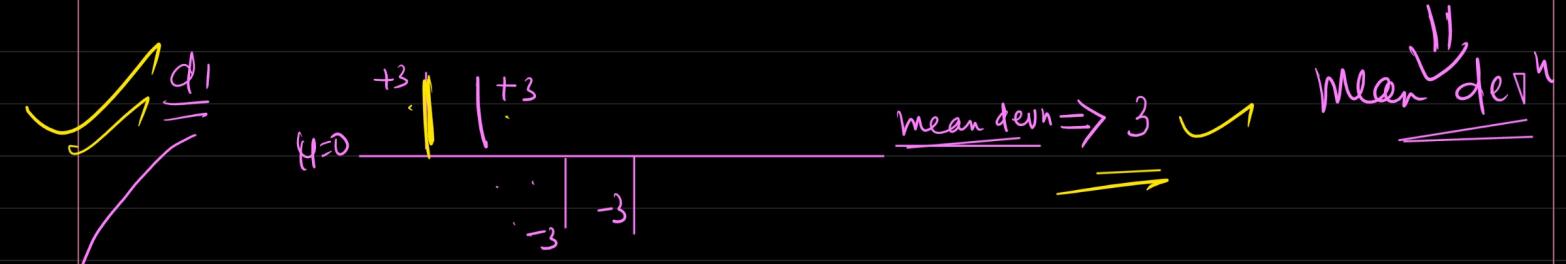
$$\sigma = \sqrt{\text{Variance}}$$

$(x - \mu)^2 \rightarrow \underline{\text{why.}}$



+ve and -ve
negative

Can I take abs → $\frac{|3| + |3| + |-3| + |-3|}{4} = \frac{12}{4} = 3$



$$\sqrt{\frac{3^2 + 3^2 + (-3)^2 + (-3)^2}{4}} = \sqrt{\frac{36}{4}} = 3$$

$$\sqrt{\frac{8^2 + 1^2 + (-2)^2 + (-1)^2}{4}} = \sqrt{\frac{64 + 1 + 4 + 1}{4}} = 4.4$$

* $\overline{\text{Var sample}} = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$ $\frac{(x_i - \bar{x})^2}{N}$

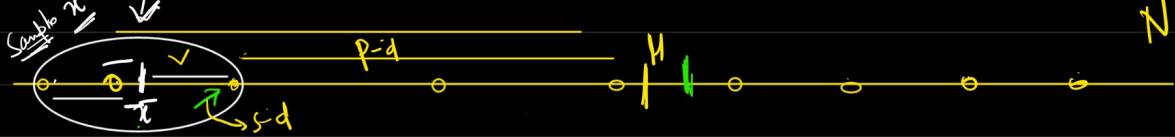
Why $n-1$??

Bessel correction / unbiased estimator



We use $n-1$ rather than n is because sample variance will be unbiased estimator.

$x_i - \mu$ X (you don't have a class to pop)



Why $- \underline{N-1}$ only

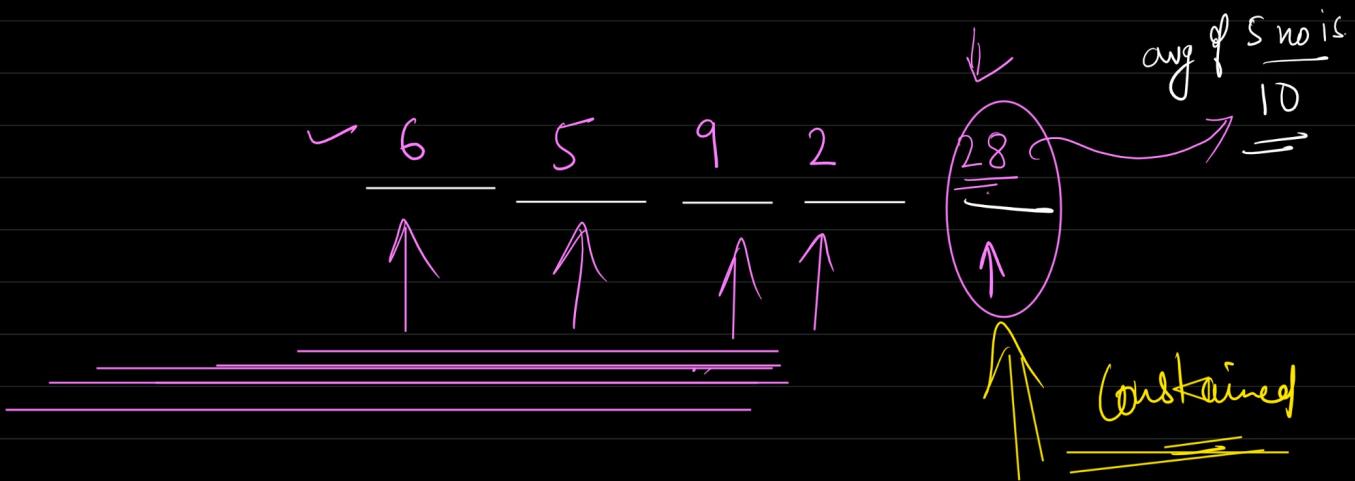
Why not $N-2$, $N-3$,

$N-4$ $N-5$

degree of freedom



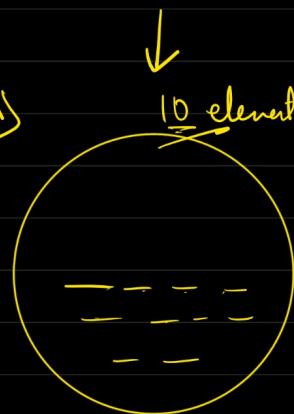
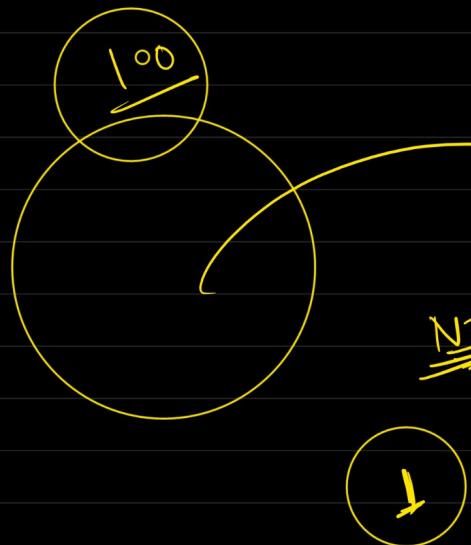
maximum no of logical independent
variables.



$N \underline{N-1}$

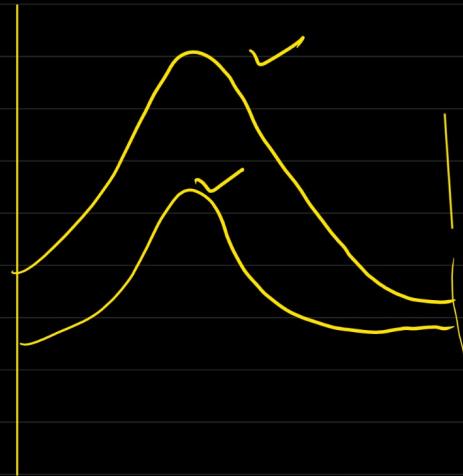
$$\frac{12}{\text{---}} \quad \begin{array}{c} 28 \\ \hline \end{array}$$

avg 20



No of
Positions
are changeable

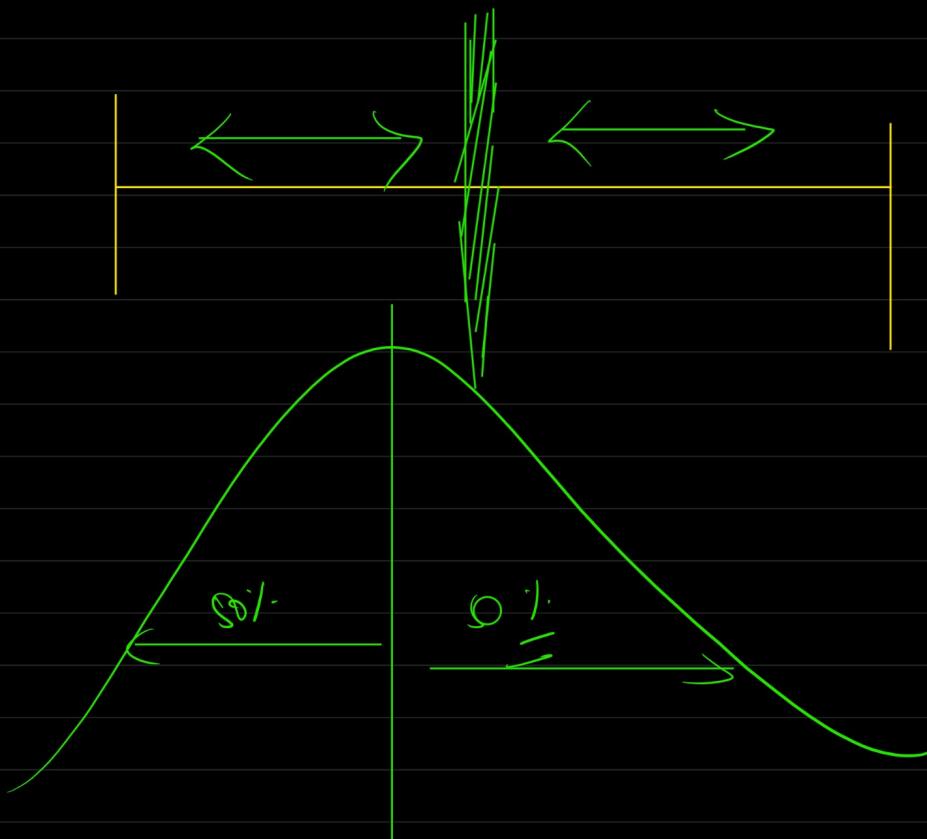
* Measures of Symmetry



Descriptive Stats

- MCT
- MD

Measure of
Symmetry

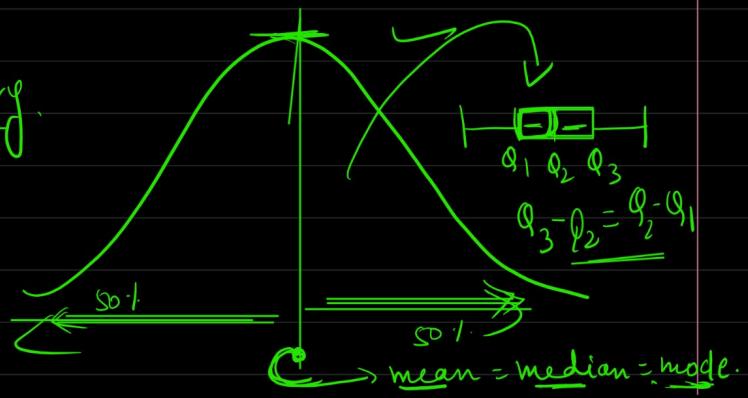


① Skewness

② Kurtosis

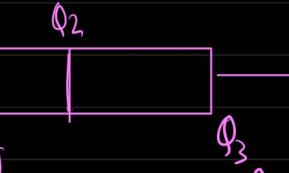
* Skewness → dataset's symmetry.

$$\text{Skewness} = 0$$



Skewed | non-symmetric

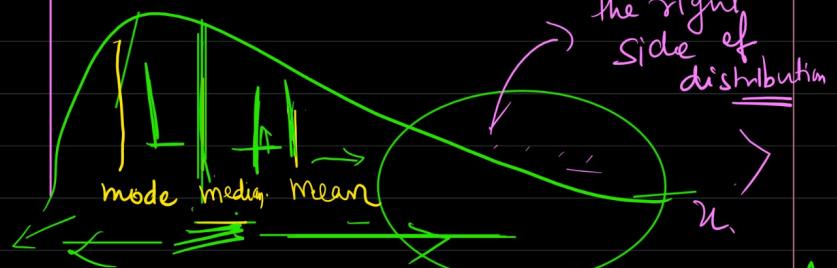
① Positive Skewed (Right Skewed data)



$$Q_3 - Q_2 \geq Q_2 - Q_1$$

tail is on

median \rightarrow physical midpoint



mean > median > mode

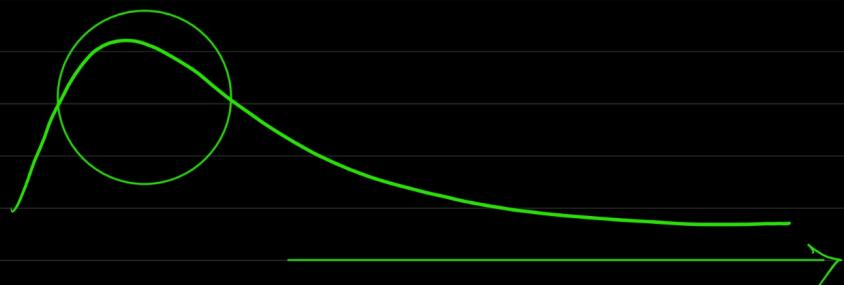
median — 50%

50%

data is incline

mean

Example \rightarrow Distribution of salary in a company.



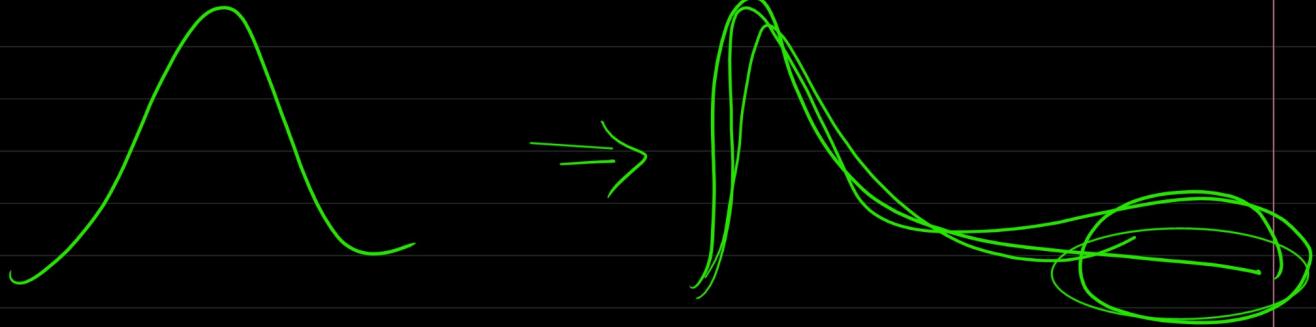
Distribution of marks in difficult exam



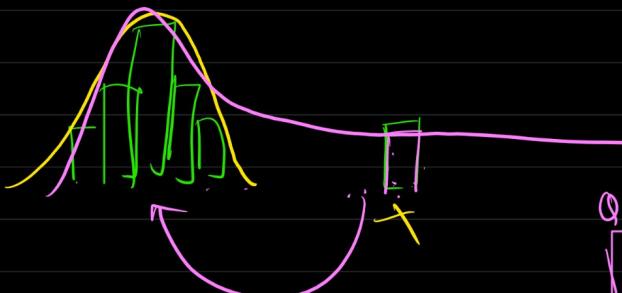
Article 340



Why?



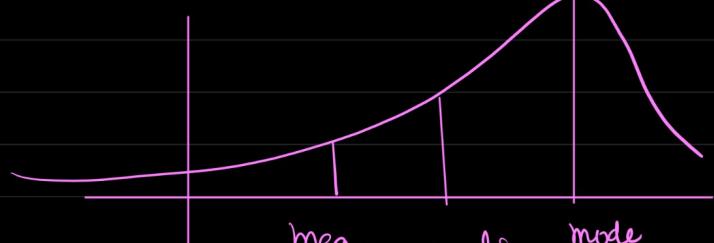
Outliers



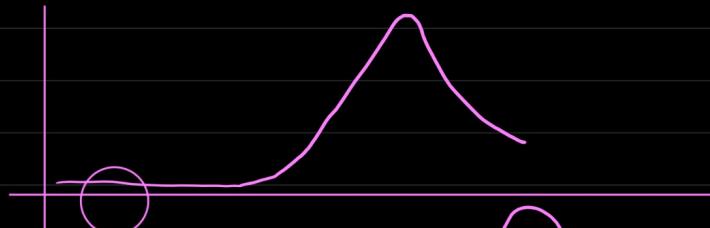
$$Q_1 \quad Q_2 \quad Q_3 \\ Q_2 - Q_1 > Q_3 - Q_2$$

② Left Skewed data

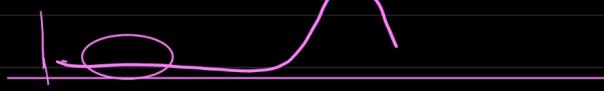
mode > median > mean



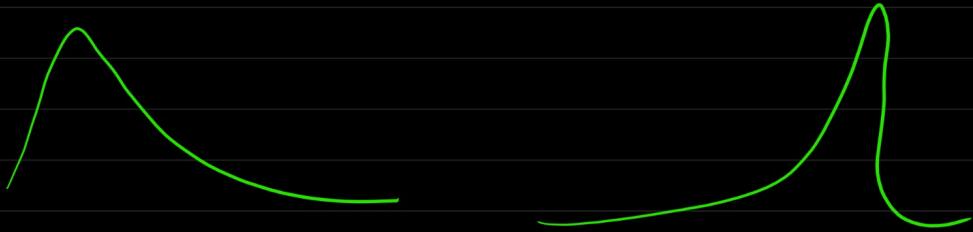
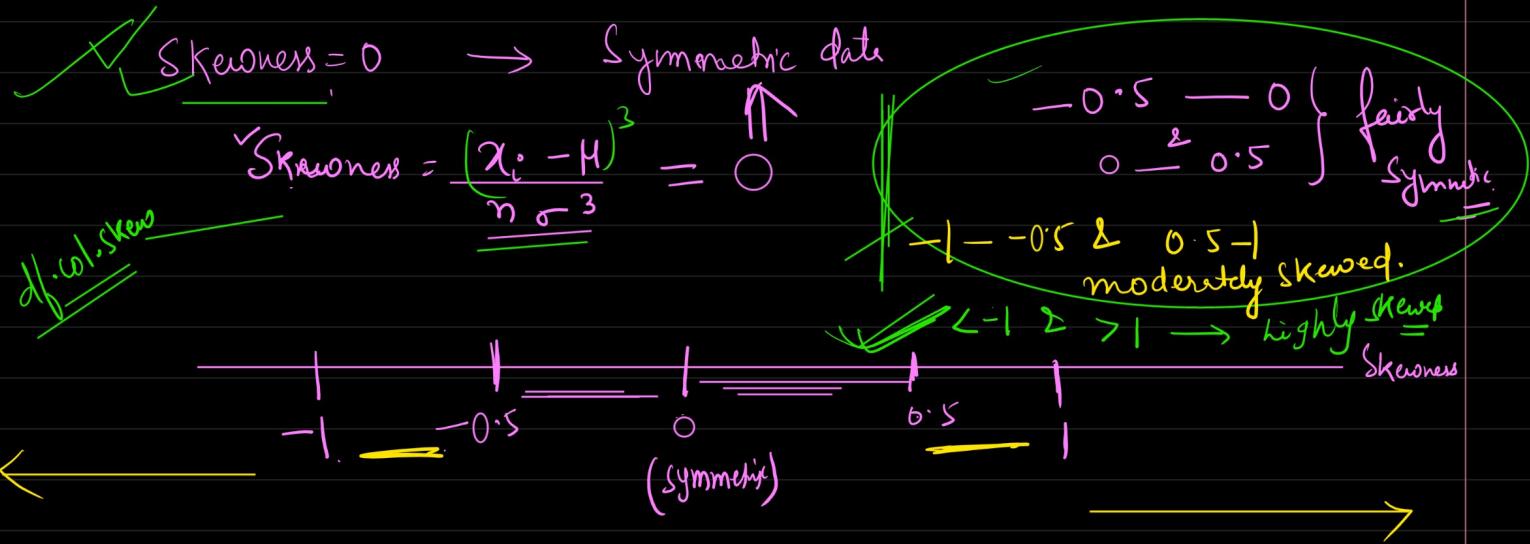
① Death rate



② Mark in easy exam



③ Wealth distib



Reasons:-

→ Many ML models wants date to be symmetric

→ Visualization → Dist plot
→ Q-Q-plot.

Statistical

→ Skewness

df-wl-skew

-1 < > +1

Transformation

- || → treat the outliers
- || → Box-Cox transformation | Yeo Johnson
- || → Exponential trans

↳ reciprocal
log transformation



(1st moment)

mean =

$$\frac{\sum (x_i - \bar{x})}{N}$$

→ mean is computed taking \bar{x} as reference

(2nd moment) Varp =

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

→ Here mean is the reference

(3rd moment) Skewness α_3

$$= \frac{\sum (x_i - \bar{x})^3}{N \sigma^3}$$

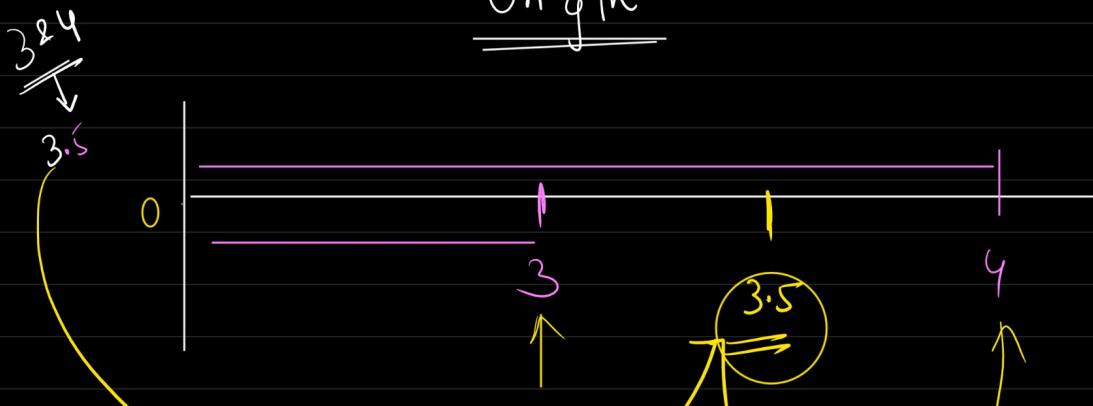
→ mean is the reference.

(4th moment) Kurtosis α_4

$$= \frac{\sum (x_i - \bar{x})^4}{N \sigma^4}$$

mean $\frac{\sum (x_i - \bar{x})}{N} \rightarrow$ On an avg how much a specific value or data point lie away from

On origin



Not necessarily

Boxplot

Skew

outlier

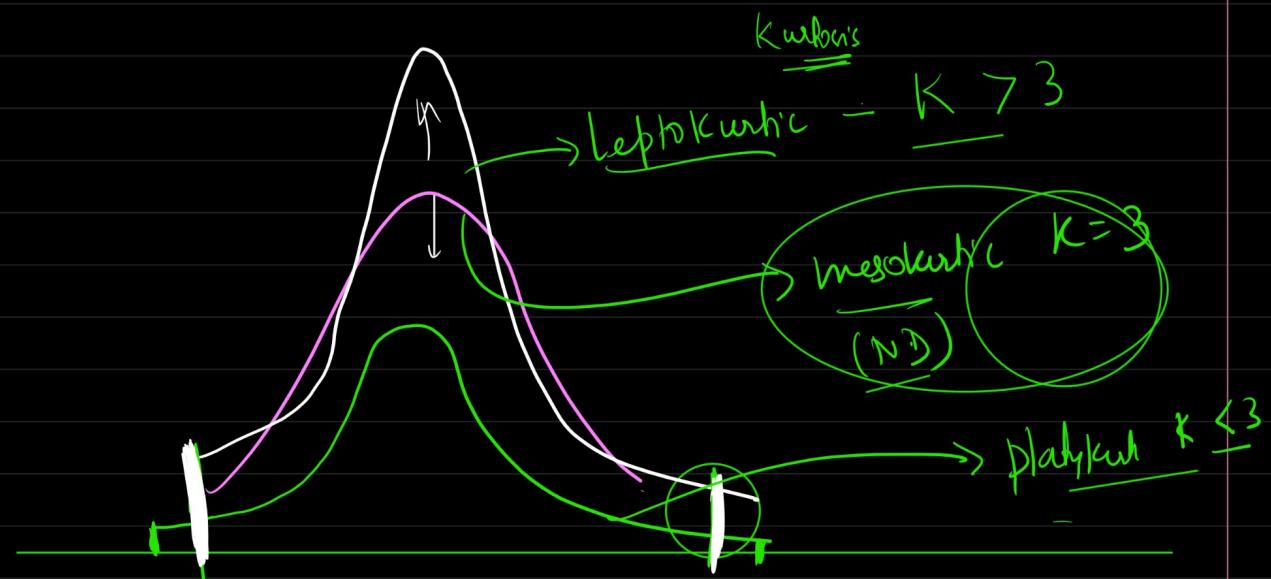
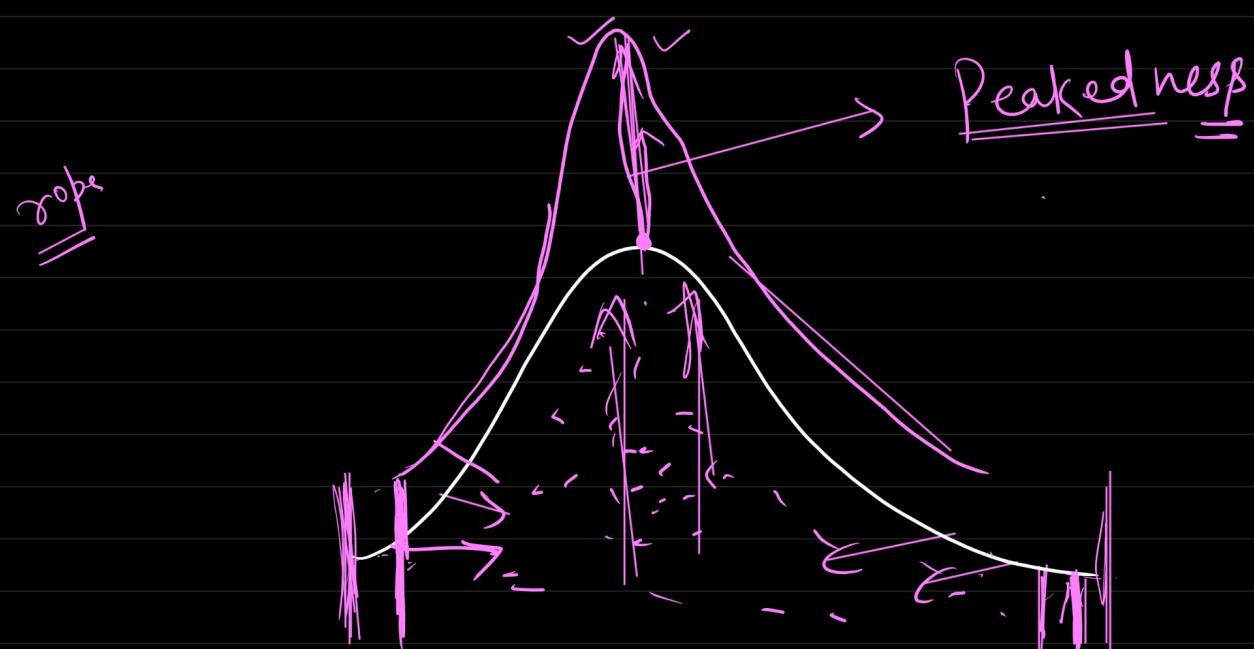
It can't have outliers ??

* Kurtosis

$$q_4 = \frac{(x_i - \mu)^4}{N s^4}$$

tail of the distribution | fattness of the tail

tail.



$K > 3 \rightarrow$ Leptokurtic \rightarrow tail is fat \rightarrow Many outliers

$K = 3 \rightarrow$ Meso

$K < 3 \rightarrow$ Platykurtic