

* Agenda

- Quartiles
- IQR
- Boxplot
- ⇒ Variance
- Std dev.
- * Measure of Symmetry
- * Correlation & Covariance
- * Prob dist.
- * CLT

* Quartiles ..

Problem with range \rightarrow if is affected by outliers

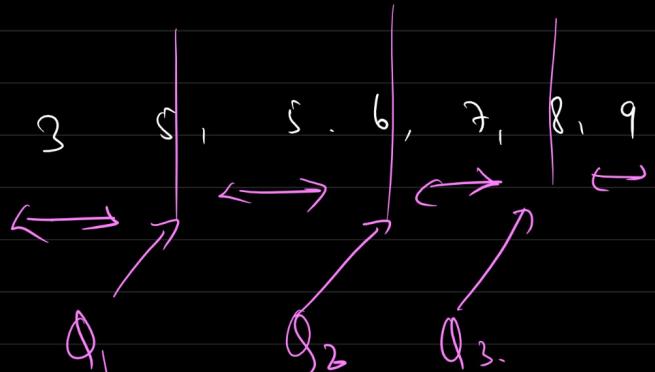
- * Quartiles are the values that divides a list of Nos into quarters
 - \hookrightarrow Put the no in order
 - \rightarrow Then cut the no into 4 equal parts
 - \rightarrow The quartiles are at the cut.

e.g. 6, 8, 5, 5, 7, 3, 9.

Order - 3, 5, 5, 6, 7, 8, 9

Cut the no. into quarters

$$\begin{aligned}Q_1 &= 5 \\Q_2 &= 6 \\Q_3 &= 8\end{aligned}$$



Ex. 1, 1, 1, 1, 2, 2, 3, 3, 3, 4
 total no = 11 Q₁
 if total no odd
 Q₂ Q₃

$$Q_1 - \frac{n+1}{4}^{\text{th}} = \frac{11+1}{4} = \frac{12}{4} = 3^{\text{rd}} \text{ no}$$

$$Q_1: \underline{1}$$

$$Q_3 = 3\left(\frac{n+1}{4}\right)^{\text{th}} = \frac{3(11+1)}{4}^{\text{th}} = \frac{3 \times 12}{4} = 9^{\text{th}} \text{ obs.}$$

$$Q_2 = \left(\frac{n+1}{2}\right)^{\text{th}} = \frac{11+1}{2}^{\text{th}} = 6^{\text{th}} \text{ no.} = \underline{2}$$

if no is even
 Q₁ = $\frac{n}{4}^{\text{th}}$ no
 Q₃ = $\frac{3n}{4}^{\text{th}}$ no
 $Q_2 = \frac{n}{2} + \left(\frac{n}{2} + 1\right)^{\text{th}}$

* Q₂ \Rightarrow Median

* Five Point Summary

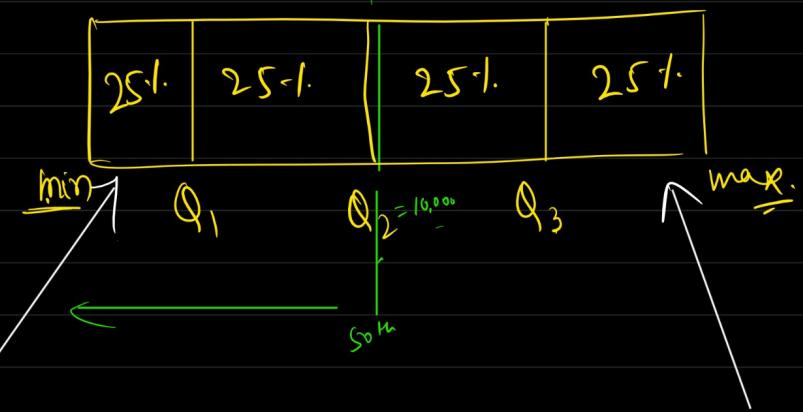
Q₀ (min) — 0th percentile
 Q₁ — 25th percentile
 Q₂ — 50th percentile
 Q₃ — 75th percentile
 Q₄ (maximum) — 100% percentile

transaction amount	
1000	
2000	
3000	
:	
!	

$$\min(Q_0) \rightarrow 1000$$

Q₁ (25th percent) — 5000 \Rightarrow 25% of transaction amount is equals to or below 5000 in the data.

Q₂ (50th percent) — 10,000 \Rightarrow 50% of transaction amount is equals to or below 10,000.



Inter quartile range

$$\underline{Q_3 - Q_1}$$

2, 3, 3, 3, 3, 3, 4, 4, 5, 5, 5, 6, 6, 7, 8, 99 outlier

$$Q_3 - Q_1$$

$$25^{\text{th}}\% = \frac{25}{100} \times 16 = \frac{1}{4} \times 16^{\text{th}} - 4^{\text{th}} \text{ no.}$$

$$Q_1 = 3.$$

$$Q_3 = 75^{\text{th}} \text{ perc} = \frac{75}{100} \times 16 = \frac{3}{4} \times 16 = 12^{\text{th}} \text{ no}$$

$$Q_3 = 6$$

$$\text{IQR} = 6 - 3 = \underline{\underline{3}}$$

$$Q_3 - Q_1$$

Outliers are extreme values.



$$\text{Lower fence} = Q_1 - 1.5 \text{ IQR}$$

$$\text{Upper fence} = Q_3 + 1.5 \text{ IQR}$$

2, 3, 3, 3, 3, 4, 4, 5, 5, 6, 6, 6, 7, 8, 9, 9



$$L.F = Q_1 - 1.5 \times IQR, U.F = Q_3 + 1.5 \times IQR$$

$$= 3 - 1.5 \times 3$$

$$= -1.5$$

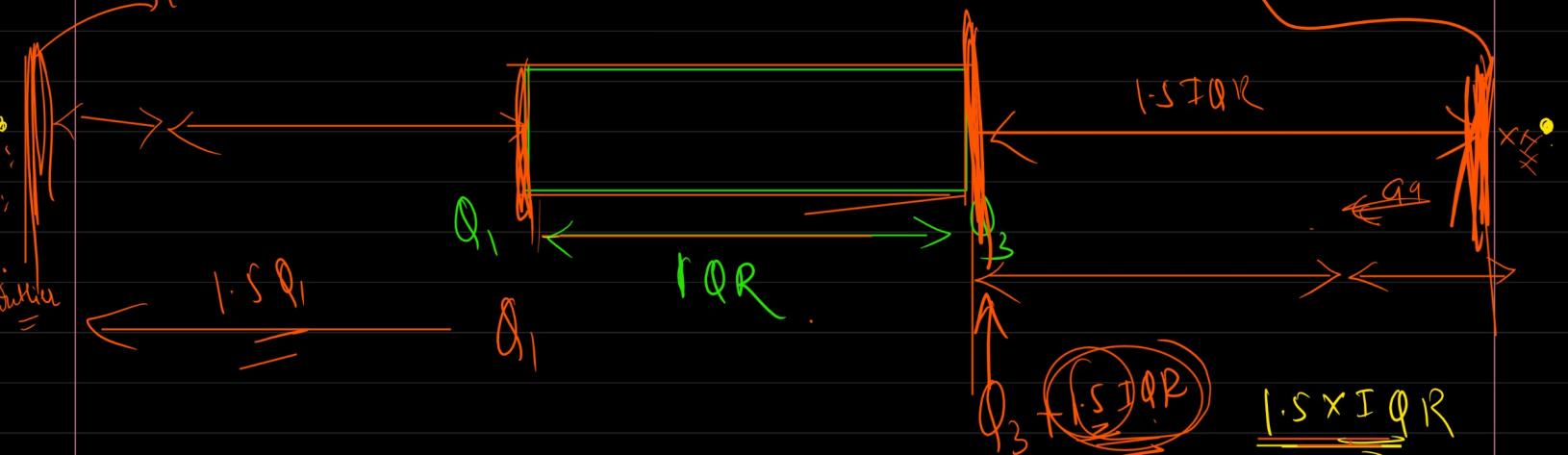
$$\downarrow$$

$$6 + 1.5 \times 3 = 4.5 + 6$$

$$= 10.5.$$

$$Q_1 = 3$$

$$Q_3 = 6$$



Quartiles \rightarrow Quarters \rightarrow 4 equal parts.

$$1.5 \times IQR$$

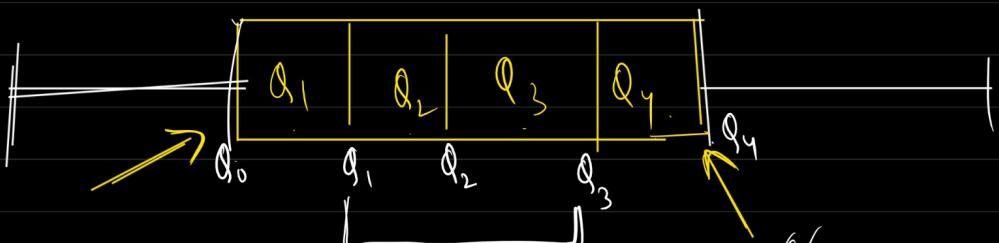
$$10 \times 1.5$$

$$\Rightarrow 15 \text{ m}$$

Measures of Spread | Dispersion

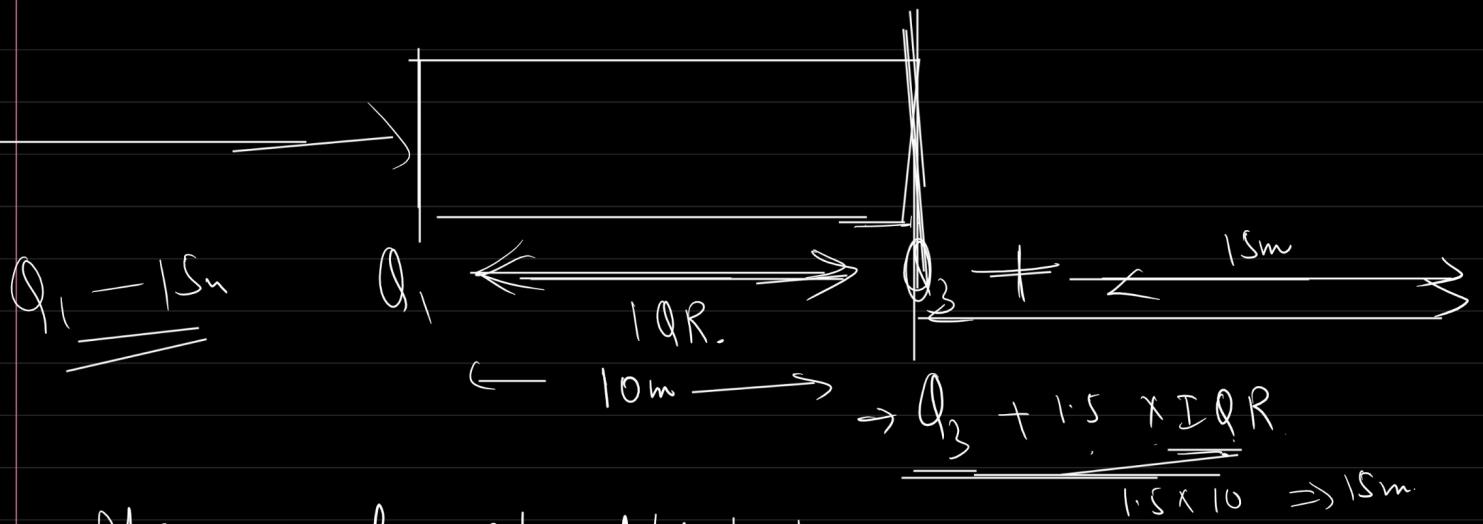
use - To know
outliers

Sns.boxplot()



$$Q_3 - Q_1 = IQR$$

$(Q_0 \text{ & } Q_4 \rightarrow \text{is max value})$



Measures of spread/dispersion

* Variance

* Standard deviation

* Mean deviation

$$\sum \frac{x - \bar{u}}{n}$$

$$\begin{array}{c}
 \text{Mean} = \bar{u} \\
 \hline
 \text{On an avg each of the data} \\
 \text{is } 1.2 \text{ units away from mean value}
 \end{array}$$

$$\text{Mean} = \bar{u}$$

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

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

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

Population variance

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

Sample Variance

* Variance :-

- Calculate mean
- for each no in data, subtract the mean and no
- Square the difference
- Calculate the avg of square of difference

→ data 1 2 3 3 4 4 (Sample)

x	\bar{x}	$x - \bar{x}$	$(x - \bar{x})^2$
1	2.83	-1.83	3.34
2	2.83	-0.83	0.68
3	2.83	0.17	0.03
3	2.83	0.17	0.03
4	2.83	1.17	1.37
4	2.83	1.17	1.37
<hr/>		<hr/>	
mean = 2.83		<hr/>	

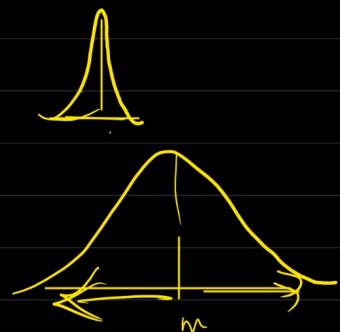
$$S^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

$$= \frac{6.82}{5} = 1.37$$

Population		Sample
mean :	μ	\bar{x}
Variance	σ^2	S^2
No of dp	N	n

Spread of data.
Variance

Variance \uparrow Spread \uparrow
Variance \downarrow Spread \downarrow



* Standard deviation = A measure of how spread the data is

Square root of Variance.
= $\sqrt{\text{Var.}}$

$$\text{Std Dev of Pop} = \sigma = \sqrt{\text{Var}_p}$$

$$\text{" " " " Samp} = S = \sqrt{\text{Var}_s}$$

$$\sqrt{\text{Var}} = \underline{\text{std}}$$

Variance \rightarrow How much spread.

Std dev \rightarrow Standard way of knowing where your data lies

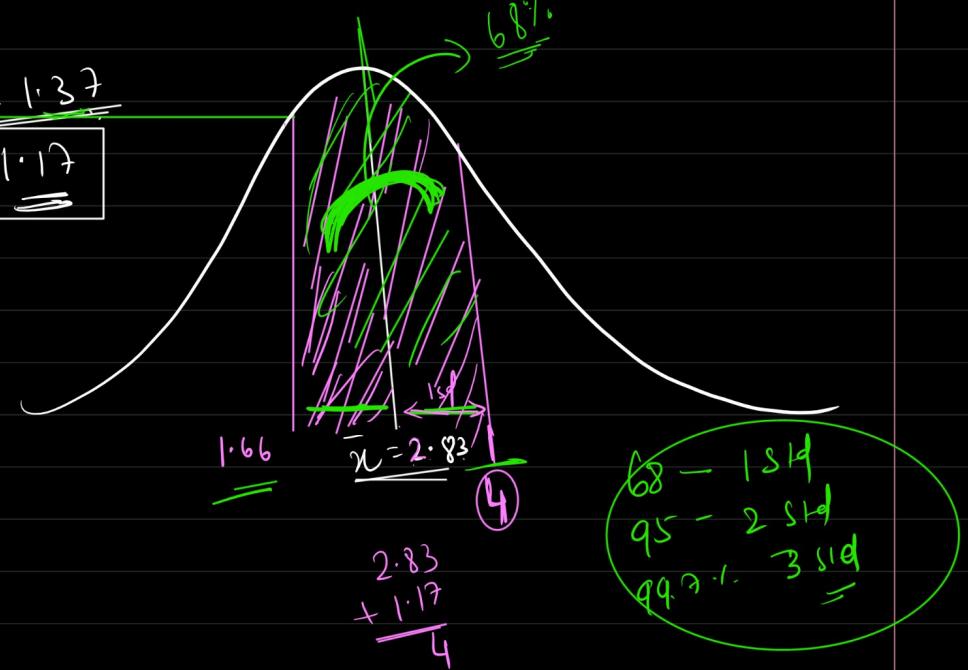
1, 2, 3, 3, 4, 4

$$\rightarrow \text{Var} = 1.37$$

$$\rightarrow \text{std} = \sqrt{1.17}$$

How much df is
 $\frac{1}{\sqrt{3}}$ of away
from mean

$$1.17 + 2.83$$



$$\text{Vana } \sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{N}$$

$+3$ \bar{x} -3 $(x - \bar{x})^2$

$$|x - \bar{x}|$$

* $\text{Var Sample} = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$ Bessel correction

* We are estimating population variance using sample.

✓ $(x-\mu)^2 > (x-\bar{x})^2$

$$\frac{(x-\mu)^2}{N} \approx \frac{(x-\bar{x})^2}{n-1}$$

$$(2) = \frac{10}{5} = \frac{8}{4}$$

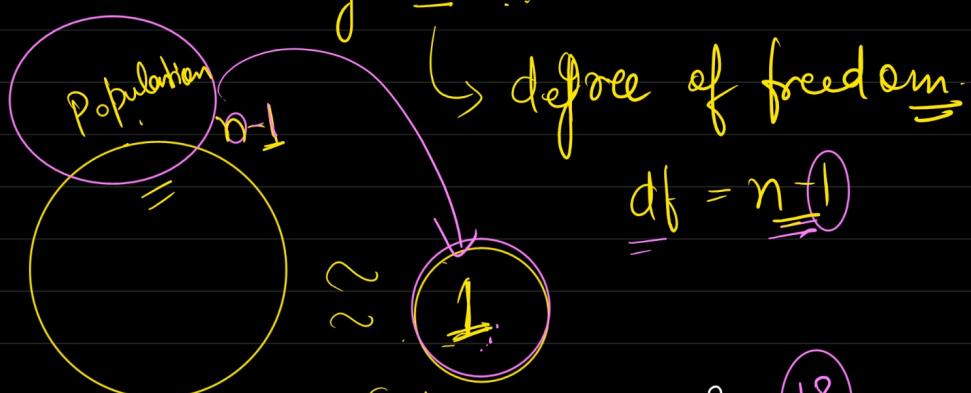


$n-1$ rather than n because sample variance will be unbiased estimator

$$\text{S}^2 = \frac{\sum (x - \bar{x})^2}{n-1}$$

----- x ----- x -----

Why $n-1$??



$$df = n-1$$

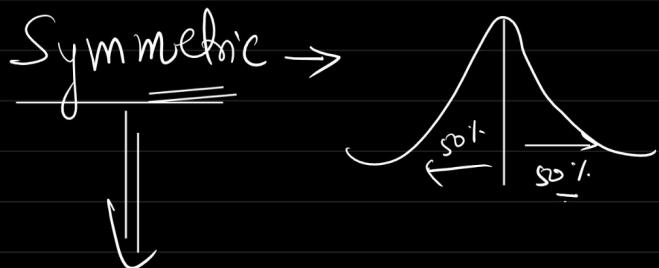
$$\frac{2}{2} + \frac{18}{18} = 10$$

$$n-1$$

* Measures of Symmetry

* Skewness

- ① MCT ✓
- ② MD ✓
- ③ M S



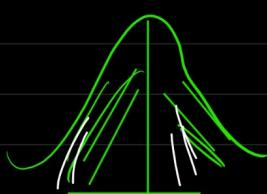
Skewness ⇒ Measure of data Symmetry.

$$\text{Skewness} = \frac{\bar{x} - \mu}{n\sigma^3}$$

No Skewness

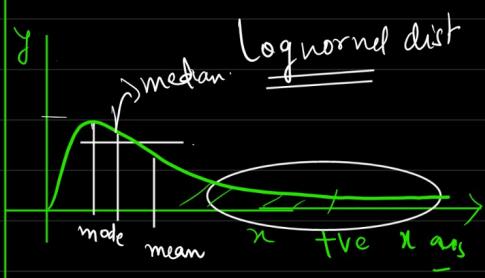
Skewness = 0

Symmetric



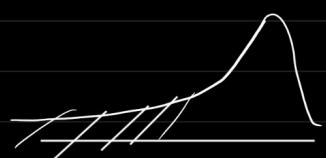
mean = median = mode

Positive Skewed

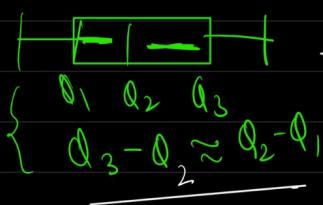


mean > median > mode

Negative Skewed.



mode > median > mean



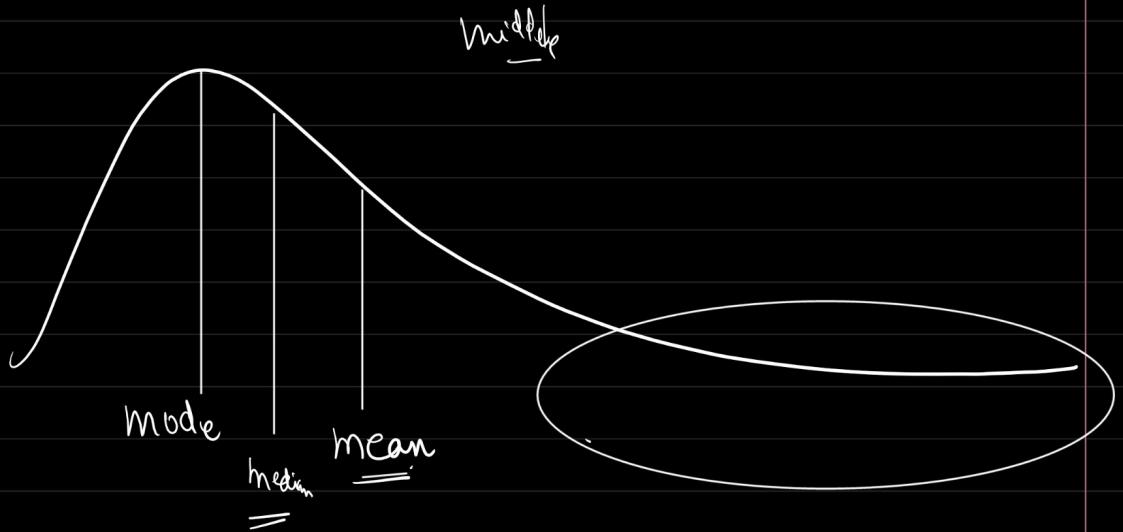
$$Q_3 - Q_2 > Q_2 - Q_1$$

Transformations

- log transform
- reciprocal
- Square
- Square root

Set

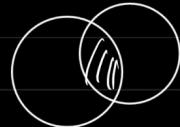
|| - Box Cox transform



x	reciprocal	log	Square	Square root	Box Cox
1	1	log(1)	1 ²	$\sqrt{1}$	=
2	$\frac{1}{2}$	log(2)	2 ²	$\sqrt{2}$	
3	$\frac{1}{3}$	log(3)	3 ²	$\sqrt{3}$	
4	$\frac{1}{4}$	log(4)	4 ²	$\sqrt{4}$	
5	$\frac{1}{5}$		5 ²		

Set

① Intersection $A \cap B$



② Union $A \cup B$



③ Difference $= A - B \rightarrow$ items which is present only in first set

$A - B$



④ Subset - All the elements of B present in A



⑤ Skewness

⑥ Symmetric differ →



* Covariance & Correlation

