

## Revision

DATE

PAGE

$\bar{X}$  (Arithmetic mean)

$$\bar{X} = \frac{\sum X}{n} \quad \xrightarrow{\text{estimate of}} \quad \mu$$

(Sample mean)  $n$  (population mean)

Median (middle number of series when ordered)

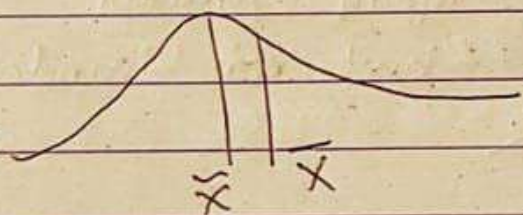
$n = \text{odd}$  (middle no. is series)

$n = \text{even}$  (mean of two middle numbers)

In symmetric distribution the mean is roughly equal to median.

But, In asymmetric distribution mean is greater than median.

$\text{mean} > \text{median}$  (due to outliers)



$$\bar{X} > \tilde{X} \text{ (median)}$$

## Mode

data having highest frequency

eg:- [28, 28, 5, 10]

mode = 28

[28, 5, 10, 15]

mode = ?

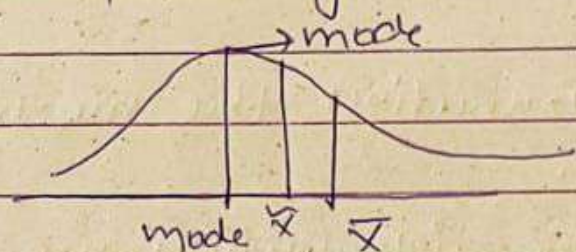
there is no significance of mode in small datasets.



mode signifies in large dataset like Census of country.

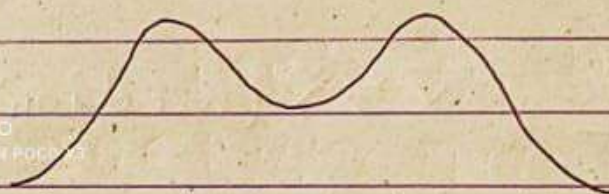
the three central tendency measure, mean, median, mode which one is good representative is totally depend upon context of problem (nature of dataset and problem).

9n positively skewed data.

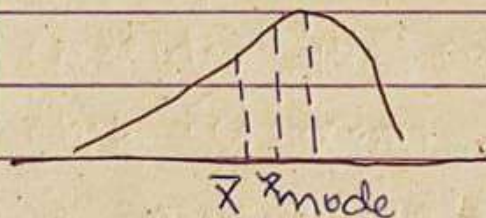


$$\text{mode} < \tilde{x} < \bar{x}$$

(median)



mean = median  
but ~~not~~ bimodal



negatively skewed data

$$\bar{x} < \tilde{x} < \text{mode}$$



## Quantiles

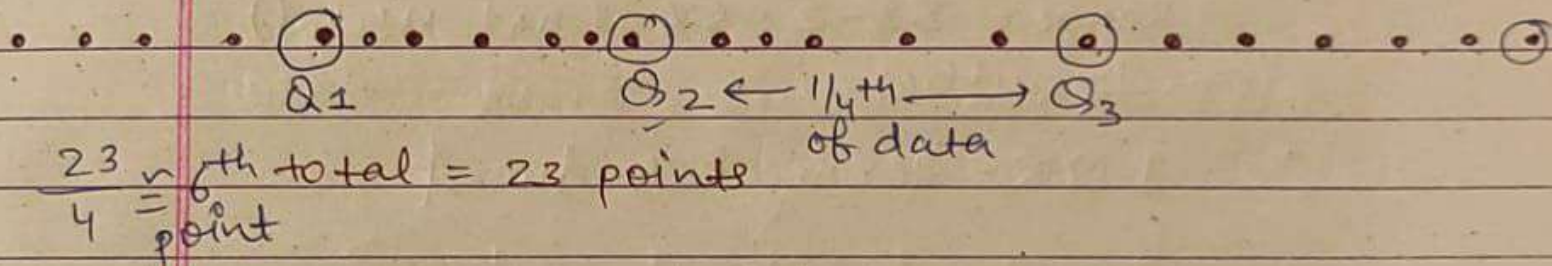
In a ordered dataset

How can you describe the useful locations on dataset.

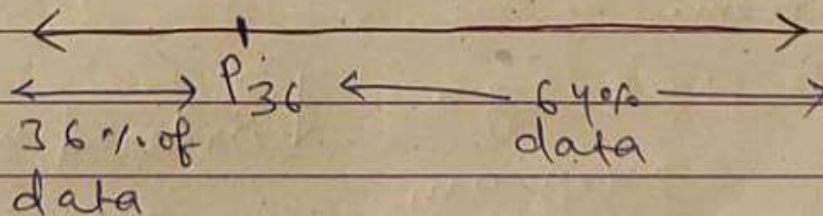
[4, 7, 8, ---, 178, 180, 189]

\* the median is the observation that is half way through the ordered dataset.

$Q_1$  (1st Quartile) :- Observation that is one through Quarter of the way through dataset.



percentile :- One percent of the way through the ordered dataset.





## Ranges

Consider an ordered dataset

$$\text{Range} = (\max - \min)$$

try to give us idea on the spreadness of data.

but it is more significant when dataset not containing any outlier.

but when data containing outlier there is no significance of Range.

2      2      5      6      9      10      58

$$\text{Range} = 58 - 2 = 56 \quad (\text{not useful})$$

58  $\rightarrow$  outlier

$\therefore$  IQR term introduced

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

IQR (Inter Quartile range)

Q. Consider the sugar content in one scoop of each icecream.

A      26.1 gm

B      24.9 gm

$\downarrow$  ordered way

—  
—  
2

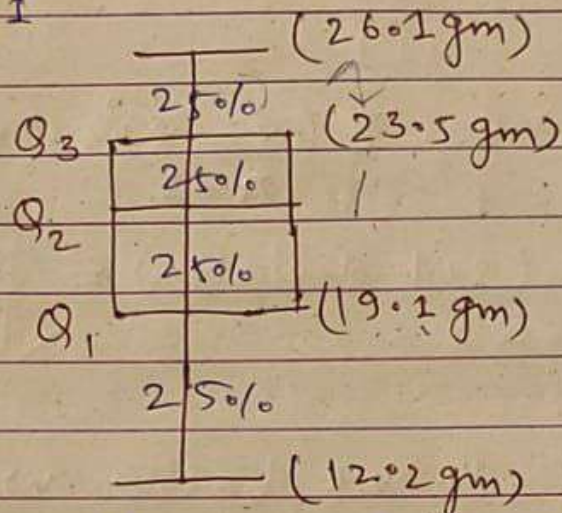
—  
—  
12.2 gm



$$\begin{aligned}\text{Range} &= 26.1 - 12.2 \\ &= 13.9 \text{ gm}\end{aligned}$$

(max<sup>m</sup> spreadness of data b/w point is 13.9 gm)

$$\begin{aligned}\text{IQR} &= Q_3 - Q_1 \\ &= 4.4 \text{ gm}\end{aligned}$$



give much more idea about spreadness.  
It gives % of data in specified range  
by drawing whisker plot.



→ variance and std deviation

consider a sample of population

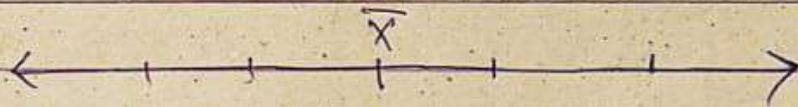
$$\bar{x} = \frac{\sum x}{n}$$

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

$$s \text{ (std-deviation)} = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$$

Q. why do we bother with variance?

Objective :- describe spread of data



Let's find avg. deviation from the mean

$$\sum (x - \bar{x}) = 0$$

for right points from  $\bar{x}$  = +ve deviation

for left point from  $\bar{x}$  = -ve deviation

If we take  $|x - \bar{x}|$  then this function is not differentiable.

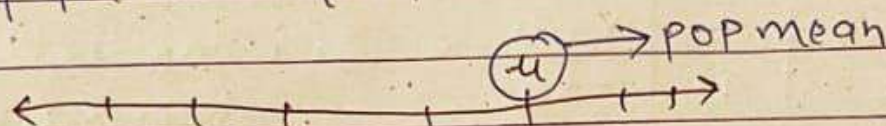
∴ Average square deviation from the mean:

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



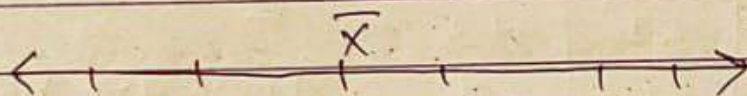
Q. why did we divide by  $n-1$ ?

the variance is the avg. square deviation from the population mean.



$$\text{variance} = \frac{\sum (X - \mu)^2}{N} \quad N = \text{no. of data point}$$

But we don't have population mean we estimate it using sample ( $\bar{X}$ )



$$S^2 = \frac{\sum (X - \bar{X})^2}{??}$$

\* the sample mean is one possible position for the true population mean.

\* At any other position, the sum of square would be larger.

$$S^2 = \frac{\sum (X - \bar{X})^2}{n-1} \quad (\text{smaller denominator adjust the variance upward})$$

\* but why  $n-1$  not other number this thing is empirically proven



## → Degree of freedom

$$\mu = 53$$

| Obs | x  | $x - \mu$ | 3 independent<br>observation<br><br>Dof = 3 |
|-----|----|-----------|---|
| 1   | 41 | -12       |   |
| 2   | 59 | +6        |   |
| 3   | 50 | -3        |   |

$$\bar{x} = 58$$

| Obs | x  | $x - \bar{x}$ |
|-----|----|---------------|
| 1   | 61 | +3            |
| 2   | 51 | -7            |
| 3   | -  | -             |

$$\text{Dof} = 2$$

these are bound to  
value to specify ~~the~~

$$\frac{\sum x}{n} = \bar{x} \text{ eqn.}$$

$$\sigma^2 = \frac{\sum (x - \mu)^2}{N} \quad \text{(pop variance)}$$

$$s^2 = \frac{\sum (x - \bar{x})^2}{n - 1} \quad \text{(sample variance)}$$



→ Coefficient of variation

$$CV = \frac{S}{\bar{X}}$$

$$X = [1, 2, 3] \quad \bar{X} = 2 \quad S_X = 1$$

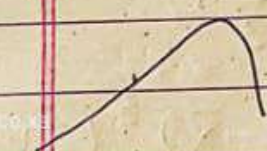
$$y = [101, 102, 103] \quad \bar{y} = 102 \quad S_y = 1$$

$$CV(X) = \frac{1}{2} = 0.5$$

$$CV(y) = \frac{1}{102} = 0.0098$$

Scaling the variance with respect to mean (dataset)  
On finding CV units don't matter for two different datasets because it gets scaled w.r.t to dataset.

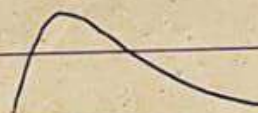
→ Skewness



(negatively skew)



(No skew)



(positive skew)

+ve skew = mode < median < mean

-ve skew = mean < median < mode

\* the greater the skew the greater distance between mode, median and mean.



## Pearson method

$$\text{mode skewness} = \frac{\text{mean} - \text{mode}}{\text{std. dev}}$$

$$\text{median skewness} = \frac{3(\text{mean} - \text{median})}{\text{std. dev}}$$

$$\text{mode} \approx 3\text{median} - 2\text{mean}$$

## Moment based calculation

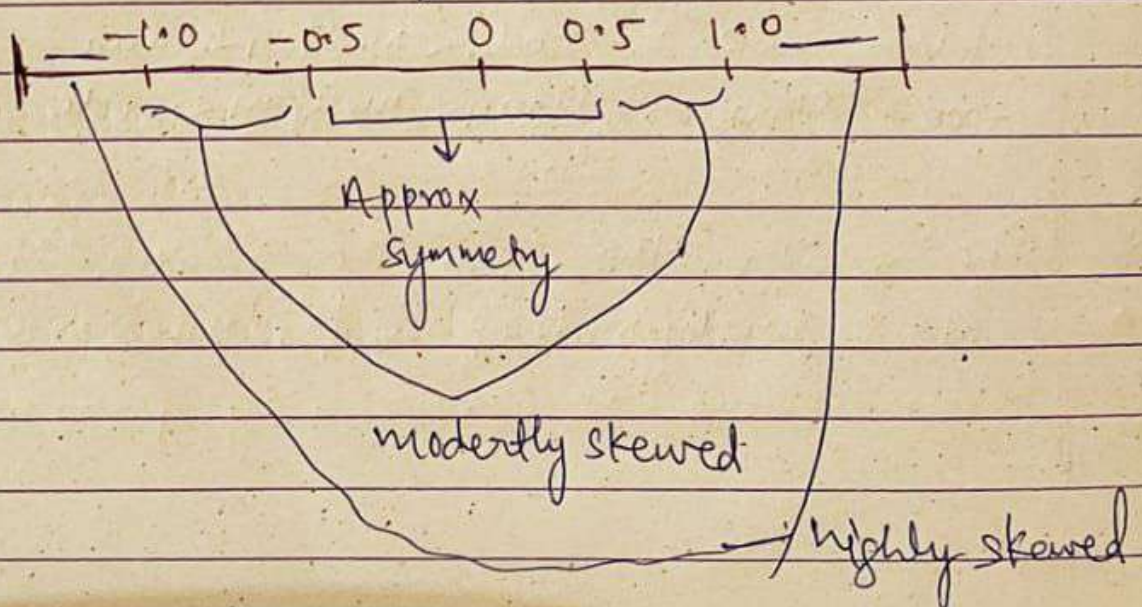
$$\text{First moment} = \sum \frac{x}{n}$$

$$\text{Second moment} = \sum \frac{x^2}{n} = \sum \frac{(x - \mu)^2}{n} \text{ (population)}$$

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

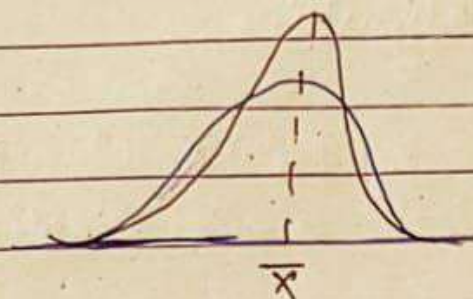
$$\text{Third moment} = \sum \frac{x^3}{n} = \frac{1}{n} \sum \frac{(x - \mu)^3}{\sigma^3}$$

$$\text{(Sample)} = \frac{n}{(n-1)(n-2)} \frac{\sum (x - \bar{x})^3}{s^3} \text{ (skew)}$$





## → Kurtosis



$$\sigma = 5$$

$$\text{Skew} = 0$$

$$\bar{x} = 0$$

the black curve is more peaked and has fatter tail.

$$\text{Fourth moments} = \frac{1}{n} \left( \frac{\sum (x - \mu)^4}{\sigma^4} \right)$$

$$= \frac{n(n+1)}{(n-1)(n-2)(n-3)} \left( \frac{\sum (x - \bar{x})^4}{s^4} \right) - \frac{3(n-1)^2}{(n-2)(n-3)}$$

A normal distribution has kurtosis of 3 called mesokurtic

Kurtosis  $> 3$  (leptokurtic)

Kurtosis  $< 3$  (platykurtic)

Kurtosis ranges from (1 to  $\infty$ )

$$(\text{Excess Kurtosis} = \text{Kurtosis} - 3)$$

(normal dist)

Excess Kurtosis ranges from  $-2$  to  $\infty$  and normal distribution is 0.

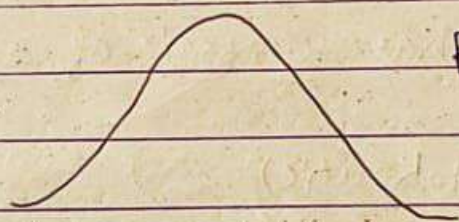


- \* Kaplanky (1945) discovered that there were certain anomalies to the higher peak = fatter tail relationship.

$$p(x) = \frac{1}{3\sqrt{\pi}} \left( \frac{9}{4} + x^4 \right) e^{-x^3} \quad (\text{plot the distribution})$$

→ "Kurtosis as peakdness" (incorrect notation)

- \* Kurtosis tells you virtually nothing about the shape of the peak - its only unambiguous interpretation is in terms of tail extremity, i.e. either existing outlier (for the sample Kurtosis) or propensity to produce outliers for the Kurtosis of probability distribution.



Fourth Standardised moment

$$\frac{1}{n} \left( \frac{\sum (x - \mu)^4}{\sigma^4} \right)$$

$$= \frac{1}{n} \left( \frac{\sum (x - \mu)^4}{\sum (x - \mu)^2 / n} \right)$$

Outliers contribute greatly to this summation.

"Kurtosis as tail extremity"



→ covariance and correlation  
Describes the relationship between two numerical variable.

"No covariance No correlation"

$$\text{COV}(x, y) = \sigma_{xy} = \frac{\sum (x - \bar{x})(y - \bar{y})}{n-1}$$

COV = +ve      positively related  
COV = -ve      negatively related

x point  $\bar{x}$  is positive deviation &  $\bar{y}$  point  $\bar{y}$  is positive deviation & then

$(x - \bar{x})(y - \bar{y}) \rightarrow +ve$       both are on +ve deviation

x point  $\bar{x}$  is negatively deviated &  $\bar{y}$  positively  $\bar{y}$  is.

$(x - \bar{x})(y - \bar{y}) \rightarrow -ve$

x  $\rightarrow$  neg      y  $\rightarrow$  pos

(-ve relation)

we can also see from var

$$\begin{aligned} s^2 &= \frac{\sum (x - \bar{x})^2}{n-1} = \frac{\sum (x - \bar{x})(x - \bar{x})}{n-1} \\ &= \frac{\sum (x - \bar{x})(y - \bar{y})}{n-1} \end{aligned}$$



covariance doesn't describe strength of relationship

$$\text{CORR}(x, y) = \rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

$$(-1 \leq \rho \leq 1)$$

$\text{corr}(x, y)$  tells about strength of relationship.

### → Standard error

the standard error of statistics (usually an estimate of a parameter) is the standard deviation of its sampling distribution.

If the statistic is the sample mean, it is called standard error of the mean.

The sampling distribution of the mean is generated by repeated sampling from the same population and recording of the sample mean obtained. This forms a distribution of different means and this distribution has its own mean and variance.

Mathematically,

the variance of the sampling mean distribution obtained is equal to the variance of the population divided by the sample size.

∴ Sample size  $\uparrow \rightarrow$  sample mean more around pop mean



Standard  
error of mean

$$= \frac{\sigma}{\sqrt{n}}$$

$\sigma$  = std. dev of  
population

in other word the standard error is a measure  
of sample means around the pop mean.