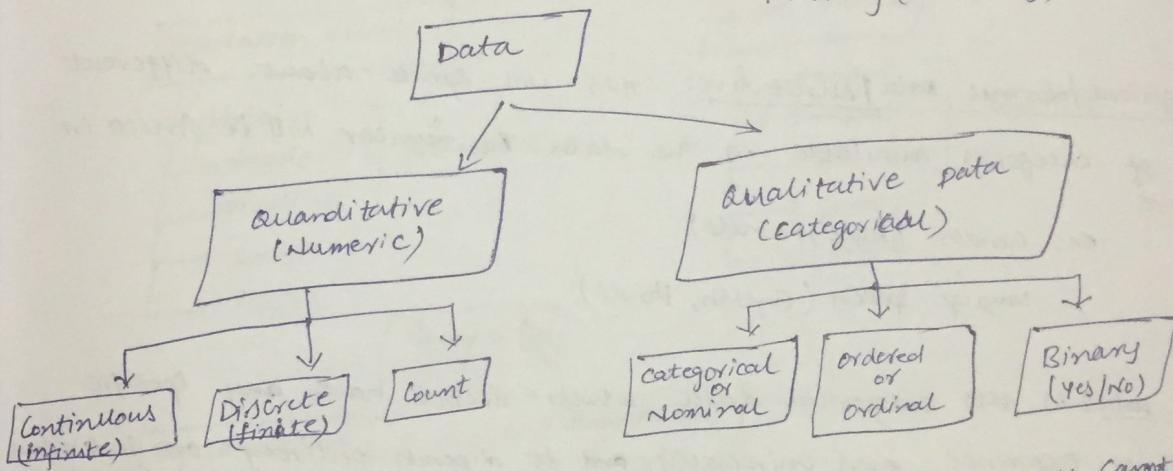


Statistics

Data: Collection of facts which are unorganised which is in raw format. This data has to be processed which will be meaningful.

Information: Collection of unorganised facts which is processed and arranged in organised way to give meaningful information.

(or)
Data that has been processed will have meaningful insight.
processing (cleaning, ETL)



Continuous: variable will have infinite values and this values cannot be

~~discrete~~: countable & usually represented in decimals.

Ex: height of students in a class

weight of students in a class

temperature in this month

Ex: → Age which is continuous if the measure of granularity is in minutes and seconds.

→ population which is continuous for every sec there is a new born or dead

→ Temperature which is continuous every hour it keeps on changing

However on circumstance this continuous variable can be discrete variable.

Ex: Dividing Students on basis of Age (categorical & discrete)

Discrete: variable will have there values as whole numbers. Such as counts of data.

Ex: No of lang a student can speak (1, 2, 3)

Political parties in a state (2, 3)

Religions in a Country (2, 3)

Different hair colors (2, 3)

Count Data: data that represents counts of data. (i.e discrete data)

categorical / Nominal data / Qualitative: this will speak about different types of categories available in the data. The number will be finite in nature.

Ex: Gender (Male, Female)

language spoken (English, Hindi)

Nominal data: is also categorical data which doesn't have any specific ranking or preference given universally and it depends entirely on individual interest. The represented data will not have any preference.

Ex: choose a color (Red, Green, Blue) you like

choose a Lang (Hindi, French, English) you like

which fruit you like (Mango, Apple, Strawberry)

choose a car (BMW, Benz, Ferrari) which you like.

Ordinal data: which will have ordered set of data.

How likely is your math teacher teaches.

Ex:

worst, bad, average, good, very good.

→ (it is ordered)

such as

Population & Sample

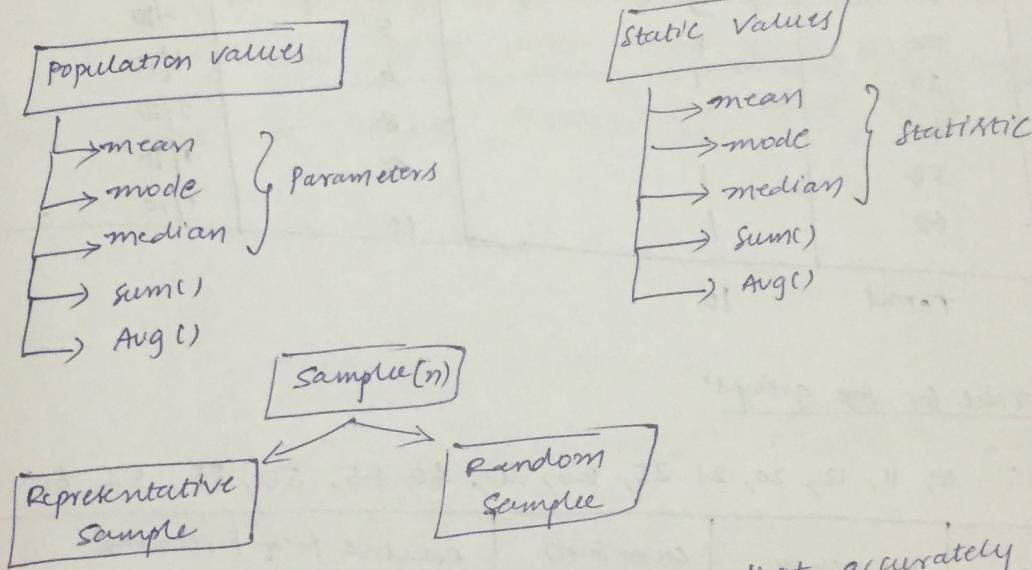
population (N): is the collection of items of interest to our experiment/study which is denoted by " N ".

Ex: All the college students (N)

Sample (n): is a sub set of population which is denoted by " n "

Ex: 10 students from each class (n)

$n = \text{Sample} = \text{size of the sample}$. = no of units taken in sample.
values refer to population are called "parameters" and values
refer to sample are called "statistic".



Representative Sample: is a subset of population that accurately reflects the actual population.

Random Sample: is picking the items in equal proportionate. This method is less likely used in real world. and much complicated to pick the sample in equal proportionate. (by) equal chance w.r.t population.

Frequency Distribution

discrete

continuous

frequency distribution: is a table/graph that displays the frequency of various outcomes in a sample.

Ex: 10, 10, 20, 50, 60; 10, 20, 30, 40, 40. (Total 10 data points)

frequency Table for above data set

Students	frequency/Count	relative cumulative sum	relative frequency Percent
10	3	3	3/10
20	2	5	2/10
30	1	6	1/10
40	2	8	2/10
50	1	9	1/10
60	1	10	1/10
Total	10		

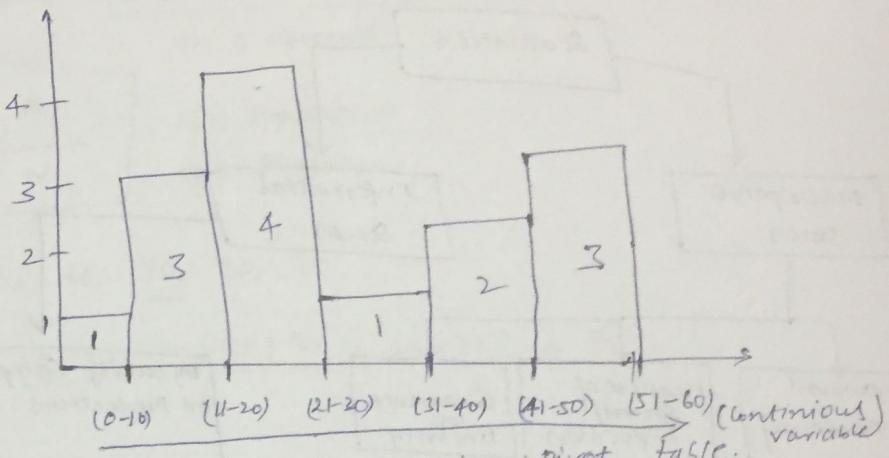
frequency Table for Age groups:

Ex: Age: 10, 11, 12, 20, 21, 25, 26, 30, 40, 45, 50, 55, 56, 60

Age	frequency	cum sum	relative freq. Percentage
0 - 10	1	1	1/14
11 - 20	3	4	3/14
21 - 30	4	8	4/14
31 - 40	1	9	1/14
41 - 50	2	11	2/14
51 - 60	3	14	3/14

Usually frequency tables are plotted in Histograms.

Histogram: Graphically plotting the distribution of items from the Sample / Population.



plot the same in Excel sheet using Pivot Table.

the frequency
Count
(cont.)

10 data points)

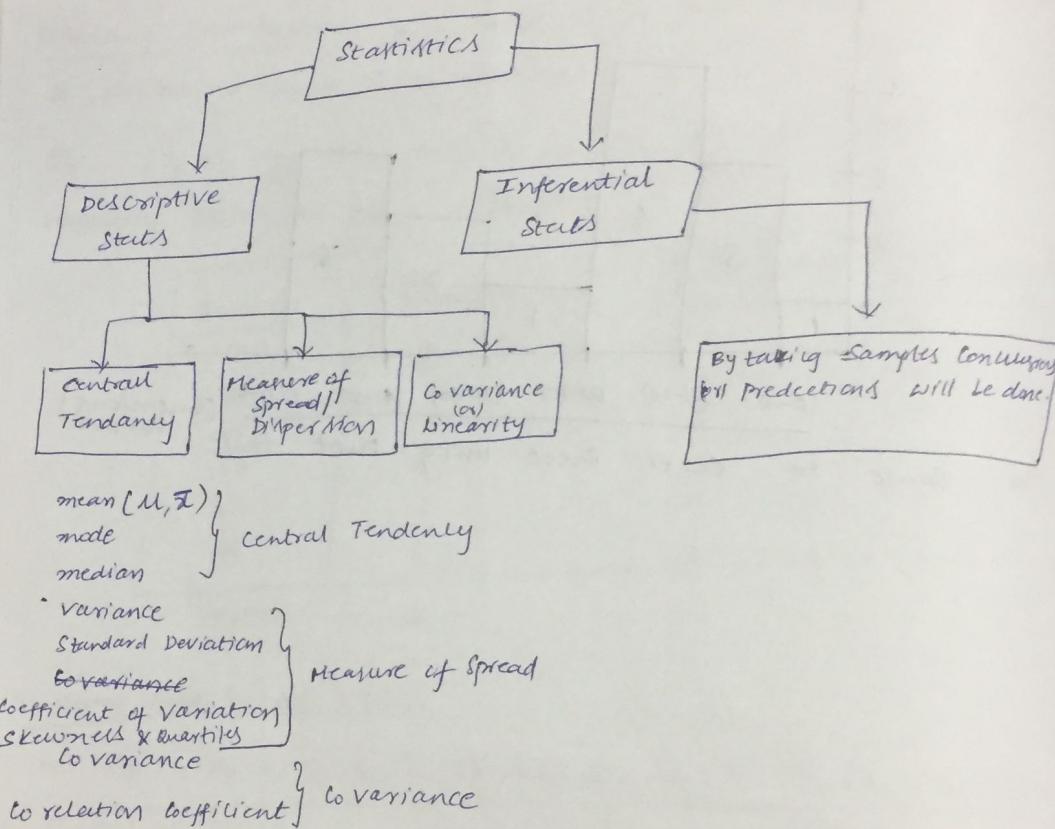
relative
frequency Percent

10
10
10
10
0
0

56, 60

total

Statistics: is a collection of data and analyse the data and interpret the data in graphs.



data

descriptive statistics: usually will have less population and hence calculation (or) stats will be applied for the entire data set.

Central Tendency: Measure of central value in a data set

can be measured with below formulae.

mean (\bar{M}) average $\frac{\text{sum of all the numbers}}{\text{no of elements}}$

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

N = Population

n = Sample

marks = 50, 60, 70, 80, 90

$$\text{Mean} = \frac{50+60+70+80+90}{5} = \frac{350}{5} = 70.$$

set = 5, 10, 100, 90, 500, 600, 700, 1400, 1500

set1 = 5, 10, 15, 20, 25, 30, 70, 80

mean = 31.875 (this mean doesn't make any sense out of the data set) this is because of data points 70, 80 which are called outliers.

outliers: will be far away from the central point (mean) hence mean value will be effected.

mean denoted by M (population) \bar{x} (sample)

* Mean is not a good measure in presence of outliers in such cases the measure of center can be calculated by median.

the mid point element of sorted data set will be median.

Median: Divide the dataset into two equal halves and the central element will be median.

The data set has to be sorted before making into equal halves.

Even data set $S_1 = 1, 2, \boxed{3, 4}, 5, 6$

$$\frac{3+4}{2} = 3.5 \text{ is the median.}$$

Odd data set $S_2 = \underline{8, 5, 4, \boxed{6}, 1, 3, 2}$

$$\begin{array}{l} \text{sort}(S_2) = 1, 2, \cancel{3, 4, 5, 6, 8} \\ \text{median} = 6. \times \end{array}$$

Sort the data set

$S_2 = \underline{1, 2, 3, 4, 5, 6, 8}$

$$\text{median} = 4 \checkmark$$

* Median is the positional measure, it doesn't matter for outliers.

* If mean is not near to median it indicates there are outliers.

Mode: The most frequently occurring data point will be mode.

$S_1 = 1, 1, 1, 2, 2, 2, 3, 3, 3, 3, 3$

mode = 3 (since '3' is the most repetitive element among other elements.)

* If mean, median & mode are near to each other then the data set is good for analysis.

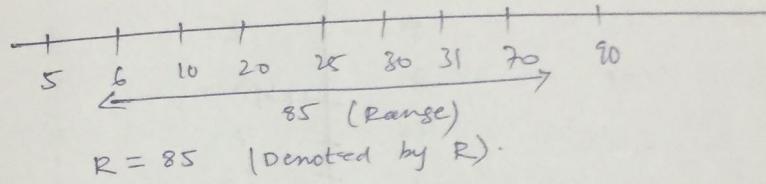
Measure of Dispersion / spread / variance

range: calculating the range $\max(\text{val}) - \min(\text{val})$

$$S = 40, 10, 20, 25, 70, 5, 6, 30, 90, 22, 31$$

$$\max(\text{val}) - \min(\text{val}) = 90 - 5 = 85$$

plot the data set:



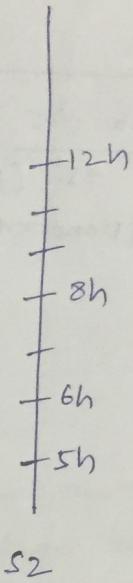
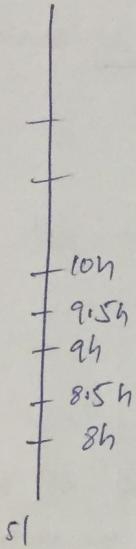
Variance: is used to calculate the distance of each data point from the mean.

Assume we have two working hours data sets.

$S_1 = 8.5h, 10h, 8h, 9h, 9.5h$ (Working hours of weekly).

$S_2 = 12h, 5h, 14h, 6h, 8h$

$$\bar{x}(S_1) = 9 \quad \bar{x}(S_2) = 9$$



Both the means of above data sets are same but there is a variance/dispersion in the data set 2.

Dataset 1 values are near to its mean value

Dataset 2 values are dispersed to its mean value

To identify the variance between the data sets we will calculate variance (σ^2) population s^2 (sample).

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

$(n-1)$ here ~~use~~ n .

a point from

weekly).

\underline{x}	\bar{x}	$x - \bar{x}$	$\frac{(x-\bar{x})^2}{n}$	$\frac{\sum x - \bar{x}}{n}$
8.5	9	-0.5	0.25	
10	9	1	1	
8	9	-1	1	0.5
9	9	0	0	
9.5	9	0.5	0.25	

$$\sigma^2(s_1) = 0.5$$

\underline{x}	\bar{x}	$x - \bar{x}$	$(x-\bar{x})^2$	$\sum(x-\bar{x})^2$
12	9	3	9	
5	9	-4	16	
14	9	5	25	12
6	9	-3	9	
8	9	-1	1	

$$\sigma^2(s_2) = 12$$

the mean = 9 for both data sets but using variance
we were able to identify the variance.

Note: Since variance is calculated by square of numbers
the variance value will be less proportionate with the
actual data set points hence standard deviation is used.
Interpreting the value of variance with mean is difficult hence Standard Deviation.

Standard Deviation (σ): Is the square root of variance

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

which is used to calculate how far the data points
are spread out from the mean point.

- * Standard deviation close to "Zero" it indicates the data points are tightly coupled/near to mean.
- * If the SD is away from "Zero"/high it indicates the data points are away from mean.

Ex: profit Percentage of Companies of different months.

$$A = 43, 44, 0, 25, 20, 35, -8, 13, -10, -8, 32, 11, -8, 21, 15 \quad \text{Mean} = 15$$

$$B = 17, 15, 12, 17, 15, 18, 12, 15, 12, 13, 18, 18, 14, 14, 15 \quad \text{Mean} = 15$$

Data set 'A' data points are dispersed away from mean = 15

Data set 'B' data points are near to its mean = 15

* In this kind of ^{tricky} situation where both means are same and thus can be identified by standard deviation

$$SD(A) = 19.47 \quad (\text{Dispersion is away from mean})$$

$$SD(B) = 2.2 \quad (\text{Dispersion is near to mean since SD is near zero})$$

* From this we can identify Company 'A' profits are zero & -ve in few months and it is not constantly performing.

Company 'B' is constantly performing profits and thus is proved by SD.

Coefficient
to its mean
→ this notion

CV is used

Mean
SD
CV

In order
it doesn't
So here

Coefficient of Variation: is the percentage of standard deviation to its mean.

→ This nothing but how much % of $SD(S)$ w.r.t mean (\bar{x})

$$CV = \frac{S}{\bar{x}} \times 100$$

CV is used to compare ~~the~~ two or more SD's in meaningful way.

	Test 1	Test 2
Mean	50	100
SD	10	15
CV	20.0	15.0

In order to compare the variability of SD b/w two data sets it doesn't make sense since means are different and far away. So here comes CV to compare the variance of SD's.

Skewness: is the measure of asymmetry which indicates

whether the data points are concentrated on which side.

* Skewness will help in identifying the shape of distribution.

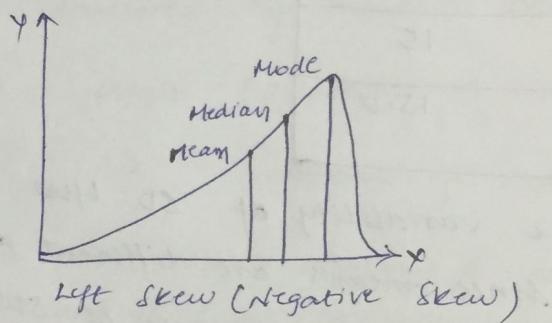
Right skew / Positive skew

Left skew / Negative skew

No / Zero skew

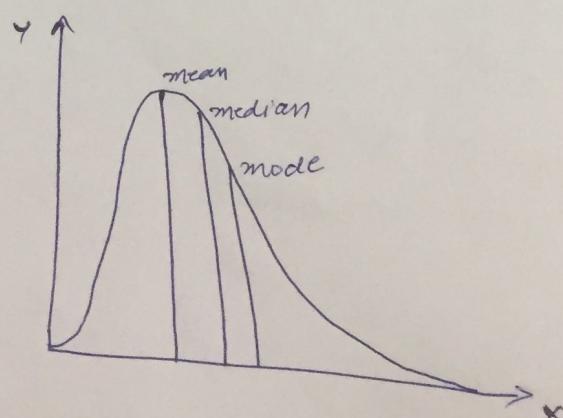
Left Skew: In Left skew the outliers are away from the y-axis.

$$\text{mean} < \text{median} < \text{mode}$$



Right Skew: In Right skew the outliers/distribution is towards the y-axis.

$$\text{mean} > \text{median} > \text{mode}$$



indicates
side.
distribution.

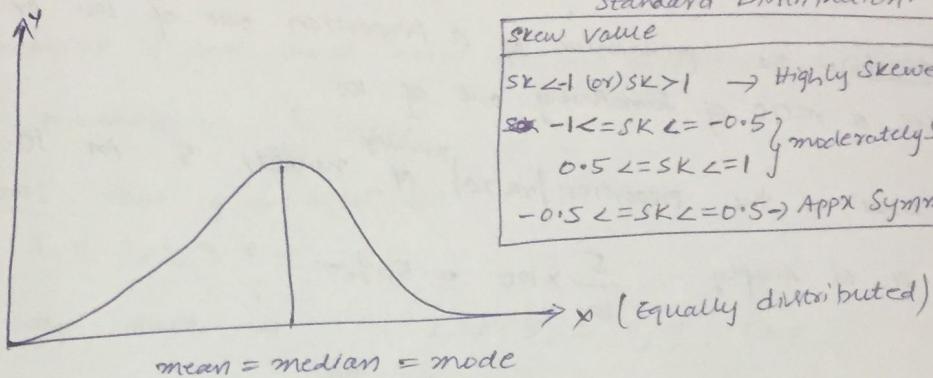
from the y-axis.

No-Skew: If there is zero skew then the distribution is a normal distribution curve (Bell curve)

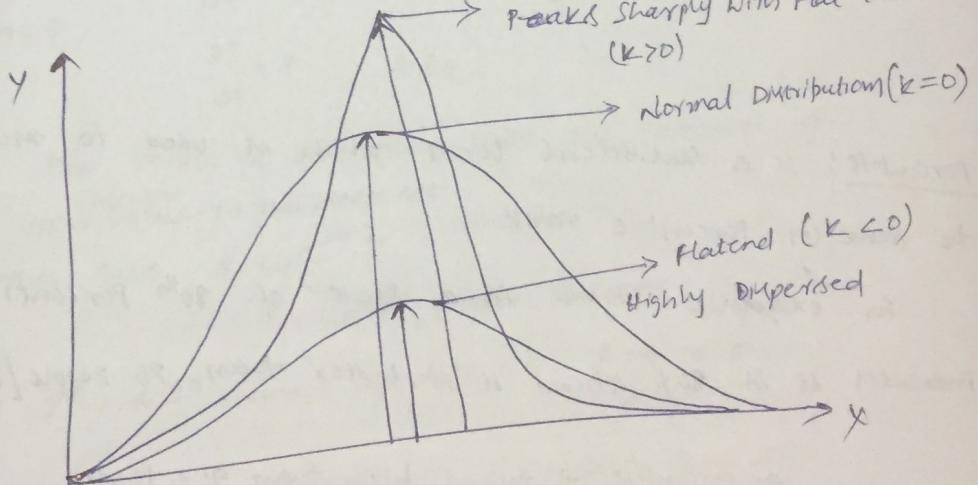
$$\boxed{\text{mean} = \text{mode} = \text{median}}$$

Skewness Table w.r.t Standard Distribution.

Skew Value
$SK < -1$ (or) $SK > 1$ \rightarrow Highly Skewed
$-1 < SK < -0.5$
$0.5 < SK < 1$
$-0.5 < SK < 0.5 \rightarrow$ Appx Symmetric



Kurtosis: is used to measure whether the curve is normal, flat or peaked



Kurtosis is measured with ' K '.

- If the value of ' K ' is near to "Zero" it is normally distributed
- If the value of ' K ' is -ve then the curve is flat.
- If the value of ' K ' is +ve or $K > 0$ then curve is peak.

Covariance & Correlation Coefficient

Percentage & Percentile:

Percentage (%) is the representation of a proportion out of 100. or it depicts a ratio of something out of 100.

Ex: What is the proportion/ratio^{percentage} of number 5 in 10.

It is simply $\frac{5}{10} \times 100 = 50\%$.

What is the proportion of 4 in 80

$$\frac{4}{80} \times 100 = 5\%$$

Percentile: is a statistical term which is used to measure the rank (or) percentile rank.

For example there is a score of 80th percentile it indicates it is 80% ahead it is better than 80 people/observations.

90 Percentile \rightarrow means better than 90% people.

By using percentile rank can be determined.

$$90 \text{ Percentile} = 10 \text{ rank}$$

$$95 \text{ Percentile} = 5 \text{ rank}$$

$$100 \text{ Percentile} = 1 \text{ rank}$$

Percentile is used to calculate/identify at which percentage the data point is standing when compared to other data points.

Definition: A percentile is defined as a point below which certain percentage of observations lie.

Percentile is calculated based on all observations.

10.

$$\text{Percentile} = \frac{P}{100} (n+1)$$

Calculation: Score of 8 students.

$$n = 5, 4, 2, 6, 9, 8, 1$$

① Sort the data $\rightarrow 1, 2, 4, 5, 6, 8, 9 \quad (n=7)$

② I want to calculate $\rightarrow 33$ percentile

③ $P = 33 \quad \frac{P}{100} (n+1)$

$$n = 7$$

$$\frac{33}{100} \times 8 = 2.64$$

Case 1: If the result is decimal take the average of lower and upper integer value in the data set.

In this case $2.64 \rightarrow 2$ (lower Integer) $\rightarrow 3$ (upper Integer)

Avg of 2nd and 3rd element of data set $\frac{2+4}{2} = 3$

Take the data in excel sheet and calculate percentile.

Quartiles: divide the data in to four segments according to where the number falls.

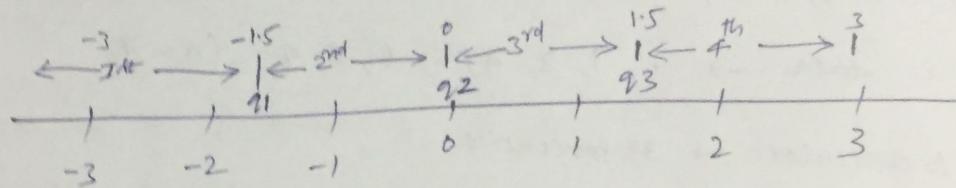
Quartiles are strictly based on Percentiles.

Quartiles are basically 3 types

(Lower Quartile) quartile-1 (q_1): is the 25^{th} percentile data.

(Second Quartile) quartile-2 (q_2): is the 50^{th} Percentile data

(Upper Quartile) quartile-3 (q_3) is the 75^{th} Percentile data.



q_1 = least 25% of data.

q_2 = least 50% of data

q_3 = least 75% of data / Highest 25% of data

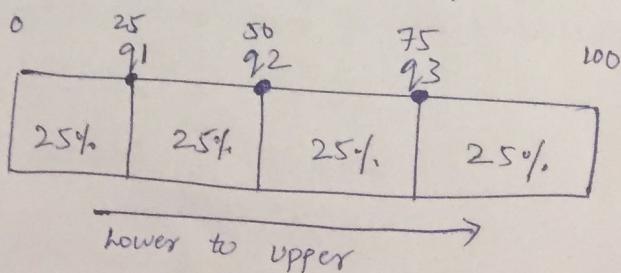
~~q_4~~ = highest 25% of data

q_1 is the median of lower half of data set.

q_2 is the median of whole data set

q_3 is the median of upper half of data set.

Interquartile: Upper quartile - Lower quartile.



$$\text{Interquartile} = q_3 - q_1$$

ents according into
Percentiles.

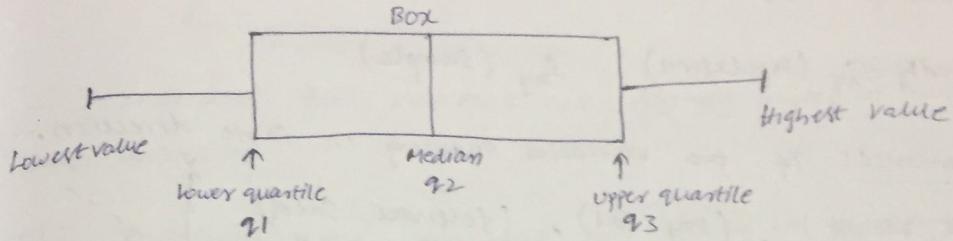
ntile data.

entile data.

centile data.

→
1
—
3

Box and Whisker plots: quartiles are basically represented graphically using box & whisker plots.



Find the quartile in excel sheet using quartile(A1:A10:1)

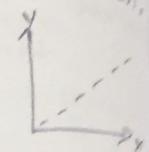
Covariance & Correlation Coefficient

Covariance: is the measure of relationship between two random variables. i.e. measure of linearity b/w variables.

Denoted by σ_{xy} (Population) say (sample)

Positive Covariance: If two variables moving in same direction.

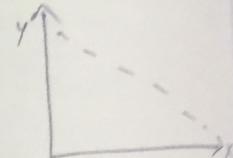
Ex: (height, weight), (x , y), (exp, sal), (festival, sales)



Negative Covariance: variables moving in opposite direction.

Ex: (no of claims of Insurance Company, profits)

(Age, hair color), (Age, fitness)



Zero Covariance: If two variables are independent of each other (Or) there is no relation b/w variables then zero covariance.

Ex: (petrol price, water bottle price)



Dis Advantage of Covariance: Interpretation of covariance is difficult.

Ex: +ve covariance can be (0.005, 5, 10, 33, 44, 555)

With the above values we wouldn't be able to interpret.