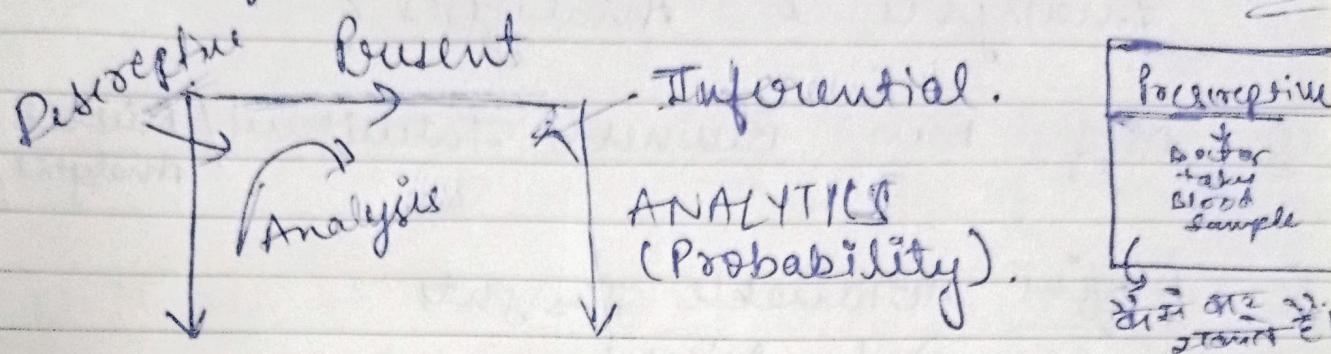


# ① Inferential & Descriptive Statistics



Variable → which tells characteristics of a particular item.

Q what is Statistics? Different parts of Statistics.

Q how we can use Stats for Business Model?

Q How ML & Deep Learning is used?

## Types of Variables

### Qualitative

(Categorical)

↓  
Object / factor

### Quantitative

Numerical

Integers

### Discrete

### Continuous

(float or decimal)

- Q what is the difference between Analysis & Analytics?
- Q Diff b/w Business Intelligence / Business Analytics?
- Q Diff Actionable Insights  
Data Science  
Machine Learning  
Central Tendency

Q Different type of Data types?

Ans → Statistics simply means numerical data, or we can say, a field of math that deals with collection of data, tabulation, and interpretation of numerical data.

### Types of Statistics

↓  
Descriptive Stats

↓  
Inferential Stats

↓  
Measure of central Tendency

↓  
Measure of variability

→ Mean  
→ Mode  
→ Median

→ Range  
→ Variance  
→ Dispersion

(i) Mean → Average of all values in Sample Set.

$$\text{Mean } (\bar{x}) = \frac{\text{Sum of all terms}}{\text{Total no. of terms}}$$

(ii) Median → Data set is ordered from lowest to highest and then find exact middle.

$$\begin{aligned}\text{Median} &= a \ b \ c \ d \\ &= \frac{b+c}{2}\end{aligned}$$

(iii) Mode → The value repeated most of time in central set.

$$\begin{aligned}\text{Mode} &= a \ b \ c \ a \ f \ a \\ &= a\end{aligned}$$

(iv) Range → Given measure of how to spread apart values in sample set or data set.

$$\text{Range} = \text{Max} - \text{min.}$$

(v) Variance → Describes how much a no. deviates from exp. value, also computed as square of deviation.

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

$n \rightarrow$  total data points  
 $\bar{x} \rightarrow$  mean of "  
 $x_i \rightarrow$  individual data points

(vi) Dispersion → Measure of dispersion of set of data from its mean.

$$\sigma = \sqrt{(1/n) \sum_{i=1}^n (x_i - \mu)^2}$$

\* inferential statistics :-

- It makes inference & prediction about population based on a sample data taken from population.
- It generalizes a large dataset & applies probabilities to draw a conclusion.
- Type :-

- ① one sample test of difference
- ② confidence Interval
- ③ Contingency table & Chi-Square stat.
- ④ T-test or Anova
- ⑤ Pearson Correlation.
- ⑥ Bi-variate Regression
- ⑦ Multi-variate Regression

Auf 2

Mistakes in Business can cost millions of dollars, Statistics can explain why a past event happened and predict what's likely to happen in future.

Ans 4

ML can be used as :-

- (1) Real-time chatbot agents.
- (2) Decision Support.
- (3) Image Recognition
- (4) Speech Recognition
- (5) NLP
- (6) Sentiment analysis.
- (7) Dynamic pricing tactics
- (8) fraud detection

Ans 5

Analysis

(1) Collection → Manipulation  
→ Examination of  
data for deep insights

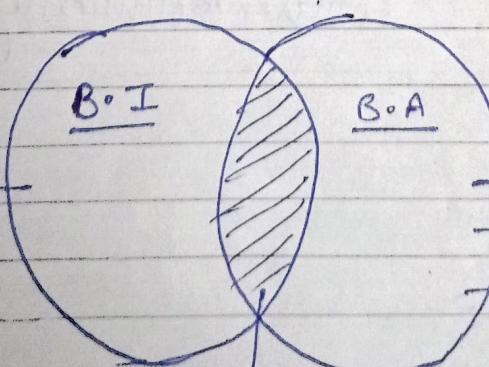
→ OpenRefine, Rapid  
Miner, NodeXL, Tableau  
etc.

Analytics

(1) Taking Analyzed  
data & working  
on it in a  
Meaningful way.

→ Python, SAS, Tableau,  
Excel etc.

Ans 6



Descriptive Analysis

- Diagnostic Analytics
- Predictive Analytics
- Prescriptive Analytics

- Collect, Analyze & visualize data
- Identify Pain Points
- Generate reports .

AJF

- Actionable Insights → Conclusions drawn from data that can be turned directly into an action or response.
- Central Tendency → the statistical measure that identifies a single value as representative of an entire distribution

Median → Most imp. component,

Skewness → left or right standard deviation

$\downarrow$   
-ve skew      +ve skew

- row → observations

- column → variables / parameters

Outlier → Which are exceptional with other data sets or values.

- ○ Skewed → Normal Distribution  
(Posion Distribution)  
(symetric distribution)
- Correlation → Multi collinearity (other name).

Count
mean
std
min
25%
50%
75%
max

To revise

- \* Outliers  $\rightarrow$  formulae.
- \* Quartiles & Interquartile
- \* Correlation & covariance  $\rightarrow$  Difference.
- \* Boxplot.
- \* Scatterplot

Codes :-

- ① mt['mpg'].var()
- ② mt.head()  $\rightarrow$  first five observations.
- ③ mt.tail()  $\rightarrow$  Last 5 observations.
- ④ mt['yr'].count()  $\rightarrow$  Total.
- ⑤ mt['cyl'].value\_counts()  $\rightarrow$  no. of counts.

(2)

8	4
4	11
6	7

## ⑥ Type of Analysis :-

**Univariate**  
Analysis of single variable

- # Best representation
- barplot
  - boxplot
  - histogram.

**Bivariate**  
Analysis of 2 var.

- Histogram / Bar chart
  - Line Graph
  - Scatterplot
- # Best Statistical Method
- correlation
  - chi square test

**Multivariate**

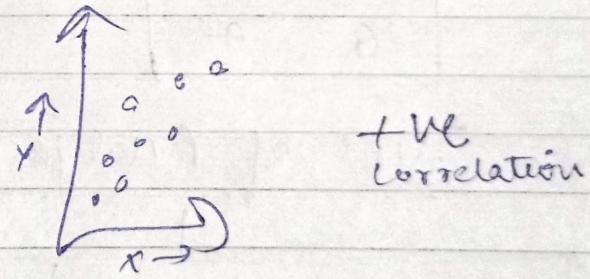
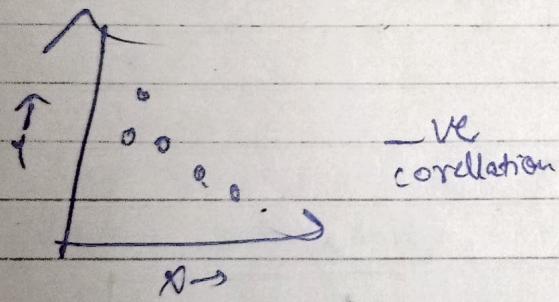
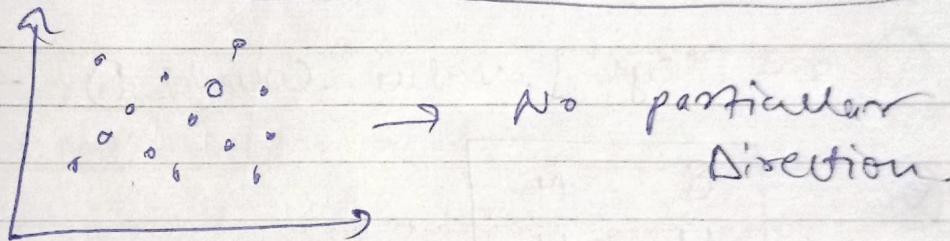
- Scatterplot matrix
- hyperbox
- Trellis Display
- Parallel coordinate
- Starplots

## \* Descriptive Statistics

- Scatter plot
- Pearson's correlation coefficient
- Simple linear.

relationship b/w 2 variables,  
 +Correlation  $\rightarrow (-1, +1) \rightarrow$  close to  $(1, 1)$   
 are considered  
 (Pecision Matrix) as strong.

$\pm 0.4$  se needed we consider No-relation,



## \* Pearson's Coefficient.

## \* Simple Linear Regression :-

$$Y = \alpha + \beta X + e$$

or

$$Y = mx + c$$

dependent variable  $\downarrow$  coefficient  
 $\downarrow$  slope  
 $\downarrow$  constant.

## 8 Simple Linear Regression.

### Multivariate:

But  
Visualization

# mt['cyl'].value.counts().plot.bar().

# mt['cyl'].plot.box().

Skewness: It is a measure of the asymmetry of a distribution.

can have :-

① Right (or +ve) Skewness

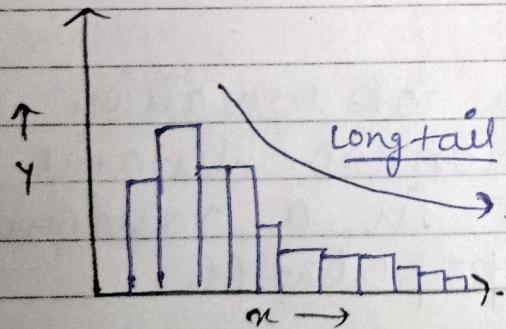
② Left (-ve) Skewness

③ zero Skewness

$\downarrow$   
mean = median

→ Normal distribution  
or  
poisson distribution

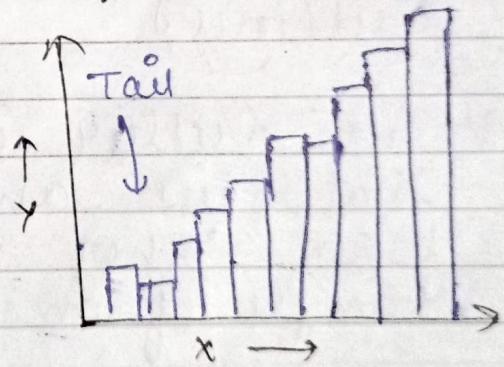
Right skewed (+ve)



mean > median

→ Longer on the right side of its peak.

Left skewed (-ve)



mean < median

→ Longer on the left side of its peak.

Pearson's Median Skewness

$$= 3 \times \frac{(\text{Mean} - \text{Median})}{\text{Stand. deviation}}$$

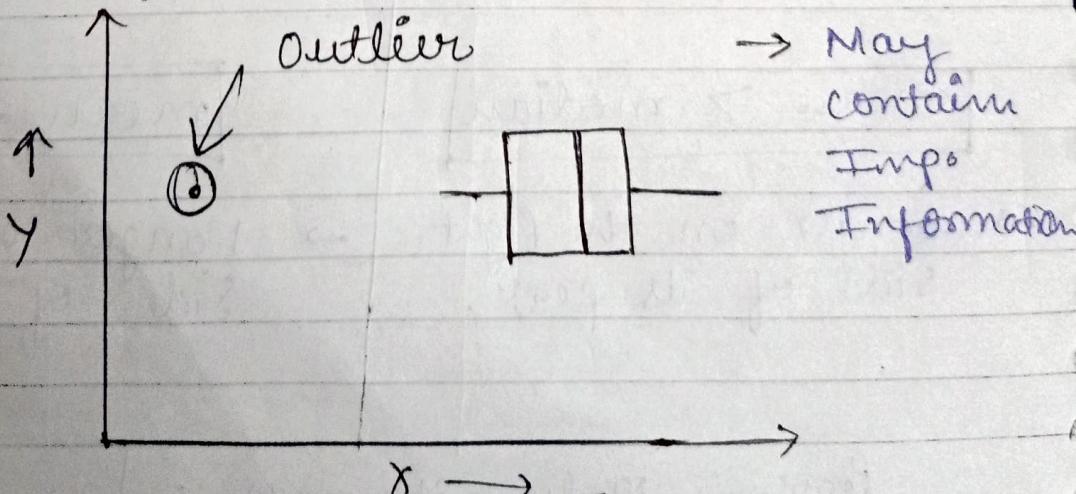
Transformation based on type of skew,

Type of Skew	Intensity of skew	Transformation
Right	Mild Moderate Strong very strong	Do not transform Square root Natural log Log base 10
Left	Mild Moderate Strong Very strong	Do not transform Reflect * Sq root Reflect * Nat log Reflect * Log <sub>10</sub>

Reflect means, take largest observation k  
subtract each observation from ( $k+1$ )

## ★ Outliers

→ An outlier is an observation that lies an abnormal distance from other values in a random sample from a population.



Boxplot representation

## Correlation

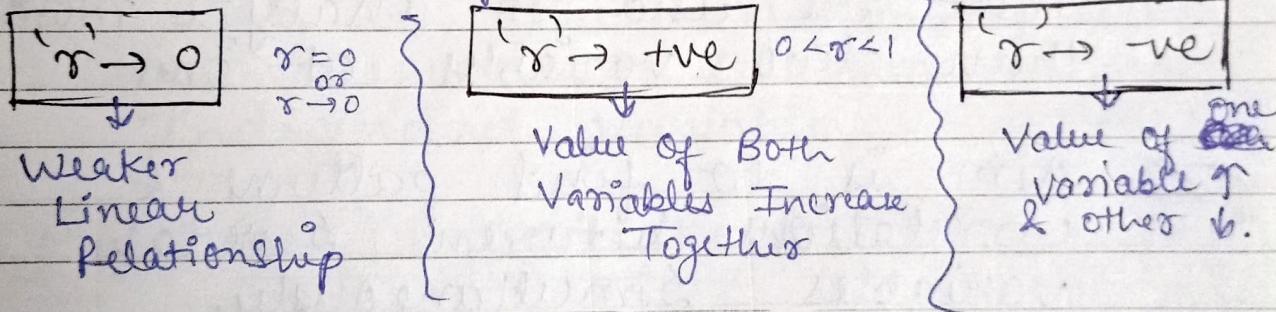
$A \propto a$  Diagonal Matrix  
 $(-1, +1)$

- It is a statistical measure that expresses the extent to which two variables are linearly related.
- Either they are directly proportional ( $A \propto B$ ) or indirectly proportional ( $A \propto \frac{1}{B}$ ) or Maybe they don't have any relationship.

Q How is it measured?

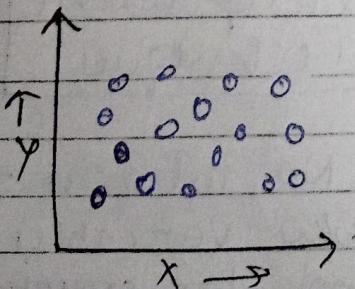
The Sample Correlation coefficient,  $r$ , quantifies the strength of the relationship.

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

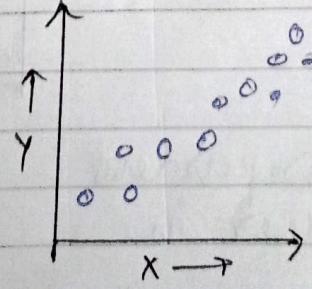


Code :  $\rightarrow \text{mt. corr}()$

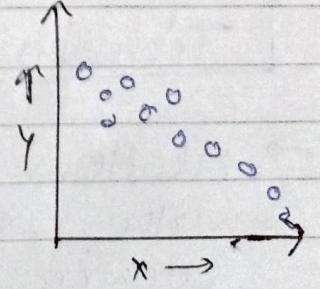
$\therefore$  If value below  $\pm 0.4$  we consider No Relation.



No Relation



+ve correlation



-ve correlation

## Pearson Correlation Coefficient

Definition : The Pearson corr. coeff. ('r') is the most common way of measuring a linear correlation.

- It is a Number between -1 & +1 that measures the strength and direction of the relationship between two variables.

## Multi-Variable Analysis

- Includes all Statistical techniques that are used to analyze more than two variables at once.
- Aim is to find patterns & correlations between several variables simultaneously.

### Analysis Techniques

↓  
**Dependence Techniques**

↓  
**Interdependence Techniques**

→ Relations are Dependent with each other in certain way

→ No relation b/w variables  
→ look for grouping them

## → Techniques :-

- ① Multiple Linear Regression →  
1 Depen Var  
2 or more Ind. Var.  
b/w 2
- ② Multiple Logistic Regression →  
spam  
not spam
- ③ Multi-variate analysis of Variance (MANOVA)
- ④ factor Analysis
- ⑤ Cluster Analysis.

## Multiple Linear Regression :-

→ It is dependence Method which looks at the relationship b/w one dependent & two or more Independent variables.

## Multiple Logistic Regression :-

- It is used to calculate (predict) the probability of a binary event occurring.
- Only two possible outcomes, either (1) or (0). like spam & Not Spam.

## MANOVA :-

→ MANOVA is used to measure the effect of multiple independent variables on two or more dependent variables.

### factor Analysis :-

- It is an interdependence technique which seeks to reduce the number of variables in a dataset.
- Model with too many variables can suffer over-fitting in which model fits too closely, making it less distinguishable to future datasets.

### Cluster Analysis :-

- It is an interdependence technique used to group similar items within a dataset into clusters.
- Clusters can be classified in 2 types :-

① Intracluster distance ,  
Distance b/w data points within  
One cluster.

② Intercluster distance ,  
Distance b/w data points in diff.  
clusters.

# Machine Learning

[Algo.]

## \* Simple Linear Regression :-

- It is a statistical method that allows us to summarize and study relationships between two continuous (quantitative) variables.
- $x \rightarrow$  'predictor', 'explanatory' or ('feature')  
'Independent Variable'
- $y \rightarrow$  'response', 'outcome', 'dependent variable'
- ~~Approach~~ Approach for predicting a response using a single feature.

### Steps of Doing it :-

- ① Analyze the correlation.
- ② Keep changing  $\theta_0$  &  $\theta_1$  or coefficient & intercept.
- ③ Minimal Learning Rate.
- ④ Obtain Global Minima.
- ⑤ Have best-fit line.

Residual Errors :-  $y_i - \hat{y}_i$   
Equation Basic :-  $y = mx + c$

$$y' = w_0 x_0 + b$$

or

Squared Error or Cost Function

$$h_{\theta}^{(x)} = \theta_0 + \theta_1 x_i$$

$$J(\theta_0, \theta_1) = \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

\* Quartiles :- Set of values that divide a list of numerical data into three quarters.

$$Q_1 = \left[ \frac{(n+1)}{4} \right] \text{th item}$$

$$Q_2 = \left[ \frac{n+1}{2} \right] \text{th item}$$

$$Q_3 = \left[ \frac{3(n+1)}{4} \right] \text{th item}$$

General formulae,

$$Q_r = l_1 + \frac{r\left(\frac{N}{4}\right) - c}{f} (l_2 - l_1)$$

$Q_r \rightarrow r^{\text{th}}$  Quartile

$l_1 \rightarrow$  lower limit

$l_2 \rightarrow$  upper limit

$f \rightarrow$  frequency

$c \rightarrow$  cumulative freq. of class preceding the quartile class.

$$\text{Quartile deviation} = \frac{(Q_3 - Q_1)}{2}$$

(midspread)

$\rightarrow$  Interquartile range :- (IQR) is the diff b/w upper & lower quartile.

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

## PANDAS PROFILING LIBRARY

⑦ Changing Object dtype into Numeric,

```
In['Gender'] = In['Gender'].map({'male': 0, 'female': 1})
```

## Scale of Measurity

- ① Nominal Eg → Eye colour, size.
- ② Ordinal Eg → Customer Service (★★★★★)  
↓  
In Proper order
- ③ Interval
- ④ Ratio

Q How will you decide that whether data-set is Categorical & Numerical.

## QUARTILES

- Quartiles set or split the data into 4 equal parts, first quartile Q1, divides the smallest 25% of the values from the other 75% that are larger.
- The second Quartile Q2 is the median i.e., 50% of the values are smaller than the median & 50% are larger.
- The third Quartile Q3 divides the smallest 75% of the values from the largest 25%.

- ~~Ques 22 30~~
- Q → What is regression  
Q → ML & its types  
Q → Stand. Error, SE, deviations  
Q → Dependent & Independent Variable

We can Minimize the Skewness :- (or Balance)

- ① Scaling
- ② Log Transformation.
- ③ Pre-processing.

\* Scales of Measurement, tells us how precisely variables are recorded.

→ Nominal : Data can only be  
Eg) (Gender, car brands) ① 'Categorized'.

→ Ordinal : Data can be, 'Categorized'  
Eg) (Top 5 Olympic Medalist) ① 'Categorized'  
② 'Ranked'. ①

→ Interval : Data can be :  
Eg) (Test Scores)  
① → Categorized  
② → Ranked  
③ → Evenly Spaced.

→ Ratio : Data can be :  
Eg) (Height, Age, weight)  
① → Categorized  
② → Ranked  
③ → Evenly Spaced  
④ → Has a Natural zero.

Q How to decide working data set  
is Categorical or Numerical ?  
Mainly based on 'OBJECTIVE',  
Categorical data means, data type  
that can be stored and identified  
based on the names or labels  
given to them.

→ A process called 'MATCHING' is done,  
to draw out the similarities or  
relations between the data and  
then they are grouped accordingly.

\* Numerical data refers to data  
that is in the form of numbers  
and not in any language or  
description form.

Q. What is Regression ?

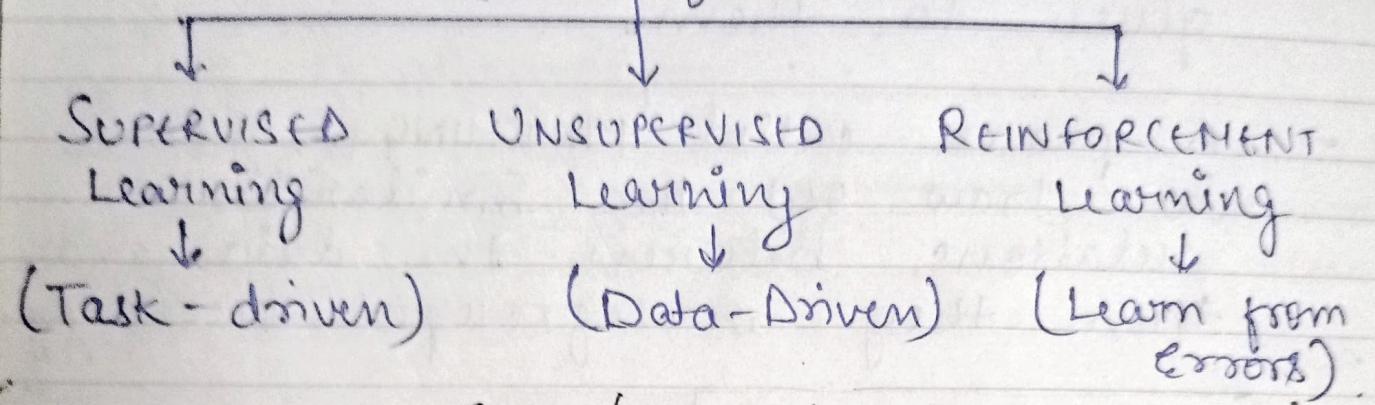
Ans It is a technique that attempts  
to determine the strength &  
character of the relationship  
between one dependent variable  
and a series of other variables.  
(aka Independent variables).

Q M.L and its types ?

Ans Machine Learning is defined as  
study of computer programs that  
leverage algorithms & statistical  
models to learn through inference

and patterns without being explicitly programmed.

## Types of M.L



Q What is 'Standard Deviation' & 'Standard Errors'?

	Standard Deviation	Standard Errors
①	Denotes variability within the Sample.	Denotes variability among multiple samples from a population.
②	Can be referred as 'Sample statistic'	Can be referred as 'Population Parameter'
③	Measures dispersion of individual value from their mean value.	measures accuracy of Sample mean as an estimate of the population mean.
④	It is type of ('Descriptive Stats')	It is type of ('Inferential Stats')

⑤ Calculated by square rooting the Variance.

Calculated as  $\frac{SD}{\sqrt{\text{Sample Size}}}$

⑥ cannot be used to calculate confidence interval of population

Used to calculate the same.

⑦ Not a function of sample size

depends on the size of sample, 3-21, S.E

$$\text{Sample Size} \propto S.E$$

Q 'Dependent' & 'Independent' Variables

Ans Effect (JOB SATISFACTION)

Cause (TYPE OF WORK ENVIRONMENT)

Independent Variables a.k.a

→ Explanatory Variables

→ Predictor Variables

→ Right-hand-side variables

(appear on right hand side of regression Eq.)

Type of I.V.

↓  
Experimental Variables.

↓  
(Can be manipulated Directly)

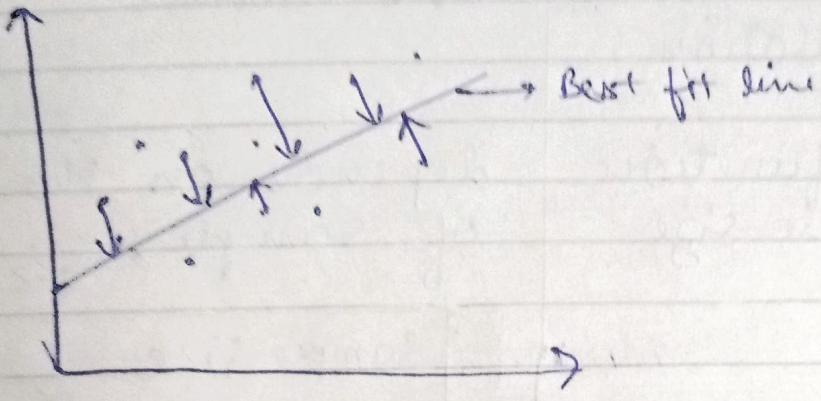
↓  
Subject Variables

↓  
(Cannot be manipulated Directly).

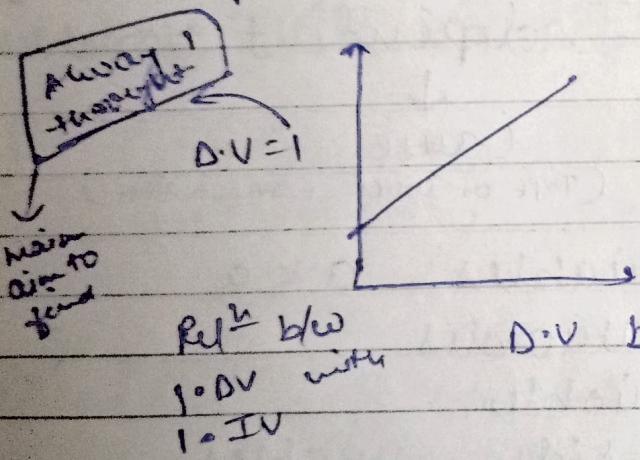
Dependent Variables  $\rightarrow$  L.H.S.

- $\rightarrow$  Response Variables
- $\rightarrow$  Outcome Variables
- $\rightarrow$  Left-hand-Side-Variables

---



Simple L.R.



D.V behaviour  $\rightarrow$  continuous/ decimal/ float-

Expression on which SRL works, = OLS  
 $\downarrow$   
 ordinary least squares.