

# \* Simple Linear Regression \*

① Supervised M/c Learning Algorithms:- Algo learns from labelled data.

↳ Labelled data - Dataset whose respective target value is already known.

↳ Supervised Learning has 2 types -

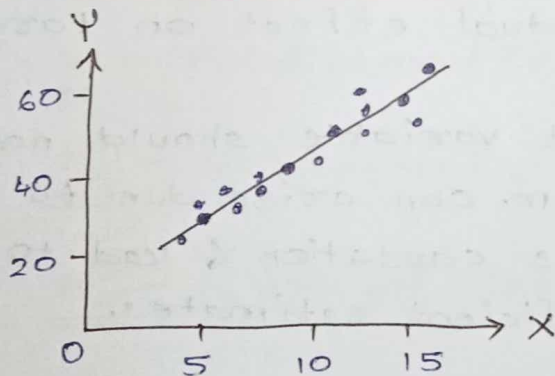
① Classification - Predict class of dataset based on independent i/p variable.

↳ Class is categorical/discrete values - like Yes/No, 0/1 etc.

② Regression - Predict continuous o/p variable.

↳ Ex. House Price Prediction.

② Linear Regression -



Used to model the relationship between a dependent variable & set of independent variables.

① Simple Linear Reg<sup>n</sup>

② Multiple Linear Reg<sup>n</sup>

• X : Dependent / Target variable.

• Y : Independent / Feature Variables.

• Regression Line : Best-fit line of the model, by which we can predict 'Y' value for new value of 'X'.

## ① Assumptions of Linear Regression:-

- ↳ Assumptions about data & its relationships.
- ↳ Violations can affect validity & reliability of regression results.

① Linearity — Relation bet<sup>n</sup> dependent & independent variable is linear. changes in both are constant.

② Independence of Errors — Errors (residuals) assumed independent of each other.

③ Homoscedasticity — Assume variance of residuals is constant across all levels of independent variables.

④ Normality of Errors — Errors normally distributed.

⑤ No or Little Multicollinearity — 2/more independent variables should not highly correlated. Difficult interpret each var individual effect on target.

⑥ No Endogeneity — Independent variable should not correlate with error term. can arise due to omitted bias / simultaneous causation & lead to biased & inconsistent coefficient estimates.

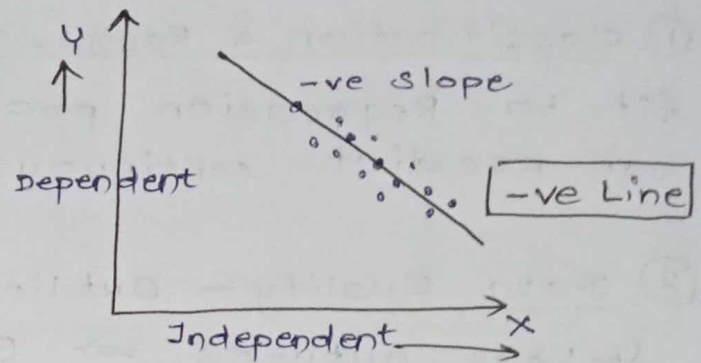
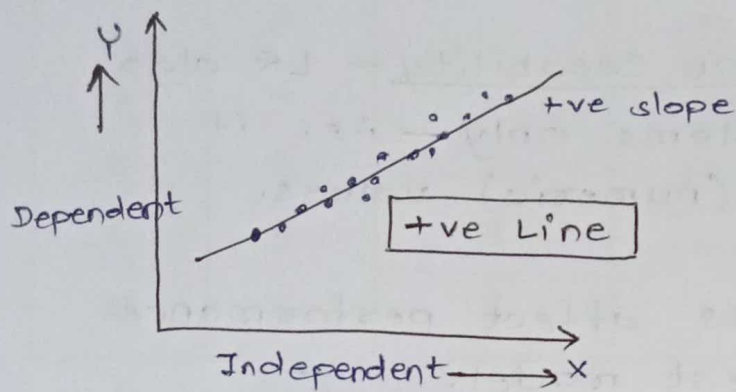
⑦ No Autocorrelation — Auto correlation (Model residuals correlated with each other). ~~observations~~ Imp in time series data, observa<sup>n</sup> dependent on previous.

⑧ No Perfect Collinearity — Perfect collinearity (one independent var perfectly predicted linearly by others independent variable).

↳ Leads to rank-deficient matrix, making it impossible to estimate unique regression coefficients.



## ① Linear Regression:



1) Simple Linear Regression: Model relationship between 1 dependent & 1 independent variable.

↳ This relationship can be represented as a linear equation.

• Eq<sup>n</sup> of Simple Linear Regression -

$$Y = \alpha_0 + \alpha_1 X + \epsilon$$

Y - Dependent (Target) variable.

X - Predictor / Feature / Independent variable.

$\alpha_0$  - Y-intercept, value of Y when  $x = 0$ .

$\alpha_1$  - Slope of line (Regression coefficient of X).

$\epsilon$  - Represents error term.

The equation is similar to eq<sup>n</sup> of Line:

$$Y = mX + C$$

Slope                      Y-intercept

## • Applications of Linear Regression:

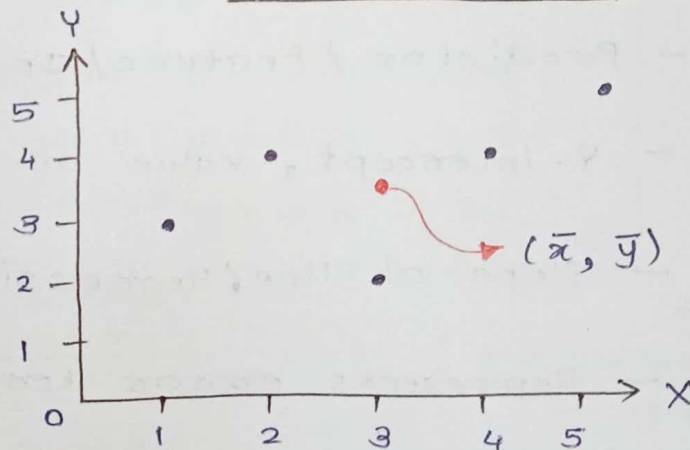
- ① Classification & Regression Capability - LR algo fit to Regression problems only - i.e. it can predict continuous (numeric) values.
- ② Data Quality - Outliers affect performance.
  - ↳ Less outliers  $\Rightarrow$  Best model.
  - ↳ So, data should not contain outliers.
- ③ Computational Complexity - LR algo is computationally less expensive than classification algo.
- ④ Comprehensive & Transparent - Easy to understand & explain. Represented using simple math expression.

## • Working of Linear Regression:

### • Data Set -

X	Y
1	3
2	4
3	2
4	4
5	5

### • Plot the Graph -



$$\bar{x} = \frac{1 + 2 + 3 + 4 + 5}{5} = 3$$

$$\bar{y} = \frac{3 + 4 + 2 + 4 + 5}{5} = 3.6$$

### • Eq<sup>n</sup> to Find Slope -

$$m = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$$

$x$	$y$	$(x - \bar{x})$	$(y - \bar{y})$	$(x - \bar{x})(y - \bar{y})$	$(x - \bar{x})^2$
1	3	-2	-0.6	1.2	4
2	4	-1	0.4	-0.4	1
3	2	0	-1.6	0	0
4	4	1	0.4	0.4	1
5	5	2	1.4	2.8	4
3	3.6			4.0	10

$$\therefore m = \frac{4}{10} = 0.4$$

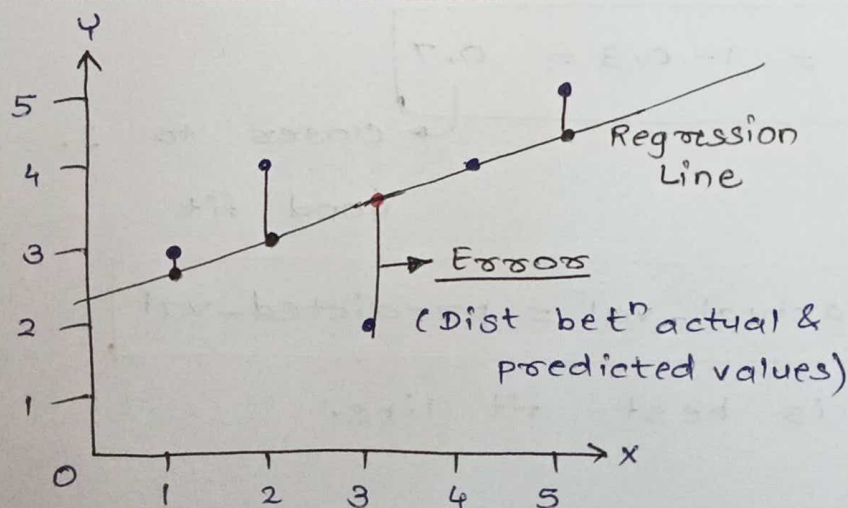
$$c = \bar{y} - m\bar{x} = 3.6 - 0.4 * 3 = 2.4$$

For  $m = 0.4$  &  $c = 2.4$ ,

Predict  $y$  for  $x = \{1, 2, 3, 4, 5\}$

$x$	$y_{\text{predicted}}$
1	$0.4 * 1 + 2.4 = 2.8$
2	$0.4 * 2 + 2.4 = 3.2$
3	$0.4 * 3 + 2.4 = 3.6$
4	$0.4 * 4 + 2.4 = 4.0$
5	$0.4 * 5 + 2.4 = 4.4$

• Plot final graph -



• Minimize Error:

Iteratively  
find  $m$ , predict  $y$   
& find error.

$m$  with min  
error  $\Rightarrow$  Regr<sup>n</sup>  
Line



## ① Find Goodness of fit:-

- R-Square - statistical measure of how close is data to fitted Regression Line.

↳ Also known as coefficient of determination / coefficient of multiple determination.

Dist (actual - mean)  
v/s  
Dist (predicted - mean)

$$R^2 = 1 - \frac{\sum (y_p - \bar{y})^2}{\sum (y - \bar{y})^2}$$

↖ Predicted
↗ Mean  
↘ Actual

x	y	(y - $\bar{y}$ )	(y - $\bar{y}$ ) <sup>2</sup>	y <sub>p</sub>	(y <sub>p</sub> - $\bar{y}$ )	(y <sub>p</sub> - $\bar{y}$ ) <sup>2</sup>
1	3	-0.6	0.36	2.8	-0.8	0.64
2	4	0.4	0.16	3.2	-0.4	0.16
3	2	-1.6	2.56	3.6	0	0
4	4	0.4	0.16	4.0	0.4	0.16
5	5	1.4	1.96	4.4	0.8	0.64
			5.2			
						1.6

$$R^2 = 1 - \frac{1.6}{5.2} = 1 - 0.3 = 0.7$$

↳ closer to 1.

Good fit

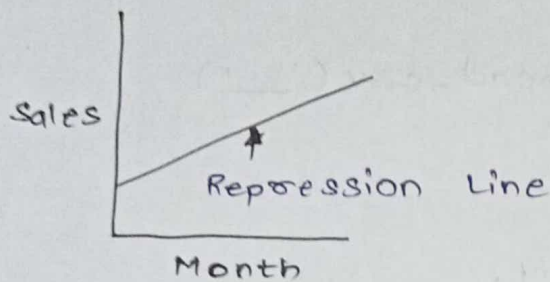
If  $R^2 = 1 \Rightarrow \text{actual-val} = \text{predicted-val}$

↳ this is best fit line.

## • Examples of Linear Regression:-

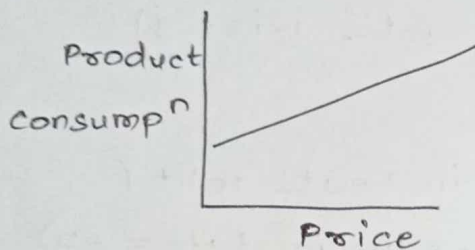
### ① Evaluate Trends & Sales estimates -

Can estimate sales in future.



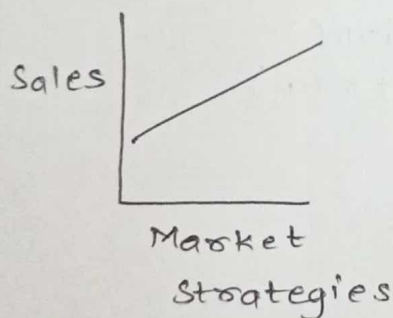
### ② Analyze Impact of Price change -

How price change affect the consumption of product.



### ③ Risk Assessment in Financial Services -

How risk management affect business finance.



## • Simple Linear Regression Implementation:-

① Import Libraries - `pd, np, plt.`



② Import Dataset - `df = pd.read_csv('___')`



③ EDA -

- Null value treatment.
- Remove duplicates
- Handle categorical data.



④ Split Dataset - `from sklearn.model_selection ⇒ train_test_split`  
`X = df.drop('target', axis = 1)`  
`Y = df['target']`

`x_train, x_test, y_train, y_test = train_test_split(X, Y, test_size = 0.2, random_state = 42)`



⑤ Training Model - `from sklearn.linear_model ⇒ LinearRegression`  
`regressor = LinearRegression()`  
`regressor.fit(x_train, y_train)`



⑥ Result Prediction -

`y_pred = regressor.predict(x_test)`



⑦ Model Evaluation - `from sklearn.metrics ⇒ r2_score, mean_squared_error`  
Find model's goodness of fit.

`r2 = r2_score(y_test, y_pred)`

`sq_err = mean_squared_error(y_test, y_pred)`