

Introduction to Machine Learning

Machine Learning (ML) is the subset of artificial intelligence that focuses on building systems that can learn from and make decisions or predictions based on the data.

In simple terms, ML enables computers to improve their performance on tasks by learning from experience (data) rather than relying on hard-coded rules.

✳ Key components of Machine Learning :

1. Data : the foundation of ML ; used for training & testing models.
2. Algorithms : Mathematical methods that process data to identify patterns or make predictions
3. Model : the representation of knowledge learned by an ML algorithm.
4. Features : Input variables used for prediction or decision making.
5. Labels (In supervised learning) : Known outcomes used to train the model.

✳ Categories of Machine Learning :

- ① Supervised Learning : The model learns from labeled data (e.g. predicting house prices based on historical data)

2. Unsupervised Learning : the model identifies

Patterns in unlabeled Data (eg clustering
customers segments).

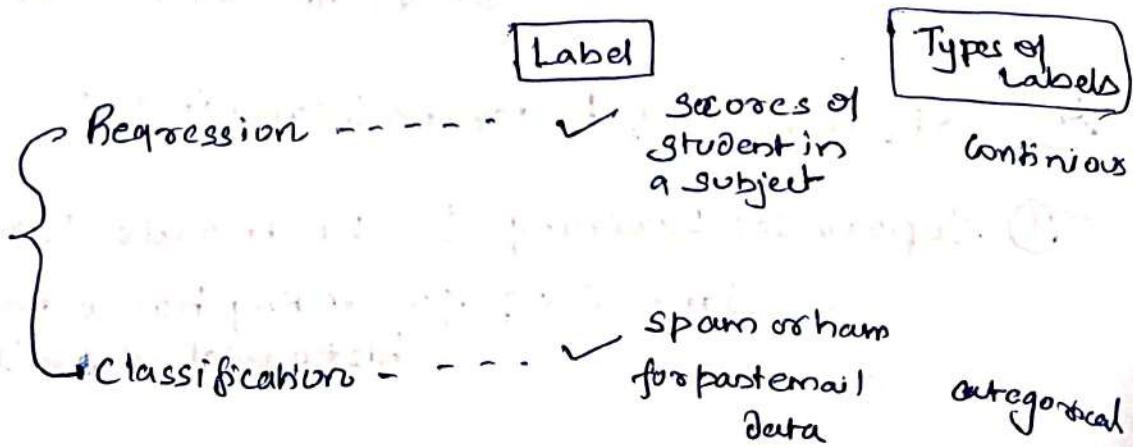
3. Reinforcement Learning : the model learns

through trial and error by receiving rewards
or penalties (eg training a robot to walk)

* Machine Learning models can be classified
into the following three types based on the
task performed and the nature of the output

- ① Regression : the output variable to be predicted
is a continuous variable eg scores of student
- ② Classification : the output variable to be predicted
is a categorical variable eg classifying
incoming emails as spam or ham.
- ③ Clustering : no predefined notion of a label is
allocated to groups / clusters formed
eg customer segmentation.

Supervised
Learning
Method



* Supervised Learning Methods

- ① Past data with labels is used for building the model.
- ② Regression and classification algorithms fall under this category.

* Unsupervised Learning methods

- ① No pre-defined labels are assigned to input data.
- ② Clustering algorithms fall under this category.

→ Supervised Learning Examples

→ Supervised learning involves labeled datasets where the outcome or target variable are already known.

① Spam email classification

→ Input: features like email text, sender's address and subject line.

→ Output: Labels like 'spam' or 'notspam'

② House Price Prediction

→ Input: features such as bedrooms, location & size.

→ Output: Predicted house price

Unsupervised Learning Examples :

It deals with unlabeled data, finding patterns or structures within data.

① Anomaly Detection

→ Input: Network traffic data, sensor data etc

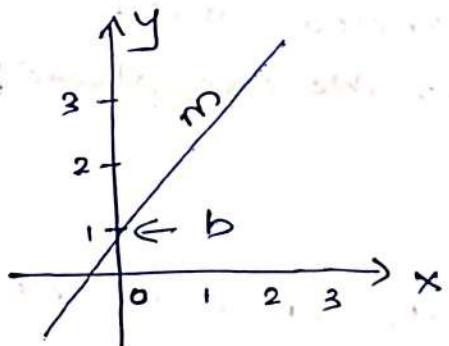
→ Output: Identifying unusual patterns that could indicate security threats.

Regression Line

first let's understand about the equation of a straight line

the equation of a straight line is usually written as $y = mx + b$ or $y = mx + c$

What does it stand for?



$$y = mx + b$$

slope

b value, when
 $x = 0$

$\Rightarrow y = \text{how far up}$

$x = \text{how far along}$

$m = \text{slope or gradient}$

(how steep the line is)

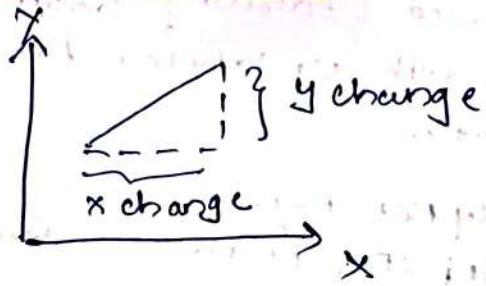
$b \Rightarrow \text{value of } y \text{ when } x = 0$

How do we find "m" and "b"?

$\rightarrow b$ is easy: just see where the line crosses Y axis.

$\rightarrow m$ (the slope) needs some calculation

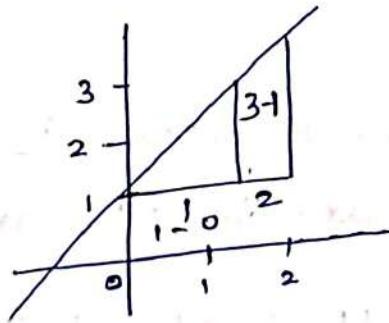
$$m = \frac{\text{change in } Y}{\text{change in } X}$$



Knowing this we can work out the eqn of the straight line:

$$m = \frac{3-1}{1-0} \therefore 2/1 = 2$$

$$b = 1$$



∴ eqn of straight line $\Rightarrow y = mx + c \Rightarrow y = 2x + 1$

∴ Slope = $m = \frac{\Delta y}{\Delta x} = \frac{y_2 - y_1}{x_2 - x_1}$

Positive Slope

- A +ve slope means that as x increases, y also increases.
- In a graph a line goes upwards from left to right.

eg: A line passing through points $(1, 2)$ and $(2, 4)$

$$\Rightarrow m = \frac{4-2}{2-1} \therefore 2/1 = 2$$

Negative Slope

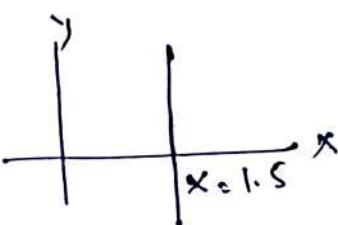
- A -ve slope means that as x increases, y decreases.
- In a graph the line goes downwards from left to right.

eg: A line passing through points $(1, 4)$ and $(2, 2)$

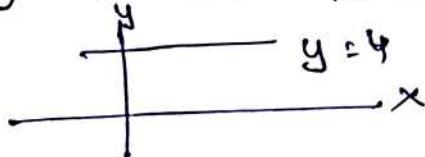
$$m = \frac{2-4}{2-1} \therefore -\frac{2}{1} = -2$$

Zero slope

- if the line is horizontal (y doesn't change) as x changes \Rightarrow the slope is zero



eg $y = 4$ or $x = 1.5$



Regression Line

A regression line is a straight line that best fits the data points on a scatter plot, it represents the relationship b/w the independent variable (x) and dependent variable (y). In statistical terms, it minimizes the sum of squared differences between the actual data points and the predicted data points on the line.

*Equation of a regression line

$$y = mx + b$$

where

y : Dependent Variable
(Predicted value)

x : Independent Variable
(Input value)

m = slope of the line

$$m = \frac{\Delta y}{\Delta x}$$

b = intercept (the value of y when $x=0$)

(eg) Predicting House Price

Data: x : Size of the house in square feet

y : Price of the house in \$1000s

Size (sq ft, x)

1000

1800

2000

2500

Price (y)

150

200

250

300

$$m = \frac{\Delta y}{\Delta x} = \frac{200 - 150}{1500 - 1000} = \frac{50}{500} = \frac{1}{10} = 0.1$$

$$y = 0.1x + 50$$

Interpretation :

- $\rightarrow m = 0.1$ for every additional square foot, the price increases by \$100
- $\rightarrow b = 50$: A house of 0 sqft would hypothetically cost \$50000.

Best fit line

In regression there is a notation of the best fit line the line which fits the given scatterplot in the best way. So...

The Best-fit line in regression is the line that most accurately represents the relationship b/w an independent variable (x) and the dependent var (y). It is determined using Ordinary Least Squares (OLS) method, which minimizes the Residual Sum of Squares (RSS), the sum of squared differences b/w the actual and predicted values.

Process of finding the Best fit line

- ① Start with Scatterplot : useful to visualize the relationship b/w independent & dependent variable.

2. Define Residuals

Residuals (e_i) represents the error ^{bias} the actual values y_i and predicted values \hat{y}_i :

$$e_i = y_i - \hat{y}_i$$

the smaller the residuals, the better the fit of the line.

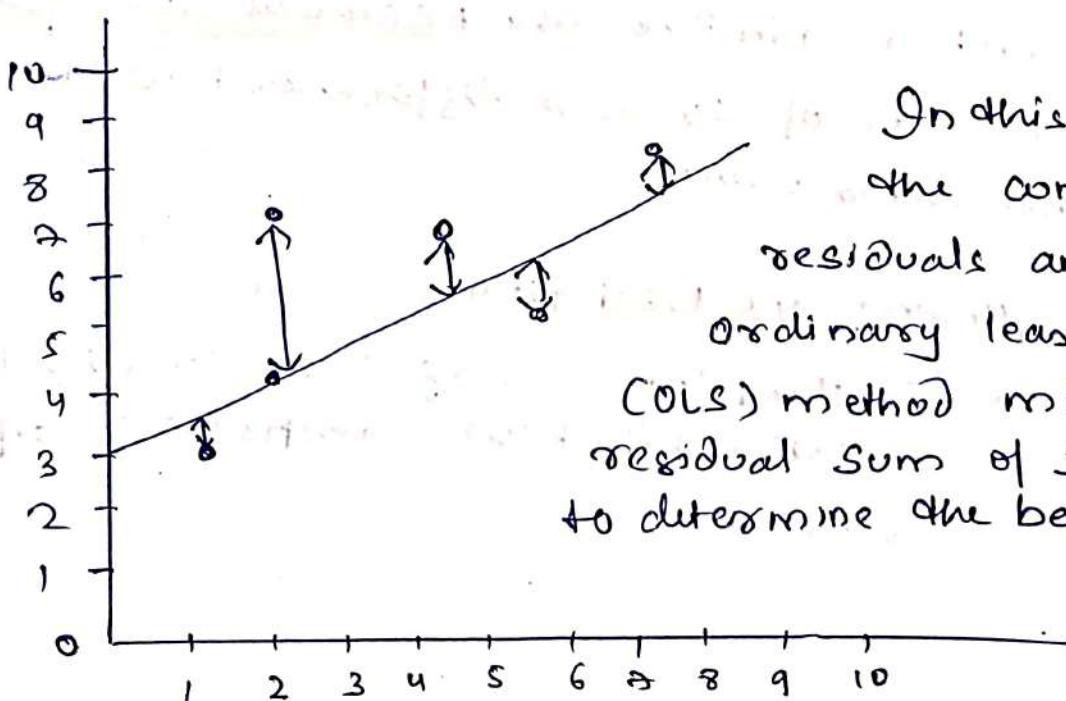
3. Residual Sum of Squares (RSS)

to evaluate how well a line fits the data calculate the RSS:

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

the goal is to find the line that minimizes RSS, the method for achieving the Ordinary Least Squares.

4. Ordinary Least Squares (OLS)



In this diagram we see the concept of residuals and how the Ordinary Least Squares

(OLS) method minimized the residual sum of squares to determine the best fit line

Key elements from the image:

① Equation of the regression line:

$$y = B_0 + B_1 x$$

→ B_0 : Intercept (value of y when $x=0$)

→ B_1 : Slope (change in y for one-unit change in x).

② Residuals (e_i)

→ the difference b/w the actual values (y_i) and the predicted value (\hat{y}_i)

$$e_i = y_i - \hat{y}_i$$

→ Residuals are shown as vertical lines b/w the observed points (dots) and the regression line.

③ Residual Sum of Squares (RSS)

→ the RSS formula calculates the total squared error:

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (B_0 + B_1 x_i))^2$$

$$\Rightarrow \sum_{i=1}^n (y_i - B_0 - B_1 x_i)^2$$

④ Goal of OLS:

Minimize the RSS to find the optimal values of B_0 and B_1 , ensuring the regression line best fits the data.

Gradient Descent is an optimization algorithm used to find the values of the parameters (coefficients) of a function (f) that minimizes a given cost function (cost)

Unconstrained and Constrained Minimization

"Constrained Minimization" is a process in which we try to optimise an objective function with respect to some variables in the presence of the constraints on the variables of the function.

↳ In our case the objective func is a cost function.

"Unconstrained Minimization" on the other hand is a process in which we try to optimise the objective func with respect to some variables without any constraints on those variables.

Let's study them in detail

★ Unconstrained minimization is an essential concept in linear regression and optimizatn where we aim to find the best fit line for a given dataset by minimizing the Cost func usually the Mean Squared Error or MSE

In linear regression, we model the relation b/w the input features (x) and the target (y) as

$$\hat{y} = \omega x + b$$

where

ω → weight or slope

b → bias (intercept)

\hat{y} → predicted value

Goal : to find ω and b that minimize the diff b/w the predicted values (\hat{y}) and actual value (y)

Objective : Minimize the cost func

The cost func measures how well the line fits the data, for linear regression we often use Mean Squared Error (MSE)

$$J(\omega, b) \underset{\downarrow}{=} \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - (wx_i + b))^2$$

depends on weight & bias

Here: $N \Rightarrow$ nos of datapoints

$y_i \Rightarrow$ is the actual target value

$x_i \Rightarrow$ is the input feature value

$wx_i + b$ is the predicted value

In Unconstrained minimization, we aim to minimize $J(w, b)$ without any constraints on w or b . This process involves

- ① Computing gradients : Derivatives of $J(w, b)$ with respect to w and b
- ② Setting gradients to zero : Solve for w and b where the gradients vanish (stationary point)
- ③ Finding the minimum : Check if stationary point corresponds to minimum.

Simple Example

Dataset

x	y
1	2
2	3
3	4

Step 1 : Write the cost function

$$J(w, b) = \frac{1}{3} \sum_{i=1}^3 (y_i - (w x_i + b))^2$$

Step 2 : expand the cost function

$$J(w, b) = \frac{1}{3} [(2 - (w \cdot 1 + b))^2 + (3 - (w \cdot 2 + b))^2 + (4 - (w \cdot 3 + b))^2]$$

The cost func $J(w, b)$ is the Mean Squared error

$$J(w, b) = \frac{1}{N} \sum_{i=1}^N (y_i - (w x_i + b))^2$$

To minimize $J(\omega, b)$ we need to compute its gradient i.e. the partial derivatives with respect to ω and b .

* Gradient w.r.t. ω

The partial derivatives of $J(\omega, b)$ w.r.t. ω measures how $J(\omega, b)$ changes as the ω changes.

$$\frac{\partial J}{\partial \omega} = \frac{\partial}{\partial \omega} \left[\frac{1}{N} \sum_{i=1}^N (y_i - (\omega x_i + b))^2 \right]$$

Bringing the derivative inside the summation

$$\frac{\partial J}{\partial \omega} = \frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial \omega} (y_i - (\omega x_i + b))^2$$

Use the chain rule for $(y_i - (\omega x_i + b))^2$, let $u = y_i - (\omega x_i + b)$ then

$$\frac{\partial}{\partial \omega} u^2 = 2u \frac{\partial u}{\partial \omega}$$

$$\Rightarrow \frac{\partial}{\partial \omega} (y_i - (\omega x_i + b))^2 = 2 \cdot (y_i - (\omega x_i + b)) \cdot (-x_i)$$

(Note: the derivative of $y_i - (\omega x_i + b)$ w.r.t. ω is $-x_i$)

$$\frac{\partial J}{\partial \omega} = \frac{1}{N} \sum_{i=1}^N [2 \cdot (y_i - (\omega x_i + b)) \cdot (-x_i)]$$

$$\frac{\partial J}{\partial \omega} = -\frac{2}{N} \sum_{i=1}^N x_i (y_i - (\omega x_i + b))$$

Gradient w.r.t to b

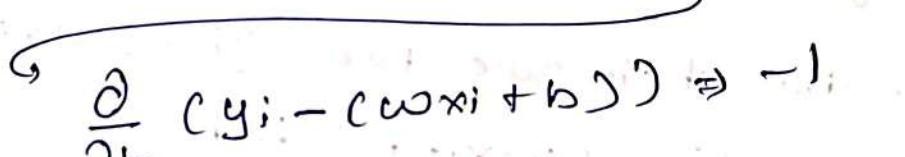
$$\frac{\partial J}{\partial b} = \frac{\partial}{\partial b} \left[\frac{1}{N} \sum_{i=1}^N (y_i - (\omega x_i + b))^2 \right]$$

Bringing the derivative inside summations

$$\frac{\partial J}{\partial b} = \frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial b} (y_i - (\omega x_i + b))^2$$

Using the chain rule

$$\Rightarrow 2 \cdot (y_i - (\omega x_i + b)) \cdot \frac{\partial}{\partial b} (y_i - (\omega x_i + b))$$



$$\frac{\partial}{\partial b} (y_i - (\omega x_i + b)) \Rightarrow -1$$

$$\Rightarrow \frac{\partial J}{\partial b} = \frac{1}{N} \sum_{i=1}^N [2 \cdot (y_i - (\omega x_i + b))] (-1)$$

$$\Rightarrow -\frac{2}{N} \sum_{i=1}^N (y_i - (\omega x_i + b))$$

Setting $\frac{\partial J}{\partial \omega} = 0$ and $\frac{\partial J}{\partial b} = 0$ to get the optimal values of ω, b :

$$\Rightarrow \omega = 1, b = 1$$

$$\Rightarrow \boxed{y = x + 1}$$

Why unconstrained?

→ there are no restrictions on values of w and b can take.

What is Constrained Minimization?

In constrained minimization, we still aim to find the best values of the parameters (like w and b in regression) that minimize the error or. However we add rules or constraints to ensure the solution satisfies certain conditions.

Eq: Imagine you are choosing a path to minimize the travel time but are required to use certain roads only, this restriction is the constraint.

for regression, a constraint could be something like:

→ Keep the model's weight w small (to prevent overfitting),

→ Ensure the sum of weights equals a fixed value.

We know that the equation is $y = wX + b$

where

$y \Rightarrow$ predicted value

$X \Rightarrow$ input feature

$w \Rightarrow$ is the weight (slope)

$b \Rightarrow$ bias (intercept)

The goal is to minimize the error?

Error (Cost function):

$$J(w, b) = \frac{1}{N} \sum_{i=1}^N (y_i - (wX_i + b))^2$$

This is called as unconstrained minimization, where we freely adjust w and b to minimize error.

Constrained Minimization Example:

Let's add a constraint:

"The weight w must not exceed a certain value (e.g. $w^2 \leq \lambda$, where λ is a small no")

Why add constraints?

- ① Prevent overfitting: if w becomes very high, the model can overfit the training data, meaning it won't generalize well to new data.
- ② Real-world requirements: sometimes weights must meet specific physical or financial limits.

How to solve constrained minimization?

To solve the constrained minimization problems, we slightly modify how we minimize the cost func.

Simple example:

Suppose we have:

$$\rightarrow \text{cost function: } J(w) = (y - wx)^2$$

$$\text{constraint: } w \leq 2$$

① without constraint

If we ignore the constraint, we minimize $J(w)$ by taking the derivative and finding where it evaluates to 0. Let's say this gives us

$$w = 5$$

But $w = 5$ violates the constraint $w \leq 2$

② with constraint

Instead of freely choosing w , we must respect the constraint, in this case $w = 2 \Rightarrow$ final soln

why? because $w = 2$ is the largest value that satisfies $w \leq 2$ and going beyond it would break the rule.

Common ways to handle constraints

- ① Hard constraint (as in the above eq)
 forcing the soln to strictly obey the rule
 eq $w \leq 2$ so the final weight is
 capped at 2.

- ② Soft constraint (Regularization)

Instead of forcing a strict rule, we penalize large weights for eq modify the cost func to

$$J(w, b) = \frac{1}{N} \sum_{i=1}^N (y_i - (wx_i + b))^2 + \alpha w^2$$

Here α is a small no that controls how much we penalize ^{large} weights.

* Optimal Weight formula Derivation

Key points:

- ① Objective: Minimize the cost func to find the best weight w .

$$\text{cost func : } J(w) = \frac{1}{N} \sum_{i=1}^N (y_i - wx_i)^2$$

bias = 0

this measures how far the predicted values wx_i are from the actual values y_i

Step 2: Expanding the cost function

→ Expand $(y_i - \omega x_i)^2$

$$J(\omega) = \frac{1}{N} \sum_{i=1}^N (y_i^2 - 2\omega x_i y_i + \omega^2 x_i^2)$$

→ Separate the terms

$$J(\omega) = \frac{1}{N} \left[\sum_{i=1}^N y_i^2 - 2\omega \sum_{i=1}^N x_i y_i + \omega^2 \sum_{i=1}^N x_i^2 \right]$$

Step 3: Differentiating with respect to ω :

take the derivative to minimize $J(\omega)$:

$$\frac{\partial J(\omega)}{\partial \omega} = \frac{1}{N} \left[-2 \sum_{i=1}^N x_i y_i + 2\omega \sum_{i=1}^N x_i^2 \right]$$

Simplifying the term:

$$\frac{\partial J(\omega)}{\partial \omega} = -\frac{2}{N} \sum_{i=1}^N x_i y_i + \frac{2\omega}{N} \sum_{i=1}^N x_i^2$$

Step 4: Setting the derivative to 0, ie to find the

minimum set $\frac{\partial J(\omega)}{\partial \omega} = 0$

$$-\sum_{i=1}^N x_i y_i + \omega \sum_{i=1}^N x_i^2 = 0$$

$$\omega \sum_{i=1}^N x_i^2 = \sum_{i=1}^N x_i y_i \Rightarrow \omega = \frac{\sum_{i=1}^N x_i y_i}{\sum_{i=1}^N x_i^2}$$

Final formula: $\omega = \frac{\sum (x_i y_i)}{\sum x_i^2}$

Interpretation

$\Rightarrow \sum(x \cdot y) \Rightarrow$ Measures how x and y vary together

$\sum(x^2) \Rightarrow$ Measures the total spread of x

(e.g.)

Given $x = [1, 2, 3]$, $y = [2, 4, 6]$

$$\text{Compute } \sum(x \cdot y) = 2 + 8 + 18 = 28$$

$$\text{Compute } \sum x^2 = 1^2 + 2^2 + 3^2 = 1 + 4 + 9 = 14$$

$$w = 28/14 = 2 \Rightarrow \text{Optimal line } \Rightarrow y = 2x$$

Note: the value of RSS would change because the units are changing because %

the RSS for any regression line is given by the expression $\sum(y_i - y_{\text{pred}})^2$, RSS is the sum of squared difference b/w the actual value and the predicted value and its value will change if the units change since it has units of y^2

for e.g. $(140\text{USD} - 70\text{USD})^{12} = 4900\text{RSS}$ whereas $(2\text{USD} - 1\text{USD})^{12} = 1$ so value of RSS is diff in both cases.

(*) TSS (Total Sum of Squares)

The total sum of squares is a way to measure how much the data values vary in the dataset

It helps us to understand the total amount of spread or variability in the dependent variable (y), which is the variable you are trying to predict or explain.

Why is TSS important?

When you build a regression model, the goal is to explain or predict this variability. TSS gives us a starting point to see how much total variation exists before any model is applied.

Breaking It Down:

Imagine you have some data points representing the exam scores

$$y = [50, 60, 70, 80, 90]$$

Step 1: find the mean (\bar{y}): the mean is the avg of all the data points:

$$\bar{y} = \frac{\text{sum of all } y}{\text{total number of } y} \Rightarrow \bar{y} = \frac{50+60+70+80+90}{5}$$

$$\Rightarrow \bar{y} = 70$$

Step 2: calculate how far each value is from the mean
Subtract the mean from each data point

$$(y_i - \bar{y})$$

for each y :

$$[50 - 70, 60 - 70, 70 - 70, 80 - 70, 90 - 70] \\ \Rightarrow [-20, -10, 0, 10, 20]$$

Step 3: Square each difference; squaring removes
-ve signs and gives more weight to
larger differences!

$$[-20^2, -10^2, 0^2, 10^2, 20^2] \Rightarrow [400, 100, 0, 100, 400]$$

Step 4: Add them up; Add these squared differences

$$TSS = 400 + 100 + 0 + 100 + 400 = 1000$$

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2$$

$y_i \rightarrow$ Actual Data Values
 $\bar{y} \rightarrow$ Mean of the data

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$\hat{y}_i \rightarrow$ Predicted values from the regression model.

R-Squared (R^2):

R^2 is a measure of how well the regression line explains the variability in the dependent variable y ; It is calculated as:

$$R^2 = 1 - \frac{RSS}{TSS}$$

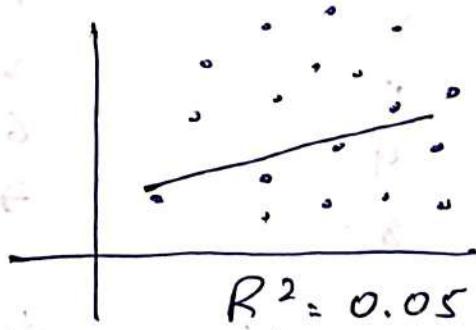
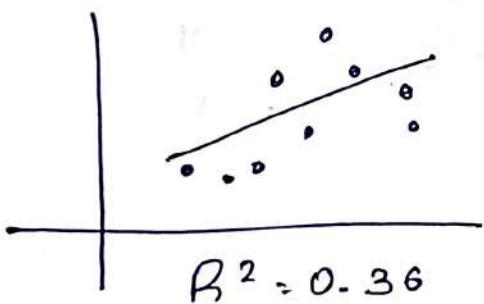
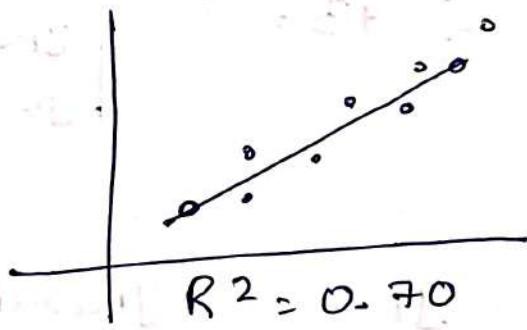
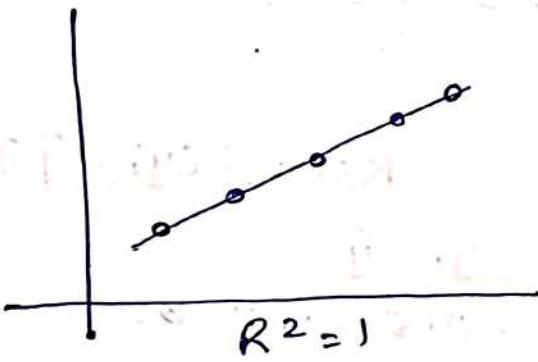
$$R^2 = \frac{MSS}{TSS} = 1 - \frac{RSS}{TSS}$$

$R^2 \rightarrow$ Proportion of the total variability explained by the model.

R_{SS}/TSS : Proportion of variability not explained by the model.

If $R^2 = 0.6$, it means the model explains 60% of the variability in the data, while remaining 40% is unexplained.

Physical Significance Of R^2



Q Given the eqn of line

$$y = \frac{x}{2} + 3$$

Q1 To find the value of RSS

$$y_{\text{Pred1}} = \frac{1}{2} + 3 = 3.5$$

$$y_1 = 3$$

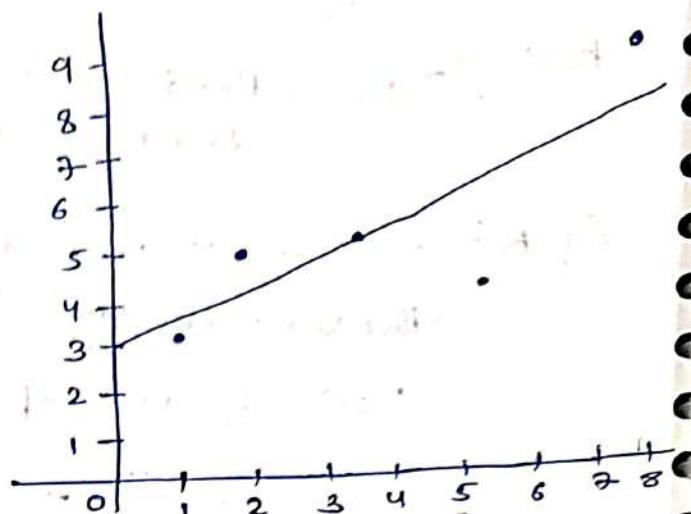
$$y_{\text{Pred2}} = \frac{2}{2} + 3 = 4$$

$$y_2 = 5$$

$$y_{\text{Pred3}} = \frac{4}{2} + 3 = 5$$

$$y_4 = 5$$

$$\left| \begin{array}{l} y_{\text{Pred6}} = \frac{6}{2} + 3 = 6 \\ y_6 = 4 \\ y_{\text{Pred8}} = \frac{8}{2} + 3 = 4 + 3 = 7 \\ y_8 = 8 \end{array} \right.$$



$$RSS = \sum (y_i - \hat{y})^2$$

⇒ table

y_i	$y_{\text{Predicted}}$	$y_i - \hat{y}$
3	3.5	-0.5 → 0.25
5	4	1 → 1
5	5	0 = 0
4	6	-2 = 4
8	7	1 = 1
RSS		6.025

Q2 To find TSS

$$y_i \quad \bar{y} \Rightarrow \frac{3+5+5+4+8}{5} = \frac{25}{5} = 5$$

$$TSS = \sum (y_i - \bar{y})^2, (4+0+0+1+9)$$

$$Q \text{ find } R^2 = 1 - \frac{\text{RSS}}{\text{TSS}} = 1 - \frac{6.25}{14} = \frac{14 - 6.25}{14} = \frac{7.75}{14}$$

Apart from R^2 , there is one more quantity named RSS (Residual Square Error)

$$\text{RSS} = \sqrt{\frac{\text{RSS}}{df}} \quad \text{Here } df = n - 2 \quad \text{where } n \text{ is no. of data points}$$

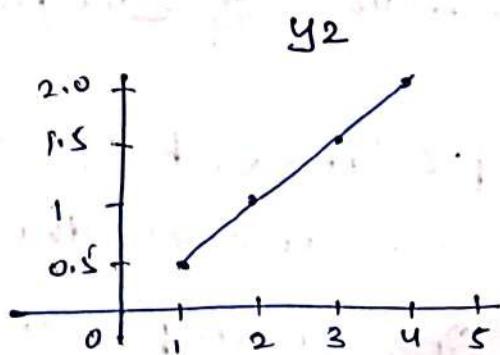
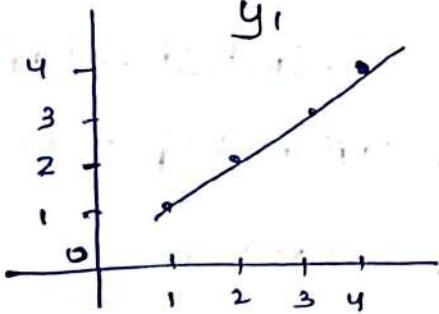
Note: the value of correlation coefficient always lies b/w -1 and +1, where a negative value implies a negative correlation and a +ve value shows a positive correlation and 0 means no correlation.

$\therefore R^2 = 1 - \frac{\text{RSS}}{\text{TSS}}$ \rightarrow the value of R^2 lies b/w 0 and 1 where 1 implies that the variance in the data is being explained by the model

0 means none of the variance values is being explained by the model, obviously it's difficult to achieve either of the extremes.

following two straight lines fitted for two sets of data

Note:



Q value of correlation coefficient (or both 1's 1 and 1) these slopes are 1 and 0.5 but \therefore all the points lie on the straight line correlation coefficient becomes equal to 1. Slope of straight line don't represent correlation.

Q Suppose that we run a linear regression on the data where $X \rightarrow$ Independent var and $Y \rightarrow$ Target var, you find that the correlation b/w y and x is -0.92 , what can you say here?

Ans The absolute value of the correlated coefficient is very high, -ve sign just implies that x and y are -vely correlated, Hence they have very strong -ve correlation.

Q The Independent var x from a linear regression is measured in miles, If we convert it to kms keeping the unit of the dependent var y the same, how will the slope coefficient change?

$$\text{Note } 1\text{ mile} = 1.6\text{ km}$$

Ans In the linear regression even, x gets multiplied by 1.6 with no change in y . So it means slope will be divided by 1.6.

Summary of Simple linear Regression

Simple Linear Regression is a statistical method used to model the linear relationship between:

- One Independent variable (X) and
- One dependent variable (Y)

The equation of a simple linear regression model is:

$$Y = \beta_0 + \beta_1 X + \epsilon$$

Where :

- $y \Rightarrow$ Dependent Variable (Predicted Value)
- $x \Rightarrow$ Independent Variable (Predictor)
- $b_0 \Rightarrow$ Intercept (Value of y when $x=0$)
- $b_1 \Rightarrow$ Slope (Rate of change of y for a unit change in x)
- $\epsilon \Rightarrow$ Error term (Unexplained variability)

② Assumptions of Simple Linear Regression

For the model to be valid, the following assumptions must hold:

- ① Linearity : the relationship b/w x and y is linear
- ② Independence : Observations are independent of each other.
- ③ Homoscedasticity : the variance of the residual (errors) is constant across all values of x .
- ④ Normality of Residuals : Residuals are normally distributed.

③ Interpretation of Coefficients

- ① Intercept (b_0) : the predicted value of y when $x=0$
- ② Slope (b_1) : the change in y for one unit increase in x

4. Goodness of fit (R^2)

- ⇒ R^2 (Coefficient of Determination) measures the proportion of variation in Y (dependent) explained by X (independent)
- ⇒ Value ranges from 0 to 1
 - ⇒ $R^2 = 0$ (X explains none of the variability in Y)
 - ⇒ $R^2 = 1$ (X explains all the variability in Y)
- Higher R^2 indicates better fit.

5. Correlation vs Regression

Correlation measures the strength and direction of the linear relationship b/w X and Y , range (-1 to 1)

Regression quantifies the exact relationship (eq slope and intercept) and allows for prediction.

6. Residual Analysis

- ⇒ Residuals ($Y_{\text{actual}} - Y_{\text{predicted}}$) represent the error in predictions.
- ⇒ A good regression model will have residuals centered around 0
- But if 0 can lead to overfit

7. Limitations

- ⇒ Sensitive to outliers
- ⇒ Only models linear relationships (non linear relationships require other methods)



Simple Linear Regression In Python

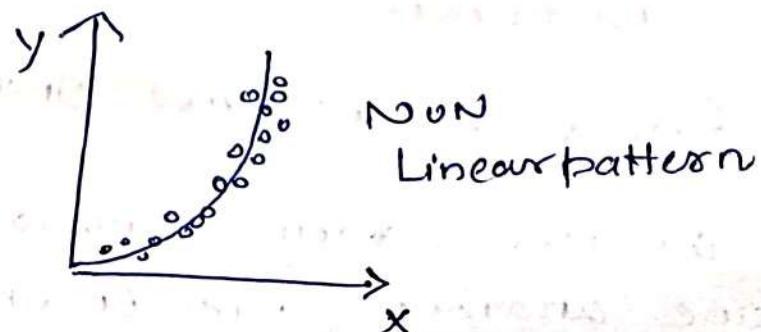
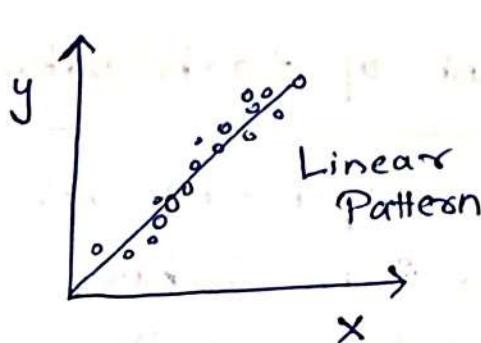
#Assumptions of Simple linear Regression

You are making conclusions on the "population" using a "sample", the assumptions that variables are linearly dependent is not enough to generalise the results you obtain on a sample to the population which is much larger in size than the sample. thus we need to have certain assumptions in place in order to make conclusions.

Let's understand the importance of each assumption one by one:

* There is a linear relationship b/w X and Y

⑥ X and Y should display some sort of linear relationship otherwise there is no use of fitting a linear model b/w them.

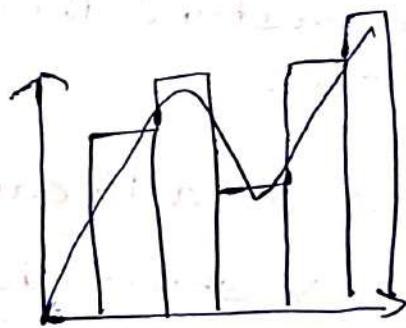


* Error terms are normally distributed with mean zero (not x, y)

- There is no problem if the error terms are not normally distributed if you just wish to fit a line and not make any further interpretation.
- But if we are willing to make some conclusions on the model that we build then we need to have a notion of the distribution of the error terms.
- The assumption of normality is made, as it has been observed that the error terms generally follow a normal distribution with the mean equal to zero in most cases.



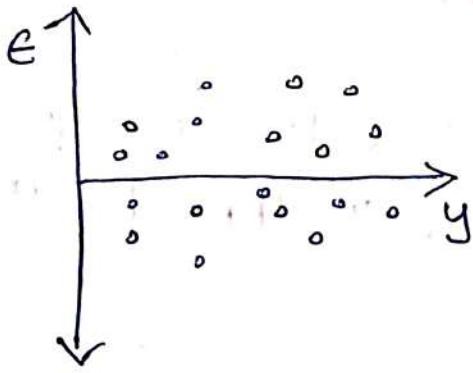
Error terms
Normally
Distributed



Error terms not
normally distributed

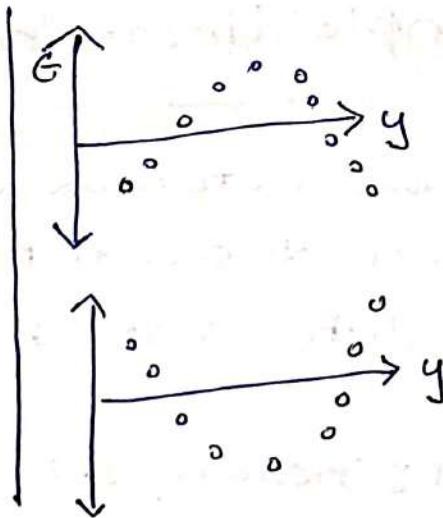
* Error terms are independent of each other

The error terms should not be dependent on one another (like in time series data where the next value is dependent on the prior one)



No visible patterns

Error terms are
Independent

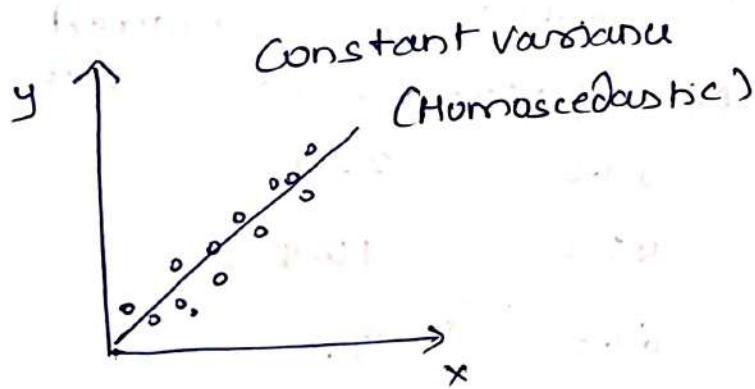


Visible pattern - Error
terms are dependent

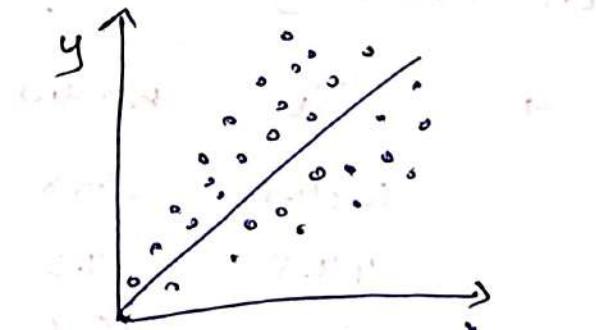


Error terms have constant variance (homoscedasticity)

- the variance should not increase (or decrease) as the error value changes
- Also the variance should not follow any pattern as the error terms change.



Constant Variance
(Homoscedastic)



changing Variance
(Heteroscedastic)

Simple Linear Regression

In this notebook we will build a linear regression model to predict Sales using an appropriate Predictor variable.

Step 1: Reading & Understanding the Data

Let's start with the following steps:

- ① Import Data using Pandas library
- ② Understanding the structure of the data

```
> import numpy as np
```

```
> import pandas as pd
```

reading the dataset

```
> advertising = pd.read_csv("advertising.csv")
```

```
> advertising.head()
```

	Tv	Radio	Newspaper	Sales	
230.1	37.8		69.2	22.1	
44.5	39.3		45.1	10.4	
17.2	45.9		69.3	12.0	
151.5	41.3		58.5	16.5	
180.8	10.8		58.4	18.9	

Here each row represents a diff market

Goal is to build the regression model which Predict sales using cols Tv, Radio, Newspaper

Checking the shape of dataset

> advertising.shape

↳ (200, 4)

checking for the missing values

> advertising.info()

↳ All 4 cols have 200 non null float64 values

Checking for Statistics value

> advertising.describe()

↳ It will show mean, median values

Visualise the dataset

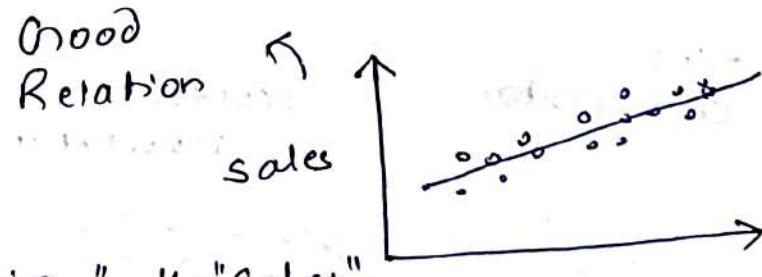
> import matplotlib.pyplot as plt

> import Seaborn as sns

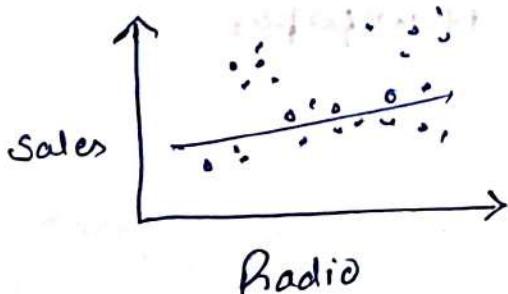
(#) Comparing "TV" & "Sales" columns

> sns.regplot(x="TV", y="Sales", data=advertising)

(#) Similarly we can
Plot



> sns.regplot(x="Radio", y="Sales",
data=advertising)

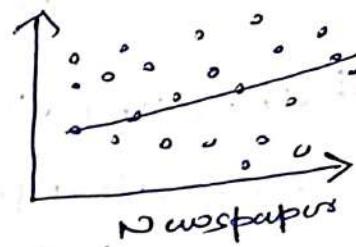


{ ↗ the data is scattered hence
Not so smooth Relation b/w
the two as compared to

TV Vs Sales

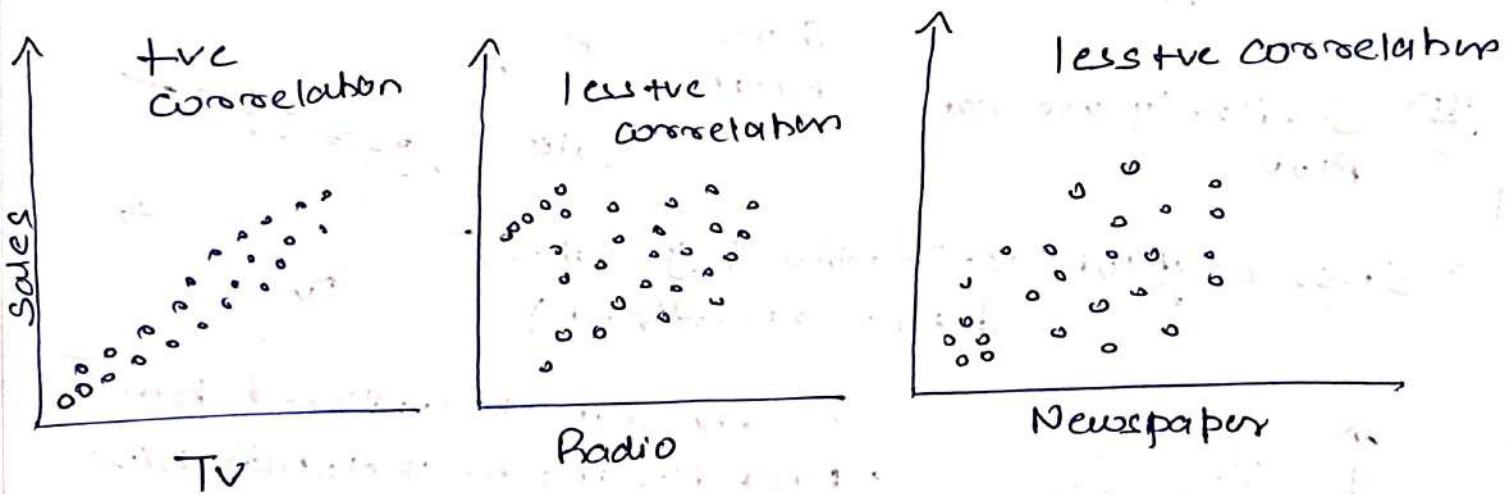
Similarly for Newspapers vs Sales column
sns. regplot (x = "Newspaper", y = "Sales",
data = advertising)

Data scattered around the line → Bad relation b/w the two.



For making a good comparison b/w the plots
↳ Pairsplots

> sns.pairplot (data = advertising, x_vars = ['TV', 'Radio', 'Newspaper'],
y_vars = 'Sales')



Another way → heatmap

Using the heatmap we want to plot the correlation

> advertising.corr()

	Tv	Radio	Newspaper	Sales
Tv	1.00	0.054	0.056	0.901
Radio	0.05	1.00	0.354	0.349
Newspaper	0.05	0.354	1.000	0.157
Sales	0.901	0.349	0.157	1.000

Corr Matrix

> sns. heatmap(advertising.corr(), annot=True)

So we read the data and visualized it using Seaborn, we also looked at the correlation b/w the target variables 'Sales' and diff predictor variables.

Hypothesis testing In Linear Model

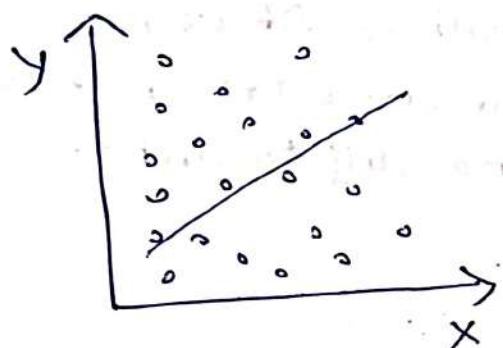
When we fit the straight line through the data, we will obviously get the two parameters of the straight line ie the Intercept (B_0) and the slope (B_1).

Now while B_0 is not of much importance right now, but there are few aspects surrounding B_1 which is needed to be checked and verified.

the first gives we ask is "Is the beta coefficient significant?" what does this mean?

⇒ Suppose we have a dataset for which the scatterplot looks like the following

Now If we run a linear regression on this dataset in Python, it will fit a line on the data which say looks like the following:



Now, we can clearly see that the data is randomly scattered and doesn't follow any linear trend or any trend or general

But Python will anyway fit a line using through the data using the least squared method, but we can see that the fitted line is of no use here

Hence everytime we perform a linear regression we need to test whether the fitted line is a significant one or not or to simply put it we need to test whether B_1 is significant or not

Here comes the idea of hypothesis testing on B_1

We start by saying that B_1 is not significant
ie there is no relationship b/w x and y.

So in order to perform the hypothesis test, we first propose the null hypothesis that B_1 is 0 and the alternate hypothesis that B_1 is not zero.

① Null Hypothesis (H_0) : $B_1 = 0$

② Alternate Hypothesis (H_A) : $B_1 \neq 0$

If we fail to reject the null hypothesis that would mean that B_1 is zero which would simply mean that B_1 is insignificant and of no use in the model.

Similarly If we reject the null hypothesis B_1 is not zero & the line fitted is a significant one.

How do we perform the hypothesis test?

↳ First we need to compute t-score which is similar to the Z score

$$Z_{\text{score}} = \frac{\bar{x} - \mu}{S/\sqrt{n}}$$

or
t score

Where $\mu \rightarrow$ Population mean
 $S \rightarrow$ Sample standard deviation

$S/\sqrt{n} \rightarrow$ Standard error

Using this, the t-score for \hat{B}_1 comes out to be
(Since the null hypothesis is that B_1 is equal to 0)

$$\frac{\hat{B}_1 - 0}{SE(\hat{B}_1)}$$

Now, In order to perform the hypothesis test,
we need to derive the P value for the given Beta,

Let's see how P value is calculated?

- Calculate the value of t-score for the mean point (in this case 0, according to the null hypothesis that we have started) on the distribution.
- Calculate the P-value from the cumulative prob for the given t-score using t-table.
- Make the decision on the basis of the p-value w.r.t to the given value of B (significance level)

Now If $P\text{value} < 0.05 \rightarrow$ reject null hypothesis
and state that B_1 is indeed significant

Q Suppose for a linear model we got B_1 as 0.5 and standard error of B_1 was found out to be 0.02 Hence t-score?

$$\text{t-score} = \frac{\hat{B}_1 - 0}{SE(\hat{B}_1)} = \frac{0.5 - 0}{0.02} = \frac{50}{2} = 25$$

Since "TV" is very strongly correlated to "Sales" let's first build a simple linear regression model with "TV" as the predictor variable.

Building a Linear Model

Building a linear model for TV as the predictor and Sales as the dependent value.

Step 2 : Performing Simple Linear Regression

Form of linear regression

$$y = C + m_1 x_1 + m_2 x_2 + \dots + m_n x_n$$

y \Rightarrow is the response var

C \Rightarrow is the intercept

$m_1 \Rightarrow$ is the coefficient for the first feature

$m_n \Rightarrow$ is the coefficient for the nth feature

In Our Case :

$$y = C + m_1 \times \text{TV} \quad \text{column}$$

the m values are called the model coefficients or model parameters

library being used "statsmodels"

> import statsmodels

> import statsmodels.api as sm

> import sklearn

we can build the model using "statsmodels" or "sklearn"

Steps In Model Building

→ Create X and Y

→ Create train and test sets (70:30 or 80:20 ratio)

→ Train model on the training set (i.e. to learn the coefficients)

→ Evaluate the model (training set, test set)

Creating X and Y

> $X = \text{advertising}[:, \text{'TV'}]$

> $Y = \text{advertising}[:, \text{'Sales']}$

Train - Test Split

↳ Sklearn comes with a built-in method for train-test-split

> from sklearn.model_selection import train-test-split

> $X_train, X_test, Y_train, Y_test = \text{train_test_split}(X, Y, \text{train_size} = 0.70, \text{random_state} = 100)$

train_size = 70%. test_size = 30%.

random_state = 100 (reshuffled data)

> $X_train.shape$

↳ (140, 1)

↳ $X_test.shape$

↳ (60, 1)

So values got split into train test

(#) Training the Model (statsmodel)

So when we are using the statsmodel library for creating a model, there is a catch there and the catch is by default the statsmodel library doesn't include "C" constant term.

So to add the constant term explicitly we need to tell the library about it

> X_train.head(3)

	Tv
4	213.4
3	151.5
185	205.0

We will create the new version of the Xtrain for the statsmodel library

> X_train_sm = sm.add_constant(X_train)

> X_train_sm.head(3)

	const	Tv
4	1.0	213.4
3	1.0	151.5
185	1.0	205.0

? If added a col called as const with all the values as 1.0 for the constant part

So our equation became

$y = C + m_1 \cdot x_1$ (original)

$y = C \cdot \text{const} + m_1 \cdot \text{Tv}$ col

So now our training set looks like

> `X_train_sm.head(3)`

	const	TV
1.0	213.4	
1.0	151.5	
1.0	205.0	

fitting the model ^{ordinary least squares}

> `lr = sm.OLS(y_train, X_train_sm)`
↳ It creates a linear regression object

> `lr_model = lr.fit()`

> `lr_model.params` # to see model parameters

	const	TV	Coefficients
	6.94	0.05	

↳ # $\text{Sales} = 6.94 + 0.05 \times \text{TV}$

↳ coefficient of TV \Rightarrow +ve

we saw that Sales \uparrow as TV \uparrow
also const is 6.94 \Rightarrow It means
if we spend 0\$ on TV \Rightarrow Sales will be
6.94

other details

> `lr_model.summary()`

↳ Only given by statsmodels

In Short

After we import the statsmodel.api we create a simple linear regression model in just few steps

we import statsmodels.api as sm

X-train-sm = sm. add-constant (X-train)

l_r = sm. OLS (y-train, X-train-sm)

l_r-model = l_r.fit ()

Here OLS stands for Ordinary least squares, which is the method that statsmodels use to fit the line, we use command add-constant so that statsmodels also fits an Intercept, if we don't use this it will ignore the const term.

> l_r-model.summary()

↳ It has info like R² values p-value

Here the top 3 things that we referred here

① Coefficients and Pvalue

② R-squared is 81.6% very high

③ P(F-statistic) is low \Rightarrow the fit is not by chance

In summary stats values

R-squared : 0.816

F-statistic : 611.2

Pprob (Fstatistic) : 1.52e-52

Coeff stderror & P > 181
for wls TV and Const

F-Statistic and Model Significance:

The F-statistic tests whether the overall regression model is significant, unlike testing the individual coefficients (β_1, β_2, \dots) for their significance, the F-statistic tests if the combination of all predictors explain the variability in dependent variable significantly.

① Purpose of F-statistic

- Determines if the model as a whole fits better than a model with no predictors (i.e. all coefficients = 0)
- Ensures that the model's overall significance isn't just due to chance.

② Prob (F-statistic)

- This is the P-value associated with F statistic
- If P-value < 0.05 → the model is significant
- If P-value > 0.05 → Model might have fit the data purely by chance.

In the eq P-value is 1.52×10^{-52} (Practically 0)
↳ Means model highly significant

Why it's useful?

- ↳ Particularly useful in multiple linear regression, where there are multiple predictors, the F-statistic helps determine if all the predictor combined are significant

R-squared

→ Measures how much variance in the dependent variable is explained by the model.

e.g. An R-squared value is 0.86 means model explains 86% of the variance, which is good.

Now we have built the model but the model evaluation is still left

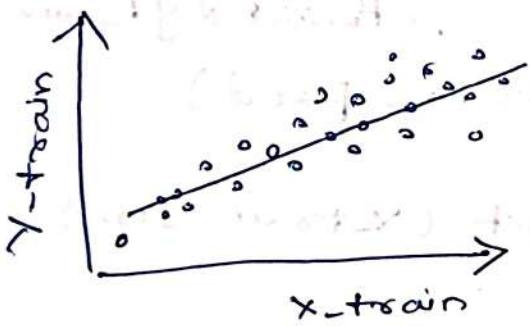
Plotting the line

> Plt. Scatter (x_{train} , y_{train})

Plt. plot (x_{train} , $6.948 + 0.054 \times x_{\text{train}}$, 'r')

($\underbrace{\phantom{0.054 \times x_{\text{train}}}_{\text{Predicted value}}}$)

Coefficients taken from summary()



Model Building Process

Residual Analysis and Predictions

Recall that one of the assumptions that you studied was the error terms should be normally distributed with mean equal to 0, so once we have built the model, you would need to verify if the model is not violating the assumptions.

for this we can plot a "histogram" of the error terms" to check whether they are normally distributed.

④ So the steps that were involved in the exercise end to end were

Step 1: Reading & understanding the Data

Step 2: Training the model

Step 3: Residual Analysis (specific to linear regression)
A fundamental assumption here is that these residuals should be normally distributed.

Step 4: Predicting and evaluating on the test set

Step 3: Residual Analysis

→ Predicted y_{train}

error = func (y_{train} , $y_{\text{train}} - \text{pred}$)

> $y_{\text{train}} - \text{pred} = \text{lr_model.predict}(X_{\text{train}} - s)$

↳ We got the predicted values

④ Now we can do error analysis

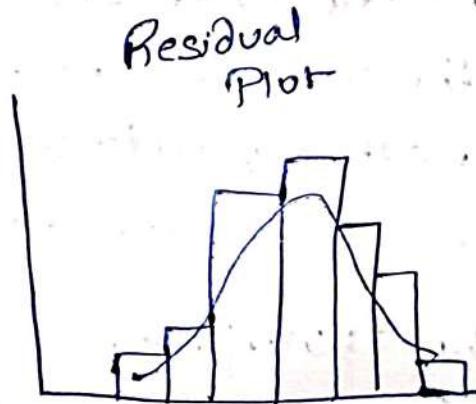
> res = $y_{\text{train}} - y_{\text{train}} - \text{pred}$

Plot the residuals (Histogram)

> sns.distplot(res)

> plt.title("ResidualPlot")

Normally Distributed

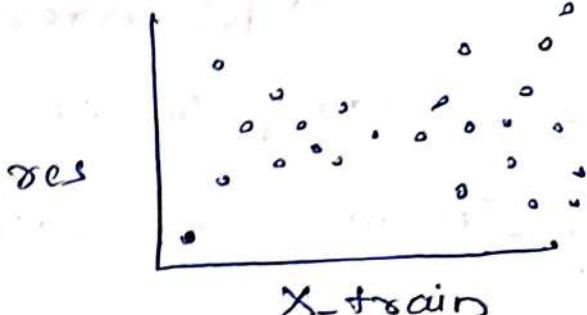


So in the first case we saw that the errors are normally distributed.

2nd thing is we look for patterns in residuals (we shouldn't be able to identify any patterns)

> plt.scatter(x_train, res)
plt.show()

(No Patterns, or clusters)



Conclusion



→ Hence the residuals are normally distributed and we couldn't identify any patterns here. So this model fit looks good.



Step 4 : Predictions and Evaluation of test set

↳ # Predictions on the test_set (y_{-test_pred})
evaluate our model, σ^2 -squared on the test

- >
- > # y_{-test_pred} is ~~fit~~ we need to add the const term
- > $x_{-test_sm} = sm.\text{add_constant}(x_{-test})$
- > $y_{-test_pred} = lr.\text{model}.\text{predict}(x_{-test_sm})$
- # Now we will evaluate the model using sklearn
- > from sklearn.metrics import mean_squared_error
- > from sklearn.metrics import r2_score

Motivation - When One Variable Isn't enough \Rightarrow Multiple Linear Regression

The term "multiple" in multiple linear regression gives us a fair idea. In itself, it represents the relationship b/w two or more independent input variables and a response variable.

Multiple Linear Regression is needed when one variable might not be sufficient to create a good model and make accurate predictions.

We know that R squared measures how well the independent variables (like TV, Radio, Newspaper) explain the variance of the dependent variable.

\rightarrow Range : 0 to 1 (or 0% to 100%)

Higher values of R^2 indicate that larger proportion of the variance in the dependent variable is explained by the independent variables.

What happens when we add more variables?

① simple linear regression (1 variable: TV)

$$\rightarrow R^2 = 0.816$$

this means 81.6% of the variability in the dependent variable (Sales) is explained by TV alone.

② Multiple Linear Regression (2 variables: TV + newspaper)

$$\rightarrow R^2 = 0.836$$

Now by adding "newspaper" as a predictor, the model explains extra 2% of the variance. So the new variable may not be as strong as TV, but it adds some predictive power.

③ Multiple Linear Regression (2 variables: TV + Radio)

$$\rightarrow R^2 = 0.910$$

→ Adding "Radio" increases the R^2 significantly. This means Radio has a strong relationship with the dependent var, helping explain 91% of variance.

* Why does R^2 increase when you add vars?

→ each new variable provides additional info that helps to explain the target.

→ R^2 always increase or stay the same when you add new independent variables.

The Multiple linear regression is just an extension of simple linear regression, hence the formulation is largely the same.

Equation of Multiple Linear Regression :

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

- ⇒ Y ⇒ Dependent variable (response / output)
- ⇒ β_0 ⇒ Intercept (value of y when all predictors are 0)
- ⇒ $\beta_1, \beta_2 \dots \beta_p$: Coefficient of predictors ($x_1, x_2 \dots x_p$)

(they measure the change in y per unit increase in the corresponding x , keeping all other predictors constant.)

- ⇒ ϵ : Error term

Interpretation of ~~constants~~ Coefficients :

- Each β represents the effect of its predictive power

Interpretation of Coefficients :

Each β represents the effect of its respective predictor X on Y , holding other variables constant
eg β_1 tells how Y changes for a unit increase in x_1 , assuming other predictors remains fixed

Apart from formulation, there are some aspects that still remain the same:

- Model now fits a hyperplane instead of a line
- Coefficients are still obtained by minimising the sum of squared errors
- Assumptions from simple linear regression still hold like zero mean, independent and normally distributed error terms with constant variance.

✳ While moving from Simple Linear Regression to Multiple Linear Regression, the new aspects to consider when moving from simple to multiple linear regression

① Overfitting

- ↳ As we keep adding the variables, the model may become far too complex.
- ↳ It may end up memorising the training data and may fail to generalize.
↳ ie. high training accuracy \rightarrow low test accuracy

② Multicollinearity

- ↳ Association b/w the predictor variables

③ Hence the feature selection becomes an important aspect.



Multicollinearity

What is Multicollinearity?

- Multicollinearity occurs when two or more independent variables in a regression model are highly correlated, in simple terms:
 - the independent variables are not truly independent
 - One variable can be predicted to a significant extent by one or more of the other variables
- this redundancy in variables leads to difficulties in interpreting the regression coefficients, which impacts the statistical inference of the model.

Why is multicollinearity a problem?

Multi collinearity affects the following:

1. Interpretation:

→ the fundamental principle of regression is the interpretation of each coefficient as the change in the dependent variable (y)

when a single independent variable (x) changes while holding others constant.

→ Multicollinearity makes the interpretation unreliable because changing one var. also causes the change in the other variable.

Conclusion from this:

- When multicollinearity is present, small changes in the data can cause large changes in the estimated coefficients.
- Multicollinearity does not affect the goodness of fit metrics (like R^2) or predictions but reduces the reliability of your regression coefficients.

How to detect multicollinearity?

1. Pairwise Correlation

- ↳ Calculate the correlation coefficients b/w pairs of independent variables.
- ↳ If two variables are highly correlated (eq correlation coefficient > 0.7) it may indicate multicollinearity.

2. Variance Inflation factor

- ↳ A more robust way to detect multicollinearity.
- ↳ VIF measures how much the variance of estimated regression coefficient is inflated due to multicollinearity.

Formula: $VIF_i = \frac{1}{1 - R_i^2}$

$$\textcircled{O} \quad VIF_i = \frac{1}{1 - R_i^2}$$

where R_i^2 is the R^2 obtained by regressing the i -th independent variable on all other independent variables.

↳ Heuristic:

$VIF < 5$: Low multicollinearity (Good to use)

$5 \leq VIF \leq 10$: Moderate multicollinearity
(Inspect further)

$VIF \geq 10$: High Multicollinearity
(Consider eliminating or transforming the variables)

Why VIF is preferred over pairwise correlation

- Pairwise correlation only evaluates relationship b/w two variables at a time.
- VIF accounts for multicollinearity caused by combination of variables providing a more comprehensive view.

Some methods that can be used to deal with multicollinearity are

- Dropping Variables → Drop vars which are highly correlated with others
- Create new vars using the combination of old vars
- PCA

* Dealing with Categorical Variables

When we have a categorical variable with say 'n' levels, the idea of dummy variable creation is to build 'n-1' variables indicating the levels, for a variable say, 'Relationship' with three levels namely → Single, In a relationship, Married we would create a dummy table like the following:

Relationship status	Single	In a relationship	Married
Single	1	0	0
In a relationship	0	1	0
Married	0	0	1

But you can clearly see that there is no need of defining those different levels, if you drop a level say "Single" you would still be able to explain the three levels. (Hence Dropping Dummy vars)

Relationship status	In a relationship	Married
Single	0	0
In a relationship	1	0
Married	0	1

How It Impacts the eqn?

Value	Indicator Variable
Gender	Female
Male	0
Female	1

$$x_i = \begin{cases} 1 & \text{if } i^{\text{th}} \text{ person is female} \\ 0 & \text{if } i^{\text{th}} \text{ person is male} \end{cases}$$

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i = \begin{cases} \beta_0 + \beta_1 + \epsilon_i & \text{if } i^{\text{th}} \text{ person is female} \\ \beta_0 + \epsilon_i & \text{if } i^{\text{th}} \text{ person is male} \end{cases}$$

In short : $N - 1$ dummy variables can be used to describe a categorical variable with N variables.

Feature Scaling

Why do we need to scale features?

- ↳ ① Ease of Interpretation
- ② Faster Convergence for gradient descent methods



Scaling in Regression Models

Why is Scaling Important?

Scaling ensures that variables with different units or magnitudes are on the same scale, which is particularly important in the following scenarios:

1. **Algorithms Sensitive to Magnitude**: Algorithms like gradient descent (used in linear regression, logistic regression, etc) or regularization techniques (Lasso and Ridge etc) converge faster when the variables are scaled.
2. **Interpretability of Coefficients**: While scaling changes the magnitude of the coefficients, it does not affect statistical measures like t-statistic, P-value, or R^2 .
3. **Multicollinearity**: Scaling does not eliminate the multicollinearity but helps normalize variables for regularization based solns.

Key Parameters Affected by Scaling

- Coefficients : Scaling changes their magnitude, but not their relationship with the dependent variable.
- Predictions remain unaffected.
- Statistical Metrics such as t-statistic, f-statistic, P-value, R² remains unchanged regardless of the scaling methods.

Methods of Scaling

① Standardization

- ↳ transforms the data into standard normal distribution where

$$\begin{aligned} \text{Mean} &= 0 \\ \text{Standard deviation} &= 1 \end{aligned}$$

$$\text{formula } \Rightarrow x' = \frac{x - \text{mean}(x)}{\text{std}(x)}$$

use case :

- ↳ Preferred when data is normally distributed

② Min Max Scaling

- ↳ Rescales the data to lie within a fixed range usually [0,1] formula

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

useful when you need the data in a specific range

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$



Model Assessment and Comparison

When building regression models, there is often a trade-off between keeping the model simple (with fewer Predictor variables) and explaining the highest possible variance (which means including as many relevant predictors as possible). To strike this balance 2 additional metrics come into play

- ↳ Adjusted R^2 and Akaike Information Criterion (AIC)

Adjusted R^2

↳ It is a modified version of R^2 that adjusts for the number of predictors in the model.

Unlike R^2 , which always increases when a new Predictor is added, Adjusted R^2 considers whether the added variable improves the model enough to justify its inclusion.

$$\text{Adjusted } R^2 = 1 - \left(\frac{(1-R^2)(N-1)}{N-p-1} \right)$$

N → No. of rows

p → No. of Predictor variables

Higher Adjusted R^2 indicates a better model. It ensures only meaningful variables are included.

Akaike Information Criterion (AIC)

→ AIC is a metric used to compare diff models considering both the goodness of fit (how well the model fits the data) and the complexity (number of predictors), Lower AIC value indicates better models.

$$AIC = n \times \log\left(\frac{RSS}{n}\right) + 2p$$

n → nos of observations

RSS → Residual sum of squares

p → nos of predictors

Eg's

Suppose we are building a regression model to predict housing prices and you have 3 candidate models:

Model	Predictors	R ²	Adj R ²	AIC
M ₁	Lotsize, beds	0.25	0.24	210
M ₂	Lotsize, beds, garage	0.38	0.36	190
M ₃	All cols	0.80	0.73	250

M₂ Model Best

$$Q \quad n=101, \quad p=10, \quad R^2 = 0.7$$

$$\text{Adjusted } R^2 = 1 - \frac{(1-0.7)(101-1)}{101-10-1}$$

$$1 - \frac{(0.3)(100)}{90} \cdot 1 - \frac{30}{90} \cdot 1 - \frac{1}{3} = \frac{2}{3}$$

In the context of Adjusted R^2

- ↳ Adding More Predictors (more independent variables) always increases R^2 , even if those predictors don't actually help the model.
- Adjusted R^2 penalize (reduces its value) if we add the predictors that don't increase model's value.

Key Difference b/w the R^2 and the Adjusted R^2

① R-squared (R^2)

$$\text{formula : } R^2 = 1 - \frac{RSS}{TSS}$$

where $RSS \Rightarrow$ Residual sum of squares

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$TSS \Rightarrow$ total sum of squares $\sum_{i=1}^n (y_i - \bar{y})^2$

$y_i \Rightarrow$ Actual value of dependent var

$\hat{y}_i \Rightarrow$ Predicted value of the dependent var

$\bar{y} \Rightarrow$ Mean of actual values

Keypoints :

- R^2 tells how much of the total variation in the dependent variable is explained by the independent variables
- Adding more variables to the model always increases or maintains the R^2 value, even if the new variable has no relationship with the dependent variable, because it ignores the number of predictors.

Adjusted R^2 (R^2_{adj})

formula : $R^2_{adj} = 1 - \left(\frac{(1-R^2)(n-1)}{n-p-1} \right)$

$R^2 \Rightarrow R^2_{adjusted}$ value

$n \Rightarrow$ number of observations (sample size)

$p \Rightarrow$ Number of Predictors (independent vars) in the model.

Keypoints :

- R^2_{adj} adjusts for the number of predictors in the model
- If we add an independent var that does not significantly improve the model $\rightarrow R^2_{adj}$ will decrease
- It is more reliable than R^2 for multiple linear regression because it accounts for overfitting

* Feature Selection

feature selection is a crucial step in building an effective linear regression model. It involves identifying the most relevant independent vars (predictors) that contribute to the dependent variable (outcome) while avoiding overfitting or redundancy.

Why is feature selection important?

① It simplifies the model

↳ It reduces the nos of predictors, making the model easier to interpret.

→ Helps in avoiding overfitting by excluding irrelevant or redundant variables.

→ Reduces Multicollinearity

↳ Prevents correlated variables from distorting the result.

Approaches to feature Selection

① Exhaustive Search :

↳ Test all the possible combinations of the Predictors (2^P models for P features).

Computationally expensive, and impractical for large datasets.

2. Manual feature elimination

↳ works for smaller nos of predictors (eg 10-20)

Steps :

① Build the model with all features

↳ fit the regression model with all the Potential Predictors.

② Drop features with high P values:

↳ Remove features that are statistically insignificant (eg $P > 0.05$)

③ Remove Redundant features

→ use correlation analysis and variance inflation factor (VIF) to identify multicollinear variables and eliminate them.

→ VIF rule of thumb : Drop features with $VIF > 5$ or 10

④ Rebuild the Model and Repeat

↳ iteratively refine the model until only significant and non redundant features remain.

3. Automated feature Selection :

→ Useful when dealing with large nos of Predictors (eg 100 or more)

Techniques :

① Recursive Feature Elimination (RFE)

- ↳ Iteratively build the model and eliminate the least important features based on a scoring metric.

② Stepwise Selection:

↳ Forward Selection:

- ↳ Start with no predictors and add them one at a time based on their significance.

↳ Backward Elimination

- ↳ Start with all predictors and remove them one at a time based on their P-values or other metrics.

↳ Bidirectional / Stepwise Selection

- ↳ Combines forward and backward selection.

RFE In Detail + Graded Ques Study

RFE stands for Recursive Feature elimination, is a feature technique used in machine learning to select the most important feature by recursively removing the least important ones, It helps in improving the model performance , reduce overfitting and speeding up the model training process by eliminating the irrelevant features.

Here's a breakdown of how RFE works

RFE works by fitting a model to the data and ranking features based on their importance, the model assigns weights to features, indicating how influential each feature is in predicting target variable.

→ Recursive Elimination

After the Model training is completed, RFE Recursively eliminates the least important features, this process is repeated until the specified number of features remains, In each iteration the least imp feature is removed and the model is retrained.

→ Ranking features:

RFE uses model's performance to evaluate the imp of each feature, features with the least weights or those contributing the least to the model's prediction are eliminated first.

Example with Linear Regression :

↳ we want to predict a target variable using linear regression

let's say we have the Dataset as follows:

feature1 (x_1)	Feature2 (x_2)	Feature3 (x_3)	Target (Y)
5	10	100	200
6	12	120	220
7	14	130	240
8	16	150	250

Step 2: Train Initial Model

↳ you first train a linear regression model using all the features (x_1, x_2, x_3), the model will assign weights to each feature based on how they affect the target variable Y.

Step 3: Train Initial Model

↳ you first train a linear regression model using all features (x_1, x_2, x_3) → the model will assign weights to each feature based on how they affect the target variable Y.

Step 3: Evaluate Feature Importance

RFG uses the weights of the features assigned by the linear regression model to rank their importance.

If x_1 has a weight of 0.8, x_2 has 1.2 and x_3 has 0.1, it means x_2 is the most important feature while x_3 is the least important.

Step 4: Eliminate Least Importance Feature

Since x_3 is the least important, RFE will remove it from the dataset, now you're left with just x_1 and x_2 .

Step 5: Retrain the Model

RFE then retrains the model with only the remaining features (x_1, x_2), the model will reassign weights to these features and repeat the process.

Step 6: Repeat until Desired Number of features remains

This process continues recursively, eliminating one feature at a time based on importance, until the specified no. of features is reached.

Advantages of RFE

- ① Improve model performance by removing unnecessary features
- ② It provides a ranking of features, helping to understand which ones are most impactful.

Disadvantages:

- Computationally expensive, Requires a model

Code example with Linear Regression

```
from sklearn.feature_selection import RFE  
from sklearn.linear_model import LinearRegression  
import pandas as pd
```

```
# eq dataset
```

```
data = pd.DataFrame({  
    'X1': [5, 6, 7, 8],  
    'X2': [10, 12, 14, 16],  
    'X3': [100, 120, 130, 150],  
    'Y': [200, 220, 240, 250]  
})
```

```
# features and target variables
```

```
X = data[['X1', 'X2', 'X3']]
```

```
y = data['Y']
```

```
# Initialize the model
```

```
model = LinearRegression()
```

```
# Initialize Model in RFE and add the nos of features  
to select
```

```
rfe = RFE(model, 2) # Select top 2 features
```

```
# fit RFE
```

```
rfe.fit(X, y)
```

Get the features selected by RFE

Selected-features = X.columns [xfe. support_]
Point (selected_features)

Notes On Multiple Linear Regression

Summary

Multiple Linear Regression (MLR) is a statistical technique used to predict the values of a dependent var Y based on multiple independent variable X_1, X_2, \dots, X_n , It extends simple linear regression which involve one independent variable.

$$Y = B_0 + B_1 X_1 + B_2 X_2 + \dots + B_n X_n + \epsilon$$

where : $Y \Rightarrow$ Dependent var (Target)

$X_1, X_2, \dots, X_n \Rightarrow$ Independent var

$B_0 \Rightarrow$ Intercept

$B_1, B_2, \dots, B_n \Rightarrow$ Coefficients of independent variables

$\epsilon \Rightarrow$ Error term (residual)

Assumptions of MLR :

- ① Linearity : the relationship b/w dependent and independent var is linear.
- ② errors should be independent
- ③ Homoscedasticity : variance of residuals should be const across all levels

④ Independent variables should not be too highly correlated.

⑤ Residuals should be normally distributed.

* Solving MLR Using an example

Problem statement : Predict the Price of house (Y) using its size in sq feet (x_1) and number of bedrooms (x_2)

Step 1 : Data

House Size (x_1)	Number of Bedrooms (x_2)	Price (Y)
1500	3	40,00,000
2000	4	55,00,000
2500	4	62,00,000
1800	3	48,00,000
2300	5	65,00,000

Step 2 : Model

$$y = B_0 + B_1 x_1 + B_2 x_2 + \epsilon$$

Assume $B_0 = 20,000$, $B_1 = 10$, $B_2 = 50,000$

Step 3 : Prediction

To predict the price of a house with size $x_1 = 2100 \text{ sq ft}$ and $x_2 = 4 \text{ bedrooms}$

$$y = 20,000 + 10(2100) + 50000(4)$$

$$\Rightarrow \boxed{y = 2,41,000}$$

Metrics to evaluate Model Performance

- ① R-squared (R^2): Proportion of variance in Y explained by x_1, x_2, \dots, x_n
- ② Adjusted R-squared: Penalizes for the addition of irrelevant predictors
- ③ Mean Squared Error (MSE): Avg Squared diff b/w predicted and actual values
- ④ Root Mean Squared Error (RMSE): Square root of MSE, same units as Y .

✳ Model Performance Metrics

To evaluate how well a multiple linear regression model fits the data, we use the following metrics

① R-squared (R^2)

It measures the proportion of variance in the dependent var(Y) that is explained by the independent vars (x_1, x_2, \dots, x_n)

↳ Range from 0 to 1, closer to 1 indicates a better model fit.

$$\text{formula} \Rightarrow R^2 = 1 - \frac{\text{RSS}}{\text{TSS}} \quad \begin{array}{l} \rightarrow \text{Residual Sum of Squares} \\ \downarrow \end{array}$$

↳ Total sum of squares

eq

We have the following values for the dependent variable and it's corresponding predicted values (\hat{Y})

Actual (Y)	Predicted (\hat{Y})
50	48
60	63
55	54
70	69

$$RSS = \sum (y - \hat{y})^2 \quad | \quad TSS = \sum (y - \bar{y})^2$$

Here \bar{y} is the mean of Y $\Rightarrow (50+60+55+70)/4$

$$TSS = (50 - 58.75)^2 + (60 - 58.75)^2 + (55 - 58.75)^2 + (70 - 58.75)^2 = 219.75$$

$$RSS = (50 - 48)^2 + (60 - 63)^2 + (55 - 54)^2 + (70 - 69)^2 = 15$$

$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{15}{219.75} = 0.9318$$

So R^2 explains 93.18% of the variance in Y.

② Adjusted R squared

Adding irrelevant predictors can artificially increase R^2 , even if the predictors have no real influence, Adjusted R^2 corrects for this.

$$R^2_{adj} = 1 - \left(\frac{(1-R^2)(n-1)}{n-p-1} \right)$$

where $n \Rightarrow$ nos of observations

$P \Rightarrow$ nos of predictors

eg

Suppose $R^2 = 0.9$, nos of predictors = 2
and the nos of observations $n=10$

$$R^2_{adj} = 1 - \frac{(1-0.9)(10-1)}{10-2-1}$$

$$\Rightarrow R^2_{adj} = 1 - \frac{(0.1)(9)}{7} = 0.8714$$

③ Mean Squared Error (MSE)

↳ Avg squared diff b/w actual and predicted values.

formula: $\Rightarrow MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$

Here $n \Rightarrow$ nos of observations

eg using the earlier table of Actual & Predicted values

$$MSE = \frac{(50-48)^2 + (60-63)^2 + (55-54)^2 + (70-69)^2}{4}$$

$$\frac{4+9+1+1}{4} = \frac{15}{4} = \underline{\underline{3.75}}$$