

# Assignment No. 1

Page No.	
Date	

Title: Linear Regression Analysis :- Driving hours vs Backache Risk.

Problem Statement:

We are given small dataset containing the number of hours spent driving ( $x$ ) and a corresponding risk score ( $y$ ) for developing acute backache.

Our goal is to:

- Analyze the data using simple linear regression.
- Derive the equation of the best-fit-line.
- Interpret the results and their real world implications.

Number of hours spent driving ( $x$ )	Risk score on a scale of 0-100 ( $y$ )
---------------------------------------	--

10	95
----	----

9	80
---	----

2	10
---	----

15	50
----	----

10	45
----	----

16	98
----	----

11	38
----	----

16	93
----	----

Independent Variable:  $x$  is the no. of hours spent driving  
dependent variable:  $y$  is the risk score for backache.

## Analysis and Consumption

### 1. Assumptions:

Before proceeding, we assume:

- A linear relationship exists between driving hours and risk score.
- No significant outliers affect the results.
- The dataset is representative of the population.

### 2. Linear Regression formula:

$$y = mx + c$$

$m$  is slope of line

$c$  is th.  $y$ -intercept

Calculation of  $m$  and  $c$

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

Calculations:

$x$	$y$	$x^2$	$y^2$	$xy$
10	95	100	9025	950
9	86	81	7396	720
2	10	4	100	20
15	50	225	2500	750
10	45	100	2025	450
16	98	256	9604	1568
11	38	121	1444	418
16	93	256	8649	1488

Now sum each column

$$\sum x = 8g$$

$$\sum y = 50g$$

$$\sum x^2 = 1143$$

$$\sum xy = 636g$$

$$n = 8$$

$$m = \frac{8(636g) - (8g)(50g)}{8(1143) - (8g)^2} = \frac{50g(12) - 4520g}{9144 - 792g} = \frac{570g}{1223} = 4.66$$

$$c = 50g - (4.66)(8g) = \frac{50g - 414.75}{8} = 94.26 \approx 11.78$$

Final Regression formula is given by  $y = 4.66x + 11.78$

$$y = 4.66x + 11.78$$

This is best fit line describing the relationship between driving hours and backache risk.

Interpretation :-

Slope (m) = 4.66 :- For every additional hour spent driving, the risk score increases by approximately 4.66 units.

Intercept (c = 11.78) :- Even at 0 hours of driving, a baseline risk score of around 11.78 exists, possibly due to other factors like posture or lifestyle.

### Prediction:

If a person drives 12 hours a day:

$$y = 4.66(12) + 11.78 = 66.7 + 11.78 = 78.48$$

So, the predicted risk score is approximately 78.5

Conclusion:

from our analysis, there is a positive linear correlation between the number of driving hours and the risk of developing acute backache. This implies that the no. of hours spent driving increases, the backache risk also increases significantly.

$$87.11 + 4.66x = y$$

# Assignment No. 2

Page No.	
Date	

Title :

Aim: Apply the principal Component Analysis for feature reduction on IRIS Dataset

Objective: The objective is to reduce the number of features in a dataset while retaining the most important information. PCA helps simplify data visualization and speeds up processing without losing significant patterns.

Theory Overview:

- PCA is an unsupervised machine learning algorithm used for dimensionality reduction.
- It transforms correlated variables into a new set of uncorrelated variables called principle components.

Why PCA ?

- High dimensional data can be noisy and computationally expensive.
- PCA helps by compressing data into fewer dimensions while preserving trends and patterns.

Key Concepts :

Covariance Matrix : Measures - how variables move together

Eigen values / Eigen vectors : Identify principle directions of variance.

Orthogonality : Principal components are perpendicular (uncorrelated)

Dimensionality : Number of variables or features in the dataset.

## PCA Algorithm Steps

1. Getting the dataset
  - Firstly we need to take the input dataset and divide it into two subparts  $X$  and  $Y$ , where  $X$  is the training set and  $Y$  is the validation set.
2. Standardize the data
  - Standardize our dataset with the particular column, the features with high variance are more important compared to the features with lower variance.
  - Make features have mean = 0 and unit variance.
3. Calculate the covariance matrix
  - Represents relationships between features.
4. Compute Eigenvalues and Eigenvectors
  - Identify directions with the most variance
5. Sort Eigenvalues and select top-k
  - keep the most significant components
6. Transform the Data
  - Multiply standardized data with top-k eigenvectors to get the reduced dataset.
7. Drop unimportant features
  - keep only significant principal components.

## Applications of principal Component Analysis :-

- PCA is mainly used as the dimensionality reduction technique in various AI applications such as computer vision, image compression, etc.
- It can also be used for finding hidden patterns if data has high dimensions. Some fields where PCA is used are finance, data mining, Psychology, etc.

### Implementing PCA on a 2-D Dataset

Step 1: Normalize the data

This is done by subtracting the respective means from the numbers in the respective column.

Step 2: Calculate the covariance matrix

Since the dataset we took is 2-dimensional, this will result in a  $2 \times 2$  covariance matrix

$$\text{Matrix (Covariance)} = \begin{bmatrix} \text{Var}[x_1] & \text{Cov}[x_1, x_2] \\ \text{Cov}[x_2, x_1] & \text{Var}[x_2] \end{bmatrix}$$

Step 3: Calculate the eigenvalues & eigenvectors

Step 4: Choosing components and forming a feature vector.

## Output:

- Standardized data
- PCA projection (2D scatter plot)
- Visualized reduced - dimension dataset

Conclusion: PCA was successfully implemented on the IRIS dataset for feature reduction.

The dimensionality of the data was reduced while preserving its important characteristics, enabling easier visualization and analysis.

$$\begin{bmatrix} [x_1, x_2]_{n \times 2} & [x_3]_{n \times 1} \\ [x_1]_{n \times 1} & [x_2, x_3]_{n \times 2} \end{bmatrix} \rightarrow (\text{min}(n)) \times n \times n$$

# Assignment No. 3

Page No.	
Date	

Title: Assignment on decision tree Classifier

Problem Statement: dataset collected in a cosmetics shop showing details of customers and whether or not they responded to a special offer to buy a new lip-stick is shown in table below.

Use this dataset to build a decision tree, with Buys as the target variable to help in buying lip-sticks in the future.

Find the root node of decision tree. According to the decision tree you have made from previous training data set, what is the decision for the test data?

[Age <21, Income = Low, Gender = female, Marital Status = Married]

ID	Age	Income	Gender	Marital Status	Buys
1	<21	High	M	Single	No
2	<21	High	M	Married	No
3	21-35	High	M	Single	Yes
4	>35	Medium	M	Single	Yes
5	>35	Low	F	Single	Yes
6	>35	Low	F	Married	No
7	21-35	Low	F	Married	Yes
8	<21	Medium	M	Single	No
9	<21	Medium	Low F	Married	Yes
10	>35	Medium	F	Single	Yes
11	<21	Medium	F	Married	Yes
12	21-35	Medium	M	Married	Yes
13	21-35	High	F	Single	Yes
14	>35	Medium	M	Married	No

## Relevant Theory:

- Decision Tree is a decision making tool that uses a flowchart-like tree structure or is a model of decisions and all of their possible results, including outcomes, input costs and utility.
- Supervised learning algorithms find the best fit for both continuous as well as categorical output variables.
- branches (edges) represent the result of the nodes and the nodes have either:
  1. Conditions [Decision Nodes]
  2. Result [End Nodes]

## Steps:

- ① Compute Entropy For Data set
- ② Which Node to select as Root
  - Find maximum Gain As Root

Age has maximum Gain therefore Age become Root node.

## Source Code:

```
#import packages
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
```

```
# reading dataset
```

```
dataset = pd.read_csv('data.csv')
```

```
X = dataset.iloc[:, :-1]
```

```
y = dataset.iloc[:, 5].values
```

```
# perform label encoding
```

```
from sklearn.preprocessing import LabelEncoder  
labelencoder_X = LabelEncoder()
```

```
X = X.apply(LabelEncoder().fit_transform)
```

```
print(X)
```

```
from sklearn.tree import DecisionTreeClassifier
```

```
regressor = DecisionTreeClassifier()
```

```
regressor.fit(X.iloc[:, 1:5], y)
```

```
# predict value for given expression
```

```
X_in = np.array([1, 1, 0, 0])
```

```
y_pred = regressor.predict([X_in])
```

```
print("Prediction:", y_pred)
```

```
from sklearn.externals.six import StringIO
```

```
from IPython.display import Image
```

```
from sklearn.tree import export_graphviz
```

```
import pydotplus
```

```
# create DOT data
```

```
dot_data = StringIO()
```

`export_graphviz (regressor, out_file = dot_data, filled = True,  
special_characters = True)`

# Draw graph

```
graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
graph.write_png('Decision-Tree.png')
```

# Show graph

```
Image(graph.create_png())
```

Output:

$$x_0 \leq 0.5$$

$$\text{gini} = 0.459$$

$$\text{Samples} = 14$$

$$\text{Value} = [5, 9]$$

True / False

$$x_2 \leq 0.5$$

$$\text{Samples} = 4$$

$$\text{Value} = [0, 4]$$

$$\text{Samples} = 10$$

$$\text{Value} = [5.5]$$

$$x_3 \leq 0.5$$

$$\text{gini} = 0.32$$

$$\text{Samples} = 5$$

$$\text{Value} = [1, 4]$$

$$x_0 \leq 1.5$$

$$\text{gini} = 0.32$$

$$\text{Samples} = 5$$

$$\text{Value} = [4, 1]$$

$$x_0 \leq 1.5$$

$$\text{gini} = 0.444$$

$$\text{Samples} = 3$$

$$\text{Value} = [1, 2]$$

$$\text{gini} = 0.0$$

$$\text{Samples} = 2$$

$$\text{Value} = [0, 2]$$

$$\text{gini} = 0.0$$

$$\text{Samples} = 3$$

$$\text{Value} = [3, 0]$$

$$x_2 \leq 0.5$$

$$\text{gini} = 0.5$$

$$\text{Samples} = 2$$

$$\text{Value} = [1, 1]$$

$$\text{gini} = 0.0$$

$$\text{Samples} = 2$$

$$\text{Value} = [0, 2]$$

$$\text{gini} = 0.0$$

$$\text{Samples} = 1$$

$$\text{Value} = [1, 0]$$

$$\text{gini} = 0.0$$

$$\text{Samples} = 1$$

$$\text{Value} = [1, 0]$$

$$\text{gini} = 0.0$$

$$\text{Samples} = 1$$

$$\text{Value} = [0, 1]$$

# Assignment No. 4.

Page No.	
Date	

Title: Assignment on Naïve Bayes Classifier Notes

Problem Statement: Implement Naïve Bayes Classification Algorithm on a given dataset

Objective: 1. Study the concept of Naïve Bayes Algorithm  
2. Implement the Naïve Bayes Algorithm on Iris.csv dataset

Relevant Theory:

- Supervised machine learning algorithm
- Based on bayes theorem and used for solving classification problems
- Mainly used in text classification that includes a high dimensional training dataset.
- Probabilistic model / classifier, it predicts on the basis of the probability of an object.  
ex. spam filtration, Sentimental analysis and classifying articles

Formula for Bayes theorem is  $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$

Where,  $P(A|B)$  → posterior probability : Probability of hypothesis A on observed event B.

$P(B|A)$  → likelihood probability : probability of the evidence given that the probability of a hypothesis is true.

$P(A)$  → Prior Probability : probability of hypothesis before observing the evidence.

$P(B)$  → Marginal Probability : Probability of Evidence.

## Algorithm:

- Data preprocessing step
- Fitting Naive Bayes to the Training set
- Predicting the test result
- Test accuracy of the result (creation of confusion matrix)
- Visualizing the test set result

## Advantages:

- One of the fast and easy ML algorithms to predict a class of datasets.
- Can be used for Binary as well as Multiclass classifications
- Performs well in Multi-class predictions as compared to the other algorithms.
- Popular choice for text classification problems.

## Disadvantages:-

- Naive Bayes assumes that all features are independent or unrelated, so it cannot learn the relationship between features.

## Applications of:-

- Used for credit Scoring
- Used in medical data classification
- Used in real-time predictions because Naive Bayes classifier is an eager learner.
- Used in Text classification such as Spam filtering and Sentiment analysis.

Output:

Correct Predictions : 28

False Predictions : 2

Accuracy of the Naïve Bayes Classification is 0.933333333333333

Conclusion:

Studied about Naïve Bayes Classification and implemented it on Iris.csv dataset for classification.

# Assignment No. 5

Page No.	
Date	

Title: Assignment on SVM on any dataset

Objective: Implement the SVM on any dataset

Input: csv file containing algorithm and its output

Relevant Theory: We want to find out two lines or boundaries

- Supervised learning algorithms is programmed to do this
- Used for classification as well as regression problem
- Goal: To create the best line or decision boundary that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future.
- The best decision line is called hyperplane.

Types of SVM :-

① Linear SVM :-

- Used for linearly separable data
- A dataset can be classified into two classes by using a single straight line, then such data is termed as linearly separable data and classifier is used called as linear SVM classifier.

② Non-linear SVM :-

- Used for non-linearly separated data
- A dataset cannot be classified by using a straight line, then such data is termed as non-linear data and classifier used is called as Non-linear SVM classifier.

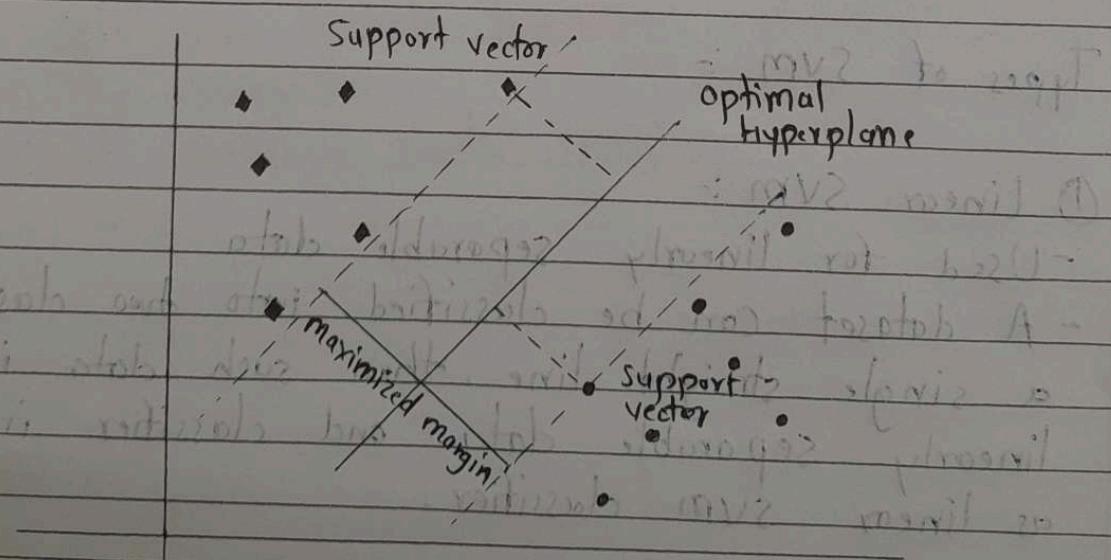
## Hyperplane and Support Vectors in SVM algorithm :-

### Hyperplane :

There can be multiple lines / decision boundaries to segregate the classes in n-dimensional space, but we need to find out the best decision boundary that helps to classify the data points.

This best boundary is known as hyperplane of SVM.

Support vectors is a set of data points that are closest to the hyperplane and which affect the position of the hyperplane are termed as Support Vector.



## Algorithm Steps

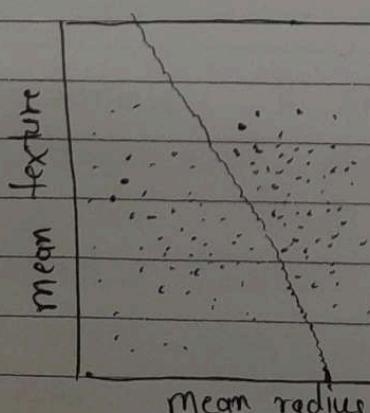
- Load the breast cancer dataset from sklearn datasets
- Separate input features and target variables
- Build and train the SVM classifiers using RBF kernel
- Plot the scatter plot of the input features
- Plot the decision boundary
- Plot the

## Applications of SVM classifier:-

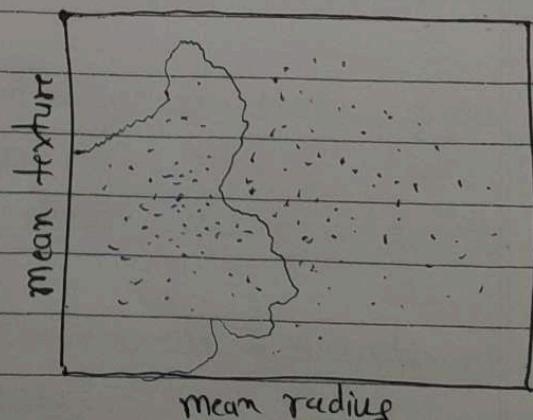
- Face detection
- Text and hypertext categorization
- Classification of images
- Bioinformatics
- Protein fold and remote homology detection
- Handwriting recognition
- Generalized predictive control (GPC)

### Output:

Linear Kernel



RBF-Kernel



## Conclusion:

studied about SVM classification and implemented it on breast-cancer dataset for classification.

# Assignment No. 6

Page No.	
Date	

Title: Assignment on k-means clustering

Problem Statement:

We have given a collection of 8 points.

$$P_1 = [0.1, 0.6] \quad P_2 = [0.15, 0.7] \quad P_3 = [0.08, 0.9] \quad P_4 = [0.16, 0.85]$$

$$P_5 = [0.2, 0.3] \quad P_6 = [0.25, 0.5] \quad P_7 = [0.24, 0.1] \quad P_8 = [0.3, 0.2]$$

Perform the k-means clustering with initial centroids as  
 $m_1 = P_1$  = cluster #1 =  $c_1$  and

$m_2 = P_8$  = cluster #2 =  $c_2$ , Answer the following

1) Which cluster does  $P_6$  belongs to?

2) What is the population of cluster around  $m_2$ ?

3) What is the updated value of  $m_1$  and  $m_2$ ?

Relevant Theory

- K-Means Clustering is an unsupervised machine learning algorithm
- Used to partition a dataset into  $k$  distinct clusters based on distance (usually Euclidean).

The algorithm works in iterative steps:-

- ① Initialize  $k$  centroids randomly or with given points.
- ② Assign each data point to the nearest centroid.
- ③ Update centroids as the mean of all points in the respective cluster.
- ④ Repeat steps 2-3 until convergence (no change in centroids or cluster assignments).

Euclidean Distance Formula :-

$$d(p, q) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Algorithm Steps :-

- ① Initialize the centroids as  $m_1 = p_1$  and  $m_2 = p_2$
- ② for each point  $P_i$ , compute the distance to both  $m_1$  and  $m_2$ .
- ③ Assign the point to the cluster with the nearest centroid.
- ④ After assigning all points, compute the new centroid of each cluster.
- ⑤ Repeat until centroids do not change (only 1 iteration is asked here)

Implementation Steps :

Initial Centroids:

$$\cdot m_1 = [0.1, 0.6]$$

$$\cdot m_2 = [0.3, 0.2]$$

Step 1 : Calculate Distance of  $p_6$  from  $m_1$  and  $m_2$

$$\therefore p_6 = [0.25, 0.5]$$

$$\begin{aligned} d(p_6, m_1) &= \sqrt{(0.25 - 0.1)^2 + (0.5 - 0.6)^2} \\ &= \sqrt{0.0225 + 0.01} \\ &= \sqrt{0.0325} \\ &\approx 0.18 \end{aligned}$$

$$d(P_6, m_2) = \sqrt{(0.25 - 0.3)^2 + (0.5 - 0.2)^2} \\ = \sqrt{0.0925} \approx 0.304$$

$P_6$  is closer to  $m_1$ , so it belongs to cluster 1 ( $C_1$ )

Step ② :- Assign All points to clusters (First iteration)

Point	Distance to $m_1$	Distance to $m_2$	cluster
$P_1$	0.108	0.5	$C_1$
$P_2$	0.1118	0.5831	$C_1$
$P_3$	0.3041	0.7071	$C_1$
$P_4$	0.254	0.6503	$C_1$
$P_5$	0.3162	0.1515	$C_2$
$P_6$	0.1803	0.3041	$C_1$
$P_7$	0.509	0.1414	$C_2$
$P_8$	0.5	0.1414	$C_2$

Step ③ : Answering the questions

① Which cluster does  $P_6$  belong to?

→ cluster 1 ( $C_1$ )

② What is the population of the cluster around  $m_2$  (cluster 2)?

→  $P_5, P_7, P_8$  → 3 points

③ Updated centroids

for  $C_1 (P_1, P_2, P_3, P_4, P_6)$ :

$$m_{1\text{new}} = \frac{1}{5} \sum P = \frac{[0.1 + 0.15 + 0.08 + 0.16 + 0.25 + 0.6 + 0.71 + 0.9 + 0.85 + 0.3]}{5} \\ = [0.148, 0.712]$$

- For  $(2 \ (P_5, P_7, P_8))$ :  $(0.2 + 0.24 + 0.3, 0.3 + 0.1 + 0.2) = (0.74, 0.6)$

$$m_{2\text{new}} = [0.2 + 0.24 + 0.3, 0.3 + 0.1 + 0.2]$$

$$= [0.7467, 0.6]$$

Conclusion:

The k-means clustering algorithm was successfully applied on the given 8 data points.

After the first iteration.

- $P_6$  belongs to cluster 1
- Cluster 2 has population of 3 points
- Updated centroids are:

$$\cdot m_1 = [0.148, 0.712]$$

$$\cdot m_2 = [0.7467, 0.6]$$

This demonstrates the fundamental working of k-means clustering in grouping points based on proximity to centroids.

# Assignment No. 7

Page No.	
Date	

Title: Assignment on Gradient Boost Classifier

Aim: Implement Gradient Boost classifier algorithm on given dataset.

Relevant Theory:

- Gradient Boosting is a functional gradient algorithm that repeatedly selects a function that leads in the direction of a weak hypothesis or negative gradient so that it can minimize a loss function.
- Gradient boosting classifier combines several weak learning models to produce a powerful predicting model.

Gradient Boosting consists of three essential parts:

① loss function

- loss function's purpose is to calculate how well the model predicts, given the available data.

② Weak learner

- It classifies the data, but it makes a lot of mistakes in doing so. Usually these are decision trees.

③ Additive Model

- This is how the trees are added incrementally, iteratively and sequentially.
- You should be getting closer to your final model with each iteration.

## ④ Steps in Gradient Boosting :-

- ① Fit the initial model
- ② Adjust hyperparameters
- ③ Make predictions
- ④ Analyze results.

### Working:

- Uses log (odds) to make predictions
- Converts log (odds) to probability - using logistic function

$$\text{Probability} = \frac{e^{\text{log odds}}}{1 + e^{\text{log odds}}}$$

- Updates predictions using:

$$\text{New log odds} = \text{Base log odds} + (\text{learning rate} \times \text{residual})$$

Formula for modifying leaf values :

$$\sum \text{Residual} / (\text{Previous Prob} \times (1 - \text{Previous Prob}))$$

Pros and Cons

Advantages :-

- High accuracy
- Handles both numerical & categorical data
- Robust to missing data
- Flexible with hyperparameters

## Disadvantages :

- Can overfit
- Requires many trees  $\rightarrow$  high computation cost
- Difficult to interpret without tools.

## Output :

{'Gradient-Boost Default Test Score': 0.8708736373507032  
'Gradient-Boost GridSearch Test Score': 0.9185505911254414}

## Conclusion:

Student a studied about Gradient Boost classification and implemented it on income\_evaluation.csv dataset for classification