

ML EXPERIMENT 1

⇒ Aim:-

Predict the price of the Uber ride from a given pickup point to the agreed drop-off location.

Perform following tasks:

a) Pre process the dataset

b) Identify Outliers

c) Check for correlation

d) Implement linear regression & random forest regression model.

e) Evaluate the model & compare their respective scores like R², RMSE etc.

⇒ Outcome :-

Apply preprocessing techniques on dataset.

⇒ SYSTEM REQUIREMENT :-

- 64 bit Open Source Linux or its derivatives

- Python, Jupyter Notebook

⇒ THEORY :-

- Linear Regression -

i) Linear Regression is a linear approach for modelling the relationships between a scalar response & one or more explanatory variables.

ii) It constructs a straight line in the data plane that maps independent variables to dependent variable.

iii) A general equation for linear regression is -

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon.$$

where $y \rightarrow$ dependent variable

~~$b_1 \dots b_n$~~ \rightarrow Coefficient of regression

$b_0 \rightarrow$ Intercept of regression

$\epsilon \rightarrow$ Noise

$x_1 \dots x_n \rightarrow$ Independent variable

→ One of the major drawbacks of linear regression is that it assumes a linear relationship between independent & dependent variables.

- Random Forest Regression:

→ Every decision tree has high variance, but when we combine all of them together in parallel then the resultant variance is low as each decision tree gets perfectly trained on that sample data, & hence the output ~~does~~ depends on multiple decision trees.

→ In case of regression problems, the final output is the mean of all the outputs.

- R² Score:

$$R^2 = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y}_i)^2}, \quad \bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$$

- RMSE scores:

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}$$

⇒ Conclusion :-

Hence, our regression models using Linear regression & random forests has been implemented & evaluated.

ML EXPERIMENT 2

=> Aim :-

Classify the email using binary classification method. Email spam detection has two states - as Spam & Not Spam.

Use k-nn & svm for classification. Analyze the performance.

=> Outcome :-

Apply & evaluate classification & clustering techniques.

=> System Requirements :-

- 64 bit Open Source Linux or its derivatives.
- Python, Jupyter Notebook.

=> Theory :-

- K-Nearest Neighbors :-

» K nearest neighbors is one of the most basic yet essential classification algorithm in Machine learning.

» It belongs to the supervised learning domain & finds intense application in pattern recognition, data mining & intrusion detection.

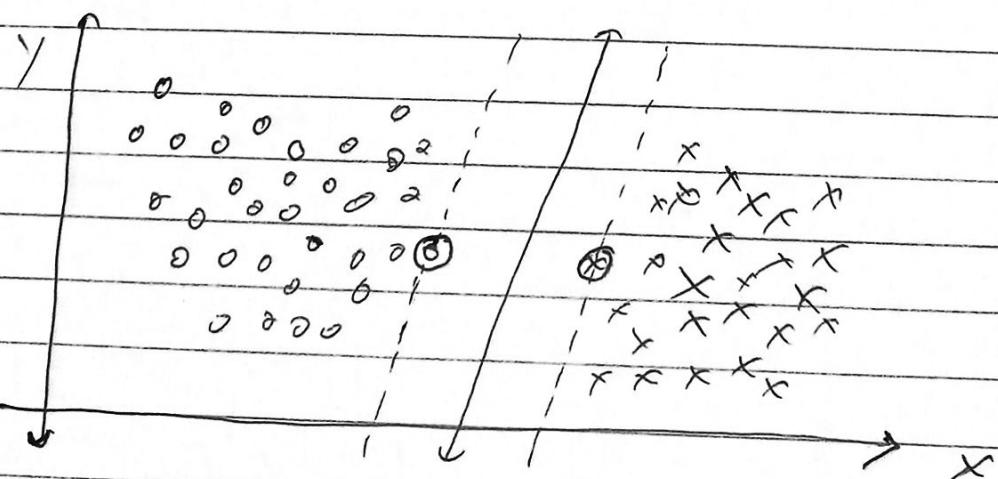
» It is widely disposable in real life scenarios since it is non-parametric meaning, it does not make any underlying assumptions about the distribution of data.

» Let m be number of training data samples. Let p be an unknown point

- a) Store the training samples in an array of data points arr .
This means each element of this array represents tuple (x_i, y_i) .
- b) Make a set S of k smallest euclidean distances.
Each of these corresponds to an already classified data points.
- c) Return the majority label among S .

~~Support Vector Machines~~

- Support Vector machines is a relatively simple Supervised Machine Learning Algorithm used for classification.
- It is more preferred for classification, it finds hyper-planes that create a boundary between the types of data.
- In 2D space, this hyper plane is a line. We plot each data item in an N dimensional plane.
- Next, we find optimal hyperplane to separate the data. So by this, you must know SVM can only perform binary classification.
- However, there are various techniques to use for multi-class problems.



(Page: 6
Date: 11/1)

⇒ Conclusion:

Hence a classification model using k-nearest neighbor & Support Vector Machine has been implemented & evaluated.

ML EXPERIMENT 3

=> Aim :-

Given a bank customer, build a neural network based classifier that can determine they will leave or not in next 6 months.

Perform the following operations -

as Read the dataset

as Distinguish the feature & target set & divide the dataset into training & testing sets.

as Normalize the train & test data

as Initialize & build the model. Identify the points of improvement & implement the same.

as Print the accuracy score & confusion matrix.

=> Outcome :-

Analyze & evaluate classification & clustering techniques.
Apply Preprocessing techniques on datasets.

=> SYSTEM REQUIREMENT :-

- 64 bit Open Source Linux or its derivatives
- Python, Jupyter Notebook.

=> THEORY :-

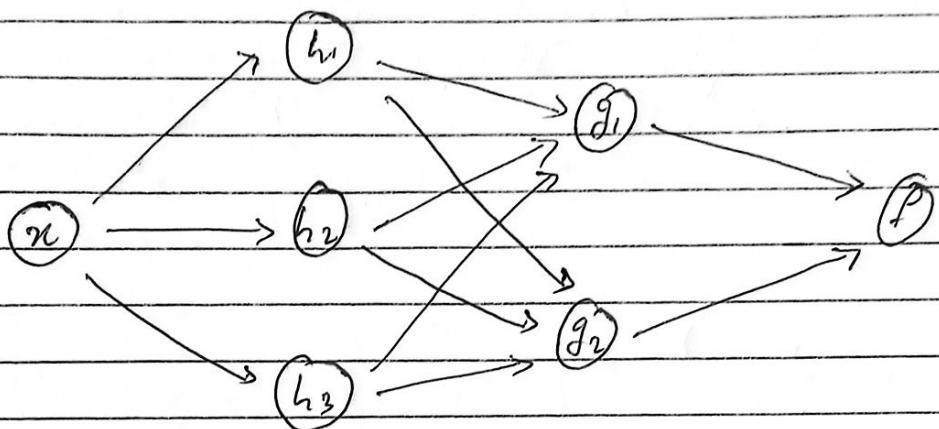
- Neural Networks :-

i) Neural networks consists of artificial neurons of functions, called parameters which allows the computer to learn, & to ~~but~~ fine tune itself, by analyzing new data.

ii) Each parameter, sometimes also referred to as neurons, is a function which produces an output after receiving one or

multiple inputs.

- iii) Those outputs are then passed on to the next layer of neurons, & so it continues until every layer of neuron have been considered, & the terminal neurons have received their input.
- iv) Those terminal neurons then output the final result for the model.



- Neural Network for Classification -

- i) A neural network can be used for classification modelling by using the correct activation functions, and correct loss function.
- ii) For a model to be used as classifier, it should use sigmoid activation function for binary classification or softmax ————— for multi-class classification.
- iii) The loss function & metrics for compiling a classification model must be "binary-crossentropy" & ["accuracy"] respectively.

Ramya DUVVU

=>

Conclusion :-

Hence a classification model using Neural Networks has been implemented & optimized.

ML EXPERIMENT 4

=> Aim :-

Implement k-nearest neighbors algorithm on Diabetes.csv dataset. Compute Confusion Matrix, accuracy, error rate, precision, recall on the given dataset.

=> Object Outcome :-

Analyze & evaluate classification & clustering techniques.

=> System Requirements :-

- 64 bit Open Source Linux or its derivatives.
- Python, Jupyter notebook.

=> Theory :-

- K Nearest Neighbors :-

i) K Nearest Neighbors is one of the most basic yet essential classification algorithm in Machine Learning.

ii) It belongs to the supervised learning domain & finds intense application in pattern ~~recognition~~ recognition, data mining & intrusion detection.

iii) It is widely disposable in real life scenarios since it is non-parametric meaning, it does not make any underlying assumptions about the distribution of data.

w) Let m be number of training data samples. Let p be an unknown point.

or store the training samples in an array of data points $a[]$. This means each element of this array represents tuple (x, y) .

b) Make a set S of k smallest euclidean distances. Each of these corresponds to an already classified data points.

c) Return the majority label among S .

- Confusion Matrix :-

		Predicted	
		Negative	Positive
Actual	Negative	True Negative (TN)	False Positive (FP)
	Positive	False Negative (FN)	True Positive (TP)

- Accuracy :-

$$\text{Accuracy} = \frac{TN + TP}{TN + FP + FN + TP}$$

- Precision :-

$$\text{Precision} = \frac{TP}{TP + FP}$$

- Recall :-

$$\text{Recall} = \frac{TP}{TP + FN}$$

- Error Rate :-

$$\text{Error Type 1} = \frac{FN}{TN + FN + FP + TP}$$

$$\text{Error Type 2} = \frac{FP}{Total}$$

DATA SET

=> Conclusion :-

Hence, a classification algorithm ~~was~~ using a K - Nearest Neighbors algorithm has been implemented and evaluated.

ML EXPERIMENT 5

19CO060

=> Aim :-

Implement k-means clustering / hierarchical clustering on sales-data-sample.csv dataset. Determine the number of clusters using the elbow method.

=> Outcome :-

Apply & evaluate classification & clustering techniques

=> System Requirement :-

- 64bit Open Source Linux or its derivatives
- Python, Jupyter Notebook.

=> Theory :-

- K - means Clustering :-

i) Clustering is one of the most common exploratory data analysis technique used to get an intuition about the structure of the data.

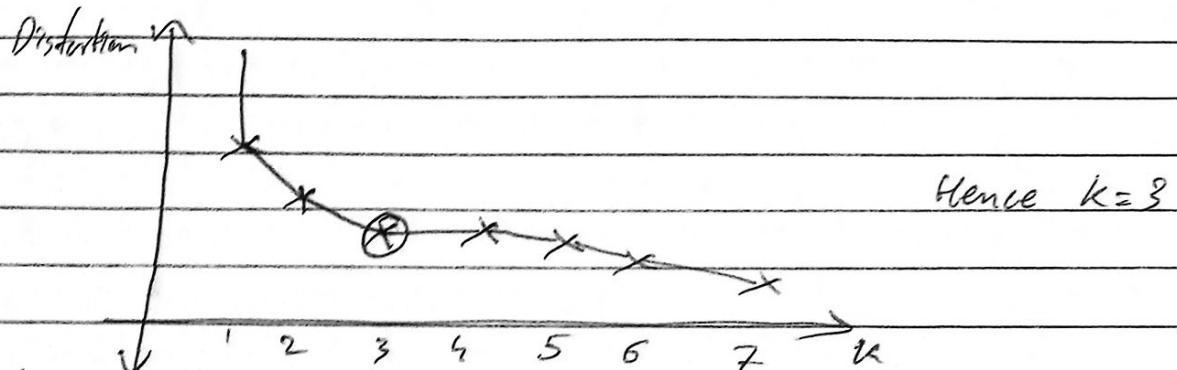
ii) K-means algorithm is an iterative algorithm that tries to partition the dataset into K pre-defined distinct non-overlapping subgroups (clusters) where each data point belongs to only one group.

iii) It tries to make the intra-cluster data points as similar as possible while also keeping the clusters as different (far) as possible.

- i) The way k-means algorithm works is as follows:
 - a) Specify number of clusters k .
 - b) Initialize centroids by first shuffling the dataset & then randomly selecting k data points for the centroids without replacement.
 - c) Keep iterating until there is no change to the centroids i.e assignment of data points to clusters isn't changing.

- Elbow Plot:

- i) A fundamental step for any unsupervised algorithm is to determine the optimal number of clusters into which the data may be clustered.
- ii) The elbow method is one of the most popular methods to determine this optimal value of k .
- iii) An elbow plot is generated by plotting ~~plotting~~ iterating the k -value from 1 to 10 and against distortion.
- iv) Distortion is calculated as the average of the squared distances from the cluster centres of the respective clusters.
- v) The k -value with the minimal distortion & an elbow like bent is chosen.



⇒

CONCLUSION :-

Hence a clustering algorithm model using a K-means clustering algorithm has been implemented.