| Answer | Coding Efficiency | Viva | Timely Completion | Total | Dated Sign of Subject Teacher |
|--------|-------------------|------|-------------------|-------|-------------------------------|
| 5 | 5 | 5 | 5 | 20 | |
| | | | | | |

Expected Date of Completion:...................... Actual Date of Completion:......................

▬-------------------------------------------------------------------------------

**Group B**

**Assignment No:2**

▬-------------------------------------------------------------------------------

**Title of the Assignment:** Classify the email using the binary classification method. Email Spam detection has two states:
a) Normal State – Not Spam,
b) Abnormal State – Spam.
Use K-Nearest Neighbors and Support Vector Machine for classification. Analyze their performance.

**Dataset Description:** The csv file contains 5172 rows, each row for each email. There are 3002 columns. The first column indicates Email name. The name has been set with numbers and not recipients' name to protect privacy. The last column has the labels for prediction : 1 for spam, 0 for not spam. The remaining 3000 columns are the 3000 most common words in all the emails, after excluding the non-alphabetical characters/words. For each row, the count of each word(column) in that email(row) is stored in the respective cells. Thus, information regarding all 5172 emails are stored in a compact dataframe rather than as separate text files.

**Link:** https://www.kaggle.com/datasets/balaka18/email-spam-classification-dataset-csv

**Objective of the Assignment:**

Students should be able to classify email using the binary Classification and implement email spam detection technique by using K-Nearest Neighbors and Support Vector Machine algorithm.

**Prerequisite:**
   1. Basic knowledge of Python
SNJB's Late Sau. K.B. Jain College of Engineering Chandwad

2. Concept of K-Nearest Neighbors and Support Vector Machine for classification.

**Contents of the Theory:**

1. Data Preprocessing
2. Binary Classification
3. K-Nearest Neighbours
4. Support Vector Machine
5. Train, Test and Split Procedure

**Data Preprocessing:**

Data preprocessing is a process of preparing the raw data and making it suitable for a machine learning model. It is the first and crucial step while creating a machine learning model.

When creating a machine learning project, it is not always a case that we come across the clean and formatted data. And while doing any operation with data, it is mandatory to clean it and put in a formatted way. So for this, we use data preprocessing task.

Why do we need Data Preprocessing?

A real-world data generally contains noises, missing values, and maybe in an unusable format which cannot be directly used for machine learning models. Data preprocessing is required tasks for cleaning the data and making it suitable for a machine learning model which also increases the accuracy and efficiency of a machine learning model.

It involves below steps:

● Getting the dataset

● Importing libraries

● Importing datasets

● Finding Missing Data

● Encoding Categorical Data

● Splitting dataset into training and test set

● Feature scaling

**Binary Classification**

Binary classification refers to those classification tasks that have two class labels.

Examples include:

- Email spam detection (spam or not).
- Churn prediction (churn or not).
- Conversion prediction (buy or not).

Typically, binary classification tasks involve one class that is the normal state and another class that is the abnormal state.

For example "*not spam*" is the normal state and "*spam*" is the abnormal state. Another example is "*cancer not detected*" is the normal state of a task that involves a medical test and "*cancer detected*" is the abnormal state.

The class for the normal state is assigned the class label 0 and the class with the abnormal state is assigned the class label 1.
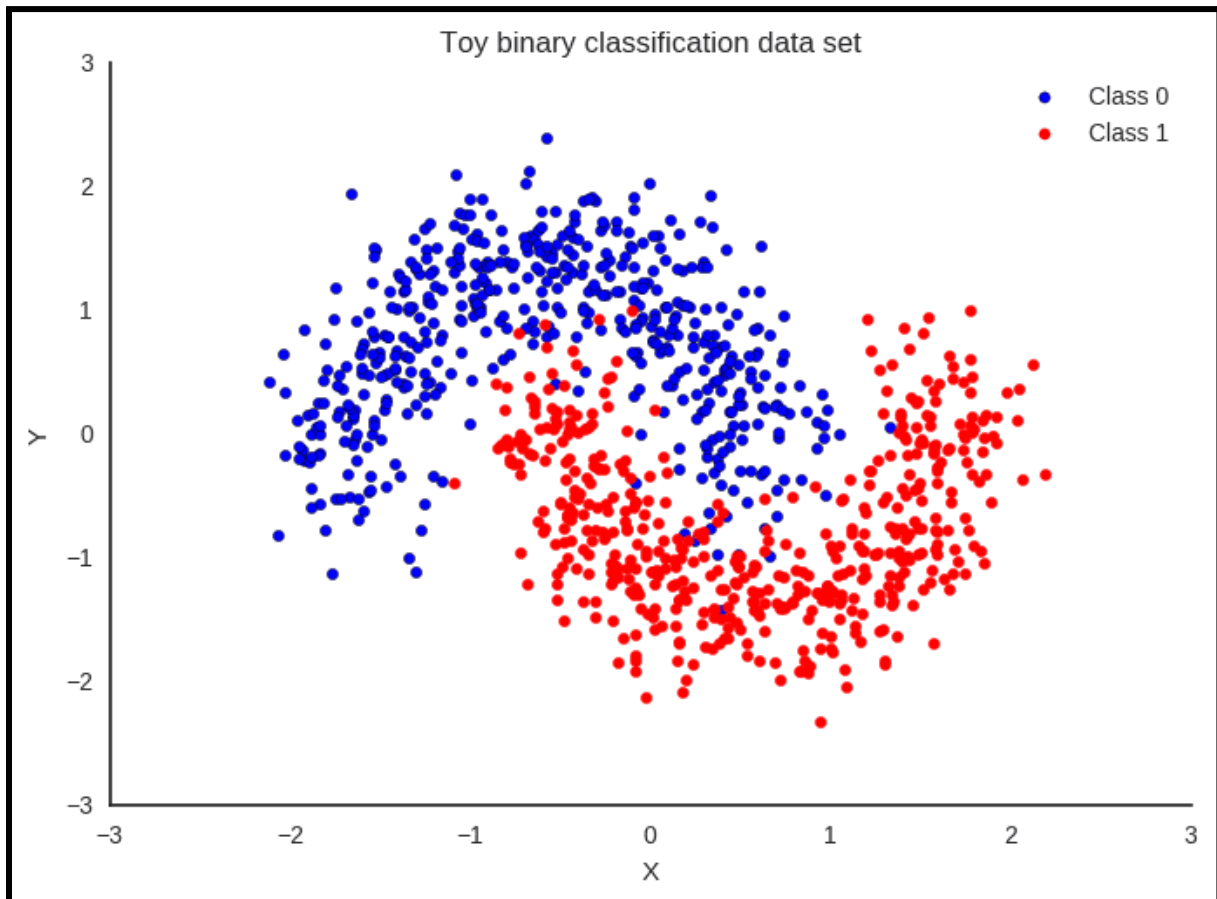
It is common to model a binary classification task with a model that predicts a Bernoulli probability distxribution for each example.

The Bernoulli distribution is a discrete probability distribution that covers a case where an event will have a binary outcome as either a 0 or 1. For classification, this means that the model predicts a probability of an example belonging to class 1, or the abnormal state.

Popular algorithms that can be used for binary classification include:

- Logistic Regression
- k-Nearest Neighbors
- Decision Trees
- Support Vector Machine
- Naive Bayes

SNJB's Late Sau. K.B. Jain College of Engineering Chandwad

Some algorithms are specifically designed for binary classification and do not natively support more than two classes; examples include Logistic Regression and Support Vector Machines.
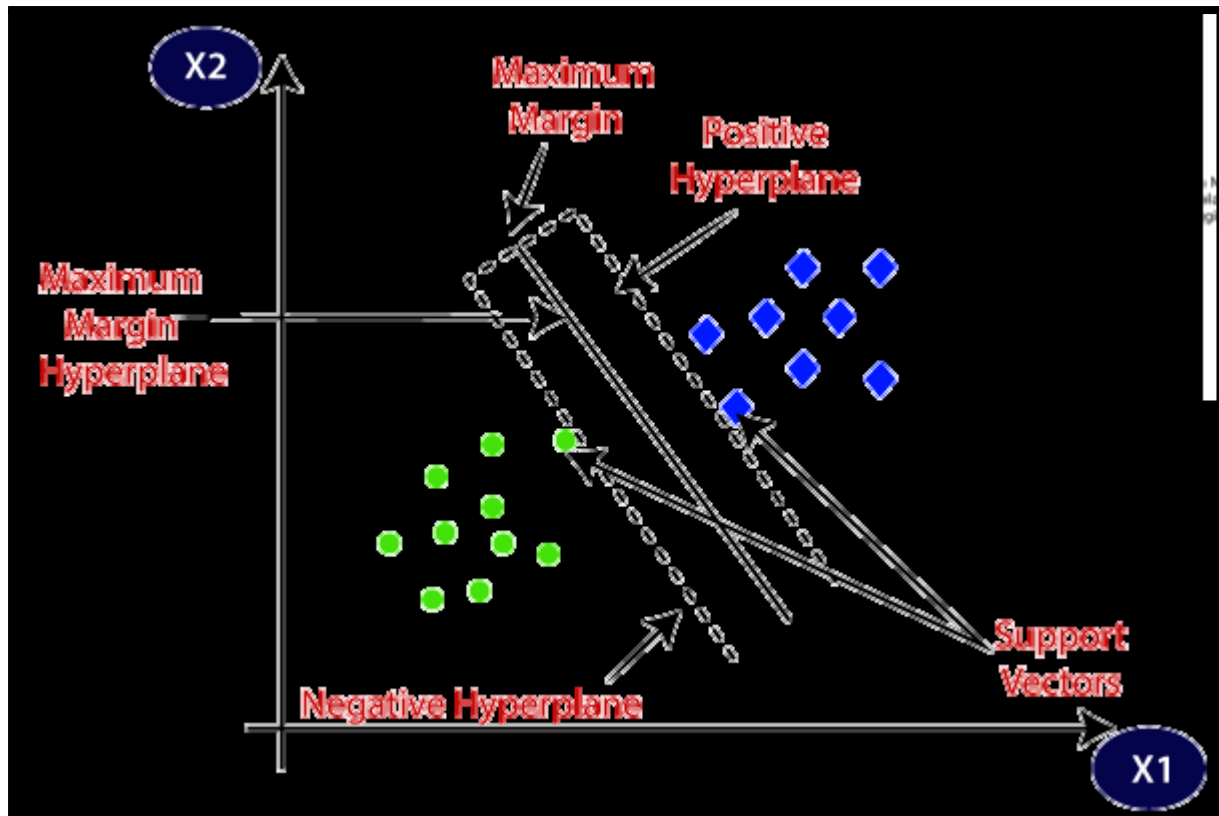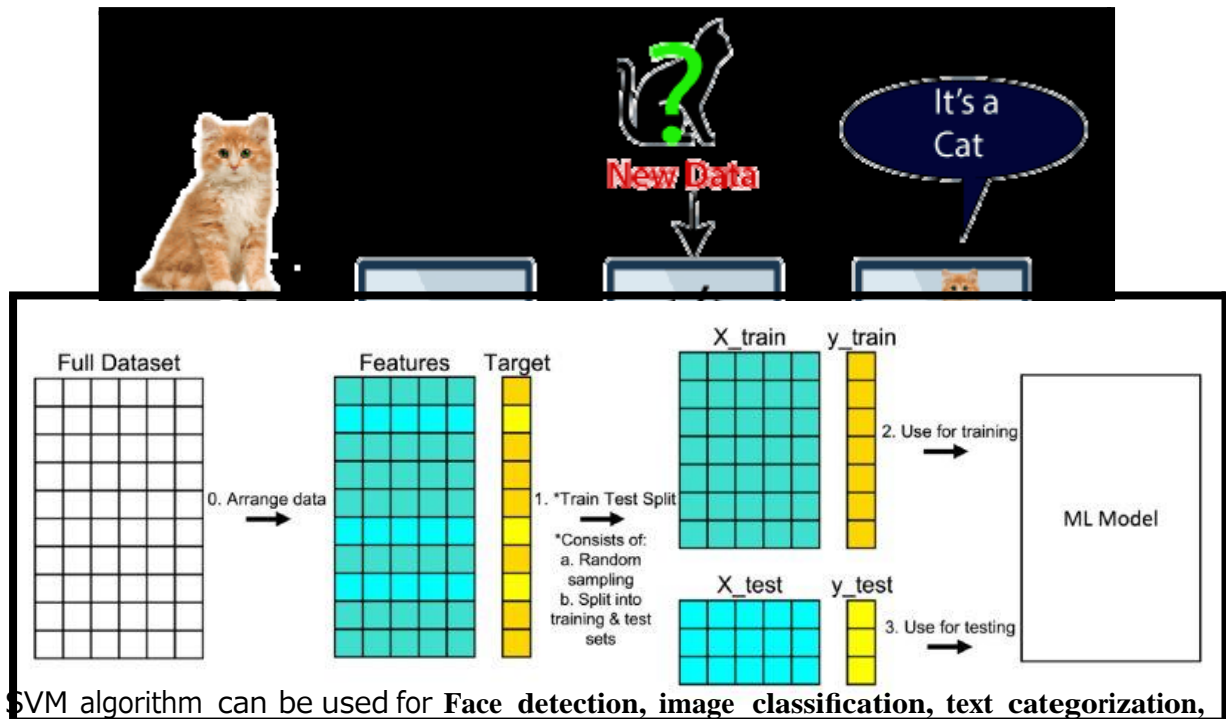


**Support Vector Machine:**

Support Vector Machine or SVM is one of the most popular Supervised Learning algorithms, which is used for Classifi cation as well as Regression problems. However, primarily, it is used for Classifi cation problems in Machine Learning.

The goal of the SVM algorithm is 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. This best decision boundary is called a hyperplane.

SVM chooses the extreme points/vectors that help in creating the hyperplane. These extreme cases are called as support vectors, and hence algorithm is termed as Support Vector Machine. Consider the below diagram in which there are two different categories that are classifi ed using a decision boundary or hyperplane:

SNJB's Late Sau. K.B. Jain College of Engineering Chandwad

**Example:** SVM can be understood with the example that we have used in the KNN classifi er. Suppose we see a strange cat that also has some features of dogs, so if we want a model that can accurately identify whether it is a cat or dog, so such a model can be created by using the SVM algorithm. We will fi rst train our model with lots of images of cats and dogs so that it can learn about different features of cats and dogs, and then we test it with this strange creature. So as support vector creates a decision boundary between these two data (cat and dog) and choose extreme cases (support vectors), it will see the extreme case of cat and dog. On the basis of the support vectors, it will classify it as a cat. Consider the below diagram:

SVM algorithm can be used for **Face detection, image classification, text categorization,** etc.

Types of SVM

**SVM can be of two types:**

- **Linear SVM:** Linear SVM is used for linearly separable data, which means if a dataset can be classifi ed into two classes by using a single straight line, then such data is termed as linearly separable data, and classifi er is used called as Linear SVM classifi er.

- **Non-linear SVM:** Non-Linear SVM is used for non-linearly separated data, which means if a dataset cannot be classifi ed by using a straight line, then such data is termed as non-linear data and classifi er used is called as Non-linear SVM classifi er.

**Train, Test, Split Procedure:**

Train test split is a model validation procedure that allows you to simulate how a model would perform on new/unseen data. Here is how the procedure works:

**1. ARRANGE THE DATA**

Make sure your data is arranged into a format acceptable for train test split. In scikit-learn, this consists of separating your full data set into "Features" and "Target."

**2. SPLIT THE DATA**

Split the data set into two pieces — a training set and a testing set. This consists of random sampling without replacement about 75 percent of the rows (you can vary this) and putting them into your training set. The remaining 25 percent is put into your test set. Note that the colors in "Features" and "Target" indicate where their data will go ("X_train," "X_test," "y_train," "y_test") for a particular train test split.

**3. TRAIN THE MODEL**

Train the model on the training set. This is "X_train" and "y_train" in the image.

**4. TEST THE MODEL**

Test the model on the testing set ("X_test" and "y_test" in the image) and evaluate the performance.

**Conclusion:**

**In this way we have explored Concept of Email Spam detection by using binary classification.**

**Assignment Questions:**

1. What is Binary Classification?
2. Explain Support Vector Machine?
3. Explain K-Nearest Neighbour algorithm for Machine Learning?