# 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 : 1for spam, 0 for not spam. The remaining 3000 columns are the 3000 most common words inall the emails, after excluding the non-alphabetical characters/words. For each row, thecount 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 asseparate 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
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 amachine learning model. It is thefirst 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 itand 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

Code :- https://www.kaggle.com/code/mfaisalqureshi/email-spam-detection-98-accuracy/notebook