

▼ Assignment 2

2. 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 link: The emails.csv dataset on the Kaggle

<https://www.kaggle.com/datasets/balaka18/email-spam-classification-dataset-csv>

```
1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 %matplotlib inline
6 import warnings
7 warnings.filterwarnings('ignore')
8 from sklearn.model_selection import train_test_split
9 from sklearn.svm import SVC
10 from sklearn import metrics
```

```
1 df=pd.read_csv('emails.csv')
```

```
1 df.head()
```

	Email No.	the	to	ect	and	for	of	a	you	hou	...	connevey	jay	valued	lay	infrastructure	military
0	Email 1	0	0	1	0	0	0	2	0	0	...	0.0	0.0	0.0	0.0	0.0	0.0
1	Email 2	8	13	24	6	6	2	102	1	27	...	0.0	0.0	0.0	0.0	0.0	0.0
2	Email 3	0	0	1	0	0	0	8	0	0	...	0.0	0.0	0.0	0.0	0.0	0.0
Email																	

```
1 df.columns
```

```
Index(['Email No.', 'the', 'to', 'ect', 'and', 'for', 'of', 'a', 'you', 'hou',
      ...,
      'connevey', 'jay', 'valued', 'lay', 'infrastructure', 'military',
      'allowing', 'ff', 'dry', 'Prediction'],
      dtype='object', length=3002)
```

```
1 df.isnull().sum()
```

```
Email No.    0
the          0
to           0
ect          0
and          0
..
military     1
allowing     1
```

```
ff          1
dry         1
Prediction  1
Length: 3002, dtype: int64
```

```
1 df.dropna(inplace = True)
```

```
1 df.isnull().sum()
```

```
Email No.    0
the          0
to           0
ect          0
and          0
..
military     0
allowing     0
ff           0
dry          0
Prediction   0
Length: 3002, dtype: int64
```

```
1 df.drop(['Email No.'],axis=1,inplace=True)
2 X = df.drop(['Prediction'],axis = 1)
3 y = df['Prediction']
```

```
1 from sklearn.preprocessing import scale
2 X = scale(X)
3 # split into train and test
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 42)
```

▼ KNN classifier

```
1 from sklearn.neighbors import KNeighborsClassifier
2 knn = KNeighborsClassifier(n_neighbors=7)
3
4 knn.fit(X_train, y_train)
5 y_pred = knn.predict(X_test)
```

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

```
Prediction [1. 1. 1. 1. 1. 0. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1.
 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 0. 1. 0. 1. 1. 1. 1. 1. 1. 0. 1. 0. 1.
 1. 0. 0. 1. 1. 0. 1. 0. 0. 1. 0. 1. 1. 0. 0. 1. 1. 1. 1. 0. 1. 1. 1.
 0. 1. 0. 0. 1. 0. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 0. 1. 0. 1. 1. 1. 0.
 1. 1. 1. 1. 1. 1. 0. 1. 1. 0. 0. 0. 0. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1.
 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 0. 1. 1. 1. 1. 0. 0. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 0. 0. 1. 1. 0. 1. 1. 1. 1. 0. 1. 1. 1. 1. 0. 1. 1. 0. 1. 1.
 0. 1. 1. 0. 0. 1. 1. 1. 1. 0. 1. 1. 0. 1. 1. 1.]
```

```
1 print("KNN accuracy = ",metrics.accuracy_score(y_test,y_pred))
```

```
KNN accuracy = 0.6038461538461538
```

```
1 print("Confusion matrix",metrics.confusion_matrix(y_test,y_pred))
```

```
Confusion matrix [[182 200]
 [ 6 132]]
```

▼ SVM classifier

```
1 # cost C = 1
2 model = SVC(C = 1)
3
4 # fit
5 model.fit(X_train, y_train)
6
7 # predict
8 y_pred = model.predict(X_test)
```

```
1 metrics.confusion_matrix(y_true=y_test, y_pred=y_pred)
```

```
array([[382,  0],
 [ 46,  92]])
```

```
1 print("SVM accuracy = ",metrics.accuracy_score(y_test,y_pred))
```

```
SVM accuracy =  0.9115384615384615
```