## **Spam Detection Classifier** Importing all necessary packages In [2]: import numpy as np import pandas as pd from sklearn.feature\_extraction.text import CountVectorizer from sklearn.feature\_extraction.text import TfidfVectorizer from sklearn.model\_selection import train test split from sklearn.naive\_bayes import MultinomialNB from sklearn.metrics import accuracy score from sklearn.metrics import classification\_report **Reading Dataset** In [3]: spam = pd.read\_csv('Spam.csv',encoding = 'latin-1') spam.head() Out[3]: v1 v2 Unnamed: 2 Unnamed: 3 Unnamed: 4 0 ham Go until jurong point, crazy.. Available only ... NaN NaN NaN ham Ok lar... Joking wif u oni... NaN NaN NaN 2 spam Free entry in 2 a wkly comp to win FA Cup fina... NaN NaN NaN U dun say so early hor... U c already then say... NaN NaN NaN Nah I don't think he goes to usf, he lives aro... NaN ham In [4]: list (spam) Out[4]: ['v1', 'v2', 'Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4'] In [6]: spam.head() Out[6]: v1 v2 Unnamed: 2 Unnamed: 3 Unnamed: 4 0 ham Go until jurong point, crazy.. Available only ... NaN NaN NaN ham Ok lar... Joking wif u oni... NaN NaN NaN 2 spam Free entry in 2 a wkly comp to win FA Cup fina... NaN NaN NaN U dun say so early hor... U c already then say... NaN NaN NaN ham Nah I don't think he goes to usf, he lives aro... NaN NaN NaN spam.rename(columns = {'v1':'Status', 'v2':'Email'}, inplace = True) In [8]: spam.columns Out[8]: Index(['Status', 'Email', 'Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4'], dtype='object') **Deleting Unnecessary Columns** In [9]: spam.drop (['Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4'], axis = 1 , inplace = True) In [10]: spam.columns Out[10]: Index(['Status', 'Email'], dtype='object') In [11]: spam.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 5572 entries, 0 to 5571 Data columns (total 2 columns): Status 5572 non-null object Email 5572 non-null object dtypes: object(2) memory usage: 87.1+ KB In [39]: spam.head() Out[39]: **Status Email** 0 ham Go until jurong point, crazy.. Available only ... ham Ok lar... Joking wif u oni... spam Free entry in 2 a wkly comp to win FA Cup fina... U dun say so early hor... U c already then say... ham ham Nah I don't think he goes to usf, he lives aro... Converting Ham = 1 and Spam = 0 In [12]: spam.loc [spam['Status'] == 'spam' , 'Status'] = 0 spam.loc [spam['Status'] == 'ham' , 'Status'] = 1 In [13]: spam.head() Out[13]: Status **Email** Go until jurong point, crazy.. Available only ... Ok lar... Joking wif u oni... 0 Free entry in 2 a wkly comp to win FA Cup fina... 1 U dun say so early hor... U c already then say... Nah I don't think he goes to usf, he lives aro... In [15]: spam\_x = spam['Email'] spam\_y = spam['Status'] CountVectorizer() Method In [17]: cv = CountVectorizer() x\_train , x\_test , y\_train , y\_test = train\_test\_split( spam\_x , spam\_y , test\_size = 0.2) In [49]: len(x\_train) Out[49]: 4457 In [50]: len(x\_test) Out[50]: 1115 In [26]: x\_traincv = cv.fit\_transform (x\_train) In [27]: x\_traincv.toarray() Out[27]: array([[0, 0, 0, ..., 0, 0, 0], $[0, 0, 0, \ldots, 0, 0, 0],$ $[0, 0, 0, \ldots, 0, 0, 0],$ $[0, 0, 0, \ldots, 0, 0, 0],$ [0, 0, 0, ..., 0, 0, 0], [0, 0, 0, ..., 0, 0, 0]], dtype=int64) In [32]: y\_train = y\_train.astype('str') In [33]: x\_testcv = cv.transform(x\_test) In [34]: model = MultinomialNB() model.fit(x traincv, y train) Out[34]: MultinomialNB(alpha=1.0, class\_prior=None, fit\_prior=True) In [35]: pred = model.predict(x\_testcv) In [36]: pred Out[36]: array(['1', '0', '1', ..., '1', '1'], dtype='<U1') Accuracy of Model build through CountVectorizer() In [45]: y test = y test.astype('str') print('Accuracy Score :', accuracy\_score(pred, y\_test)) Accuracy Score : 0.9820627802690582 TfidfVectorizer() Method In [46]: | td = TfidfVectorizer(min\_df = 1) In [47]: x\_traintd = td.fit\_transform (x\_train) x\_traintd.toarray() Out[47]: array([[0., 0., 0., ..., 0., 0., 0.], [0., 0., 0., ..., 0., 0., 0.],[0., 0., 0., ..., 0., 0., 0.],[0., 0., 0., ..., 0., 0., 0.],[0., 0., 0., ..., 0., 0., 0.],[0., 0., 0., ..., 0., 0., 0.]]) In [48]: x\_testtd = td.transform(x\_test) In [49]: perdtd = model.predict(x testtd) Accuracy of Model build through TfidfVectorizer() In [50]: print('Accuracy Score : ', accuracy\_score(perdtd, y\_test)) Accuracy Score: 0.9829596412556054 In [51]: classification\_report(y\_test, perdtd) Out[51]: ' precision recall f1-score support\n\n 1.00 0.89 0.94 168\n 1 0.98 1.00 0.99 947\n\n accuracy 0.98 1115\n macro avg 0.99 0.94 0.97 0.98 1115\nweighted avg 0.98 0.98 1115\n'

In [ ]: