



BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE, PILANI

BITS F464 Machine Learning

Name	ID No.
Donkada Vishal Dheeraj	2018A7PS0239H
Thakur Shivank Singh	2018A7PS0439H
Pranav Reddy Pesaladinne	2018A7PS0238H

Assignment 1B

Naive Bayes Classifier

Design Procedure and Implementation

The Naive Bayes classifier is a simple classifier that classifies based on probabilities of events. It is applied commonly to text classification. Though it is a simple algorithm, it performs well in many text classification problems.

We have implemented a Naïve Bayes Classifier model which classifies a given mail as spam or not using pandas and NumPy libraries. The code mainly consists of the following parts –

- Loading the Data (.txt to .csv)

Converting the given .txt file into a .csv file with mail text as one column and its label (spam or not) in next column. We read the text file and based on the format and whitespaces (label at the end of each line), make a dictionary with the column names and convert it to a .csv file using pandas library.

```
1 file = open(pathdir + '/dataset_NB.txt','r')    # opening .txt file
2 dataset = file.readlines()    # reads text file and returns a list containing each line in the file as a list item
3 dataset[-1] = dataset[-1]+'\\n'
4 file.close()
```

```
1 text = []
2 spam = []
3 for data in dataset:
4     text.append(data[:-2].rstrip())    # appending mail text to list
5     spam.append(data[-2])    # appending spam label to list
```

```
1 dictionary = {
2     'mail':text,
3     'spam':spam
4 }
5 df = pd.DataFrame(dictionary)    # convert dict to pandas dataframe
6 df.to_csv(pathdir + '/dataset_NB.csv',index=False)    # saving dataframe as .csv file
```

- Data Analysis and Pre-processing

First, we check the different statistics of the dataset like word count, number of unique/distinct words and topmost frequently occurring words. This gives us an overall idea of the text data we are dealing with and gives us an idea of the data pre-processing steps which may help the model.

```

1 print("Distribution of labels in given data")
2 print(df['spam'].value_counts()) # count of unique labels in the 'spam' column
3 # df.groupby('spam').spam.count()
4
5 count = 0
6 for mail in df['mail']:
7     count += len(mail.split())
8 print("Total number of words:", count)
9
10 word_dict = {}
11 for mail in df['mail']:
12     for word in mail.split():
13         word = word.lower()
14         if word not in word_dict:
15             word_dict[word] = 1
16         else:
17             word_dict[word] += 1
18 print("Total distinct words:", len(word_dict))
19 print("Top frequent words:-")
20 dict(sorted(word_dict.items(), key=lambda item: item[1], reverse=True)) # prints a sorted dict with topmost freq words

```

The following steps help in cleaning and filtering the text data to extract only important features (words) which may help the model in classifying the mails better.

- lower_case() : to lowercase text
- remove_punct(): to remove punctuations and other symbols
- remove_stopwords(): to remove a defined list of stopwords (and other frequently occurring words)

```

1 # function to lowercase text
2 def lower_case(text):
3     clean_text = str(text).lower()
4     return clean_text
5
6 # function to remove punctuations and other symbols
7 def remove_punct(text):
8     punctuations = '''!#$%&'()*+,-./:;<>[]^_`{|}~0123456789'''
9     clean_text = ""
10    for char in text:
11        if char in punctuations:
12            clean_text += " "
13        else:
14            clean_text += char
15    return " ".join(clean_text.strip().split())
16
17 # function to remove stopwords
18 def remove_stopwords(text):
19     tmp = []
20    for word in text.split():
21        if word not in stopwords:
22            tmp.append(word)
23    clean_text = " ".join(tmp)
24    return clean_text

```

```

1 def preprocessing(df):
2     for i,data in enumerate(df.iloc[:,0]):
3         data = lower_case(data)
4         data = remove_punct(data)
5         data = remove_stopwords(data)
6     df.iloc[i] = data
7     return df

```

- Cross Validation and Splitting Data

We use 7-fold Cross Validation and split data into training and test data based on the iterations.

```

8 k = 7 # for 7-Fold Cross Validation
9 for _ in range(0,k):
10     # K-FOLD SPLIT
11     start = int((_*len(df))/k)
12     end = int(((+1)*len(df))/k)
13     # print(start, end)
14     ### TEST data
15     X_test = df.iloc[start:end+1, :-1] # all cols except last col as features (X)
16     Y_test = df.iloc[start:end+1, -1] # last col as labels/classes/categories (Y)
17     ### TRAINING data
18     X_train = df.iloc[:start, :-1]
19     X_train = X_train.append(df.iloc[end+1:, :-1])
20     Y_train = df.iloc[:start, -1]
21     Y_train = Y_train.append(df.iloc[end+1:, -1])
22
23     X_train = X_train.reset_index(drop=True)
24     X_test = X_test.reset_index(drop=True)
25     Y_train = Y_train.reset_index(drop=True)
26     Y_test = Y_test.reset_index(drop=True)

```

- Finding unique words and Bag of words

We find unique words and store it in a dictionary, with key as the word and value as its count. Then we use a list to represent bag of words. We make a list for every mail and fill the index where the unique word is present in that mail. This list is then appended to the overall “bow” list and then converted to a NumPy array.

```

45 # UNIQUE WORDS DICTIONARY
46 uniq_dict = {}
47 for mail in X_train['mail']:
48     # print(mail)
49     for word in mail.split():
50         # print(word)
51         if word not in uniq_dict:
52             uniq_dict[word] = 1
53         else:
54             uniq_dict[word] += 1

```

```

60 # BAG OF WORDS
61 bow = [] # bag of words
62 for mail in X_train['mail']:
63     ohe_sent = [0] * len(uniq_dict) # one hot encoding of words in a sentence
64     for word in mail.split():
65         ohe_sent[list(uniq_dict.keys()).index(word)] = 1
66     bow.append(ohe_sent)
67
68 bow = np.asarray(bow)
69 print("Bag of Words shape:", bow.shape)

```

- Calculating probabilities and testing

We now implement the Naïve Bayes model and calculate probabilities based on the standard Bayes formula.

$$p(C_k|x) = \frac{p(C_k)p(x|C_k)}{p(x)} \propto p(C_k)p(x|C_k) = p(C_k, x_1, x_2, \dots, x_n)$$

We also apply Laplace Smoothing to solve the problem of zero probability for new unique words in the test data.

$$p_{\lambda}(C_k) = p_{\lambda}(Y = C_k) = \frac{\sum_{i=1}^N I(y_i = C_k) + \lambda}{N + K\lambda}$$

$$p(x_1 = a_j | y = C_k) = \frac{\sum_{i=1}^N I(x_{1i} = a_j, y_i = C_k) + \lambda}{\sum_{i=1}^N I(y_i = C_k) + A\lambda}$$

Where,

K is the number of classes/labels

A is the number of features (total number of unique words).

We have used $\lambda=1$ (ALPHA=1)

Firstly, we calculate prior probabilities and store in prior_y1 and prior_y0.

Then, we calculate the conditional probabilities for every word based on Two cases – y=0 and y=1.

- conditional_probs: A numpy array which stores the conditional probabilities of all unique words. The first row is for the case when y=0 and the second row when y=1

```

74 # CALCULATING PROBABILITIES
75 ALPHA = 1 # Laplace Smoothing (to handle missing words in dictionary)
76
77 #calculate prior probabilities => P(y=1) , P(y=0)
78 num_y1 = sum(Y_train)
79 num_y0 = len(Y_train) - num_y1
80 prior_y1 = (num_y1+ALPHA)/(len(Y_train) + (len(Y_train.unique())*ALPHA))
81 prior_y0 = 1-prior_y1
82
83 #calculate likelihoods/conditional prob => P(x/y=0) , P(x/y=1)
84 # Case Y=0:
85 pw_y0 = [0] * len(uniq_dict)
86 for j in range(len(uniq_dict)):
87     sumw = 0
88     for i in range(len(bow)):
89         if bow[i][j] == 1 and Y_train[i] == 0:
90             sumw += 1
91     pw_y0[j] = sumw
92
93 pw_y0 = np.asarray(pw_y0)
94
95 conditional_probs = np.zeros((len(uniq_dict))) # array of conditional probabilities calculated for each word in dict
96 conditional_probs = (pw_y0 + ALPHA)/(num_y0 + len(uniq_dict)*ALPHA)
97 conditional_probs = np.reshape(conditional_probs, (1,len(uniq_dict)))

```

```

99     # Case Y=1:
100     pw_y1 = [0] * len(uniq_dict)
101     for j in range(len(uniq_dict)):
102         sumw = 0
103         for i in range(len(bow)):
104             if bow[i][j] == 1 and Y_train[i] == 1:
105                 sumw += 1
106         pw_y1[j] = sumw
107
108     pw_y1 = np.asarray(pw_y1)
109     # print(pw_y1)
110
111     arr = (pw_y1 + ALPHA)/(num_y1 + len(uniq_dict)*ALPHA)
112     arr = np.reshape(arr,(1, len(uniq_dict)))
113     conditional_probs = np.append(conditional_probs, arr, axis=0)

```

Now, we begin classifying the test data. Firstly, we preprocess the data and then based on the text data, we substitute already calculated probabilities (during training) in the formula to obtain $py0_w$ and $py1_w$. Based on which probability is higher, we classify the mail as spam(1) or not spam(0).

➤ Y_pred : list which stores the predictions of the model

```

143     # TEST ACCURACY
144     X_test_pp = preprocessing(X_test)
145
146     Y_pred = [0] * len(X_test_pp)
147     for i,mail in enumerate(X_test_pp['mail']):
148         py0_w = prior_y0
149         py1_w = prior_y1
150         for w in mail.split():
151             if w in uniq_dict:
152                 py0_w *= conditional_probs[0][list(uniq_dict.keys()).index(w)]
153                 py1_w *= conditional_probs[1][list(uniq_dict.keys()).index(w)]
154             else:
155                 py0_w *= ALPHA/(num_y0 + len(uniq_dict)*ALPHA)
156                 py1_w *= ALPHA/(num_y1 + len(uniq_dict)*ALPHA)
157
158         # print(py1_w, py0_w)
159         if py1_w >= py0_w:
160             #y_pred.append(1)
161             Y_pred[i] = 1
162         else:
163             #y_pred.append(0)
164             Y_pred[i] = 0

```

- Measuring performance

Finally, we calculate the Accuracy and F1 Score comparing predictions (Y_pred) and test data actual labels (Y_test)

```

1  # calculating accuracy and f1 score
2  def score(Y_test, Y_pred):
3      # f1_score = TP/(TP+0.5*(FP+FN))
4      TP = 0
5      FP = 0
6      FN = 0
7      TN = 0
8      for i in range(len(Y_test)):
9          if Y_pred[i] == 1 and Y_test[i] == 1:
10             TP += 1
11          elif Y_pred[i] == 1 and Y_test[i] == 0:
12             FP += 1
13          elif Y_pred[i] == 0 and Y_test[i] == 1:
14             FN += 1
15          else:
16             TN += 1
17     acc = (TP+TN)/(len(Y_test))
18     f1_scr = TP/(TP+0.5*(FP+FN))
19     return acc, f1_scr

```

Results

Accuracy and F1 Score of the model over each fold and the overall average accuracy scores -

```
1. Test Accuracy: 83.22 %  
   F1 Score:      0.8537  
2. Test Accuracy: 86.81 %  
   F1 Score:      0.8742  
3. Test Accuracy: 79.17 %  
   F1 Score:      0.8171  
4. Test Accuracy: 77.78 %  
   F1 Score:      0.8  
5. Test Accuracy: 78.47 %  
   F1 Score:      0.777  
6. Test Accuracy: 77.78 %  
   F1 Score:      0.7895  
7. Test Accuracy: 76.92 %  
   F1 Score:      0.7626  
Average Test Accuracy: 80.01998001998001  
Average F1 Score:      0.8105637031831149
```

We were able to achieve 80% overall accuracy using the Naïve Bayes Classifier after Laplace Smoothing and Data Pre-processing to classify mails as spam or not.

Conclusion and Limitations

- The 'naive' nature of the Naïve Bayes Classifier comes from the strong assumption that all features are independent of each other. Although this belief, makes the model faster than most other algorithms, this case rarely happens in real life as there are a lot of interactions and relations between every feature which affect the overall outcome. The assumption that all features are independent makes Naive Bayes algorithm less accurate than other complicated algorithms. Hence, we see there is a trade off between speed and accuracy.
- The Naïve Bayes algorithm faces the 'zero-frequency problem' where it assigns zero probability to a variable in the test data set which wasn't available in the training dataset. This can be overcome using a technique called Laplace Smoothing which we have implemented in the model.
- The probability estimations may be incorrect based on the skewness of the data. Hence performance is sensitive to skewed data based on its distribution.