

# IR Assignment-4

- SUBHANI SHAIK (MT18117)

**Dataset:** [20newsgroups dataset](#)

In this dataset, used **comp.graphics, sci.med, talk.politics.misc, rec.sport.hockey, sci.space** [5 classes] which have 1000 documents each i.e., total 5000 documents.

## 1. Naive Bayes Classification for different splits:

### Pre-Processing the dataset:

- Initially read all the class document names from each folder and stored to a list.
- Read those files and done the pre-processing for both train and test data.
- For pre-processing the text, used NLTK library and done following steps:
  - Conversion to lowercase
  - Contractions
  - Removed unnecessary characters and punctuations and tokenized using RegexpTokenizer.
  - Lemmatization
  - Stop words removal.

### Procedure:

- Initially split the dataset according to our need as training and test data.
- Prepare the 'Dictionary' for this training data as dictionary structure as-  
{class\_num:[tokenized data of all docs in that class]}
- Prepare 'Vocabulary' for that split by merging all classes tokenized data and take set of them.
- Next, train the data by calculating 'class prior'(class probability) and 'word conditional probability'(these words are from vocabulary which is prepared from above) and stored them into a dictionary.

$$p(C_k | \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} | C_k)}{p(\mathbf{x})}$$

$$\begin{aligned} p(C_k | x_1, \dots, x_n) &\propto p(C_k, x_1, \dots, x_n) \\ &= p(C_k) p(x_1 | C_k) p(x_2 | C_k) p(x_3 | C_k) \dots \\ &= p(C_k) \prod_{i=1}^n p(x_i | C_k), \end{aligned}$$

Where,

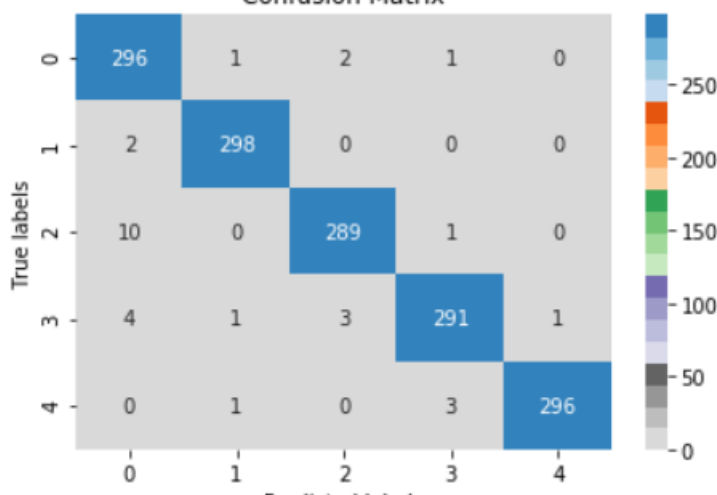
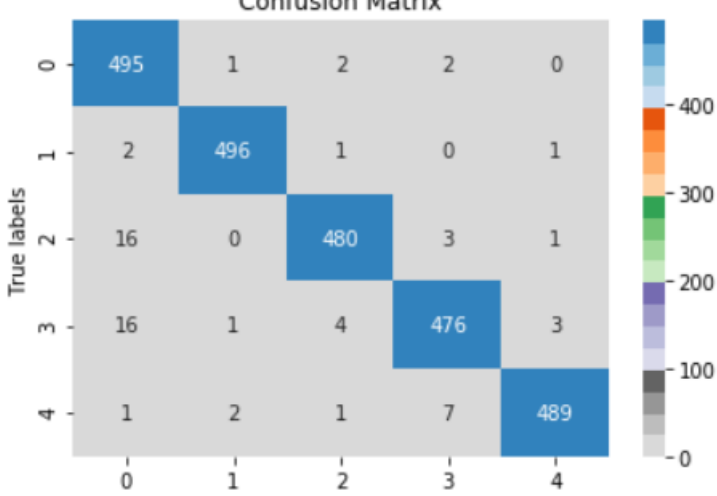
$p(C_k)$  = class prior probability,

$p(x_i | C_k)$  = conditional probability for each word.

- These probability values will be used at testing time.
- Save those variables into files by using 'Joblib' and load them for testing.

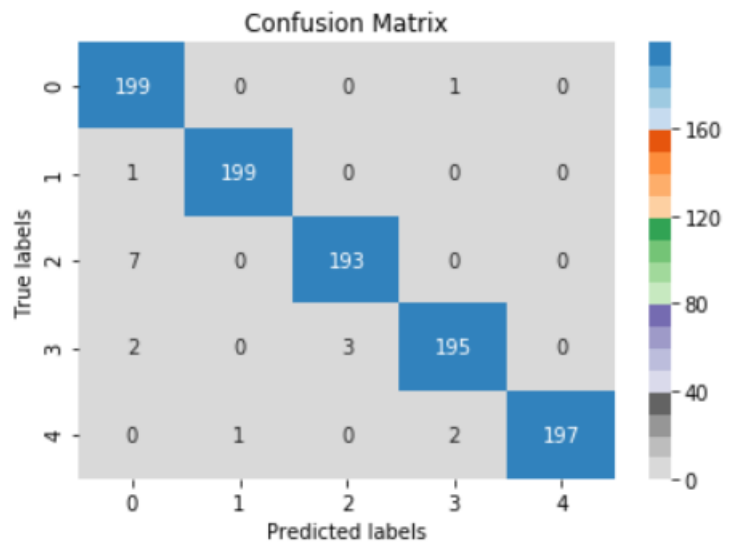
- For each class and for each word from vocabulary, take the conditional probability and class prior and calculate the logarithm of them and add them.
- Take max of all class probability for that word and predict class for it.
- Later after predicting for all the docs, calculate the 'Confusion matrix' and 'Accuracy' for that split.
- Repeated the same process for other splits.

## Results:

Split ratio	Accuracy and Confusion matrix																																				
70:30	<div><div>0.98</div><div><p>Confusion Matrix</p><table><tr><th>True \ Pred</th><th>0</th><th>1</th><th>2</th><th>3</th><th>4</th></tr><tr><th>0</th><td>296</td><td>1</td><td>2</td><td>1</td><td>0</td></tr><tr><th>1</th><td>2</td><td>298</td><td>0</td><td>0</td><td>0</td></tr><tr><th>2</th><td>10</td><td>0</td><td>289</td><td>1</td><td>0</td></tr><tr><th>3</th><td>4</td><td>1</td><td>3</td><td>291</td><td>1</td></tr><tr><th>4</th><td>0</td><td>1</td><td>0</td><td>3</td><td>296</td></tr></table></div></div>	True \ Pred	0	1	2	3	4	0	296	1	2	1	0	1	2	298	0	0	0	2	10	0	289	1	0	3	4	1	3	291	1	4	0	1	0	3	296
True \ Pred	0	1	2	3	4																																
0	296	1	2	1	0																																
1	2	298	0	0	0																																
2	10	0	289	1	0																																
3	4	1	3	291	1																																
4	0	1	0	3	296																																
50:50	<div><div>0.9744</div><div><p>Confusion Matrix</p><table><tr><th>True \ Pred</th><th>0</th><th>1</th><th>2</th><th>3</th><th>4</th></tr><tr><th>0</th><td>495</td><td>1</td><td>2</td><td>2</td><td>0</td></tr><tr><th>1</th><td>2</td><td>496</td><td>1</td><td>0</td><td>1</td></tr><tr><th>2</th><td>16</td><td>0</td><td>480</td><td>3</td><td>1</td></tr><tr><th>3</th><td>16</td><td>1</td><td>4</td><td>476</td><td>3</td></tr><tr><th>4</th><td>1</td><td>2</td><td>1</td><td>7</td><td>489</td></tr></table></div></div>	True \ Pred	0	1	2	3	4	0	495	1	2	2	0	1	2	496	1	0	1	2	16	0	480	3	1	3	16	1	4	476	3	4	1	2	1	7	489
True \ Pred	0	1	2	3	4																																
0	495	1	2	2	0																																
1	2	496	1	0	1																																
2	16	0	480	3	1																																
3	16	1	4	476	3																																
4	1	2	1	7	489																																

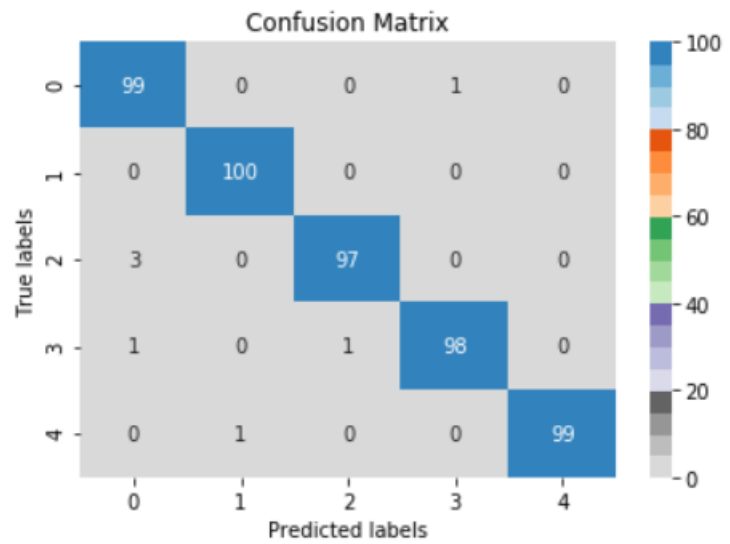
80:20

0.983



90:10

0.986



## 2. Naive Bayes Classification for different splits:

### Pre-Processing the dataset:

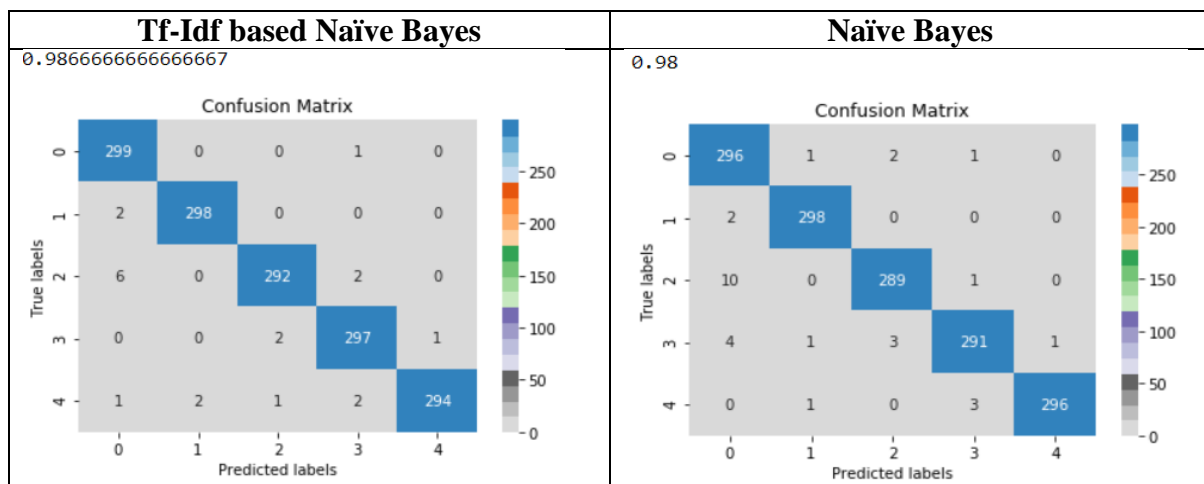
- Initially read all the class document names from each folder and stored to a list.
- Read those files and done the pre-processing for both train and test data.
- For pre-processing the text, used NLTK library and done following steps:
  - Conversion to lowercase
  - Contractions
  - Removed unnecessary characters and punctuations and tokenized using RegexpTokenizer.
  - Lemmatization
  - Stop words removal.

### Procedure:

- Same as Q1 procedure but here taken features by calculating Tf-Idf values for each word in each document of class.
- After calculating Tf-Idf values and take top 'k' features from each class and merge them all which we consider as vocabulary for this question and do same as Q1.
- Construct the confusion matrix and accuracy for this.
- Compare with Q1 70:30 split.

### Result and Analysis:

#### Accuracy and Confusion matrix for this question and Q1 70:30 split



Here we have taken features based on Tf-Idf values and done the Naïve Bayes classification. It means we are selecting best features from each class. So that our accuracy is better than without feature selection model (Q1).