

## Text preparation, classification, and model assessment

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 #import text data from csv file
6 df=pd.read_csv('/content/Text_data.csv')
7 #df.head()
8 df.tail()
9 #pd.set_option('max_colwidth', 100)
```

```
1 sns.countplot(df.ADR_label) #Imbalanced Target
```

```
1 import nltk
2 nltk.download('punkt')
3 nltk.download('stopwords')
```

```
1 import nltk
2 import string
3 import re
4 from nltk.stem.snowball import SnowballStemmer
5 stopwords=nltk.corpus.stopwords.words('english')
6 snowball_stemmer=SnowballStemmer(language='english')
7 def treat_text(text):
8     edited_text=re.sub('\W'," ",text) #replace any sumbol with whitespace
9     edited_text=re.sub("  "," ",edited_text) #replace double whitespace with single whitespace
10    edited_text=edited_text.split(" ") #split the sentence into array of strings
11    edited_text=" ".join([char for char in edited_text if char!= ""]) #remove any empty string from text
12    edited_text=edited_text.lower() #lowercase
13    edited_text=re.sub('\d+', "",edited_text) #Removing numerics
14    edited_text=re.split('\W+',edited_text) #splitting based on whitespace or whitespaces
15    edited_text=" ".join([snowball_stemmer.stem(word) for word in edited_text if word not in stopwords]) #Snowball Stemmer
16    return edited_text
17 df['Treated_Tweet']=df.Tweet.apply(lambda x: treat_text(x))
18 df.head()
```

```
1 # Most Frequent words in the dataset
2 df.Treated_Tweet.str.split(expand=True).stack().value_counts()[:10]
```

```
1 freq_words=df.Treated_Tweet.str.split(expand=True).stack().value_counts()[:10]
2 freq_words=list(freq_words.index)
```

```
1 # Rare words in the dataset appearing only once in the whole dataset
2 df.Treated_Tweet.str.split(expand=True).stack().value_counts()[-20:]
```

```
1 rare_words=df.Treated_Tweet.str.split(expand=True).stack().value_counts()
2 rare_words=list(rare_words.loc[lambda x: x==1].index)
```

```
1 #Remove Frequent and Rare words
2 def remove_noise_words(text):
3     edited_text=text.split()
4     edited_text=[word for word in edited_text if word not in freq_words]
5     edited_text=[word for word in edited_text if word not in rare_words]
6     edited_text=" ".join(edited_text)
7     return edited_text
```

```
1 df['Final_Treated_words']=df.Treated_Tweet.apply(lambda x: remove_noise_words(x))
```

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

```
1 X=df.Final_Treated_words
2 y=df.ADR_label
```

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test=train_test_split(X, y, test_size=0.20, random_state=10)
```

```
1 print(X_train.shape)
2 print(X_test.shape)
3 print(y_train.shape)
4 print(y_test.shape)
```

```
1 #Creating a list of Pipeline with well-known ML models
2 from sklearn.pipeline import make_pipeline
3 from sklearn.naive_bayes import MultinomialNB,ComplementNB
4 from sklearn.linear_model import LogisticRegression, RidgeClassifier
5 from sklearn.svm import SVC
6 from sklearn.ensemble import RandomForestClassifier
7 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
8 from sklearn.tree import DecisionTreeClassifier
9
10 pipelines=[]
11 for model in [DecisionTreeClassifier(), MultinomialNB(), ComplementNB(),
12              LogisticRegression(solver='saga'), RidgeClassifier(solver='auto'), SVC(),RandomForestClassifier()]:
13     pipeline=make_pipeline(TfidfVectorizer(), model)
14     pipelines.append(pipeline)
```

```
1 #Training the model
```

```
2 import time
3 training_time=[]
4 for pipeline in pipelines:
5     start=time.time()
6     pipeline.fit(X_train, y_train)
7     stop=time.time()
8     training_time.append(stop-start)
```

```
1 #Prediction from test dataset
2 from sklearn.metrics import classification_report, confusion_matrix, f1_score, precision_score, recall_score
3 model_name=[]
4 precision_array=[]
5 recall_array=[]
6 f1_array=[]
7 test_time=[]
8 print("Classification Report\n")
9 print("*****")
10 for i, pipeline in enumerate(pipelines):
11     start=time.time()
12     y_pred=pipeline.predict(X_test)
13     stop=time.time()
14     test_time.append(stop-start)
15     print(pipelines[i].steps[1][0].upper())
16     model_name.append(pipelines[i].steps[1][0].upper())
17     f1_array.append(round(f1_score(y_test, y_pred, average='weighted'),2))
18     precision_array.append(round(precision_score(y_test, y_pred, average='binary'),2))
19     recall_array.append(round(recall_score(y_test, y_pred, average='binary'),2))
20     print("\n",classification_report(y_test, y_pred))
21     print("*****")
```

```
1 #Plotting the various performance metrix of all models
2 training_time=np.array(training_time)/np.max(training_time)
3 test_time=np.array(test_time)/np.max(test_time)
4 score_df=pd.DataFrame({'F1 Score(Weighted)':f1_array,
5                        'Precision Score of Class 1':precision_array,
6                        'Recall of Class 1':recall_array,
7                        'Training Time': training_time,
8                        'Test Time':test_time}, index=model_name)
9
10 f=plt.figure(figsize=(15,10))
11 plt.title('Comparing Performance of various ML Models on Text Classification Dataset', color='black',
12          fontdict={'fontsize':23})
13 score_df.plot(kind='barh', ax=f.gca(), cmap='gnuplot')
14 plt.legend(loc='center left', bbox_to_anchor=(1.0, 0.5))
15 plt.show()
```

Text summary

```
1 import logging
2 logging.basicConfig(format='%(asctime)s : %(levelname)s : %(message)s', level=logging.INFO)
3 from gensim.summarization import summarize
4 # load data
5 filename = '/content/Europe Covid.txt'
6 file = open(filename, 'rt')
7 text = file.read()
8 file.close()
9 print ('Summary:')
10 print (summarize(text))
```

2020-10-25 15:33:47,048 : INFO : 'pattern' package not found; tag filters are not available for English  
2020-10-25 15:33:47,065 : INFO : adding document #0 to Dictionary(0 unique tokens: [])  
2020-10-25 15:33:47,068 : INFO : built Dictionary(171 unique tokens: ['affect', 'case', 'class', 'contin', 'countri']...) from 23 documents (total 257 cc  
Summary:  
European countries are reporting record numbers of Covid-19 cases as the continent prepares for the pandemic to intensify through winter.Those affected i  
Polish President Andrzej Duda tested positive for Covid-19 on Friday, according to a tweet from Presidential Minister Blazej Spychalski Saturday.  
Duda's diagnosis comes as the country reported 13,632 new cases Friday, the highest daily tally since the pandemic began.  
On the same day France reported 42,032 new cases in 24 hours, a new record, according to the French Health Agency.