

Natural Language Processing (COMM061)

Part 1

GROUP 21

TEAM MEMBERS: KRISHNAKUMAR KANNAN, ASHITHA REJITH, ARAVIND RAJU

Model Comparision

```
1 In this python notebook we are combining all the models that
  was done individually, namely:
2
3 ◆ Logistic Regression
4
5 ◆ Decision Tree
6
7 ◆ Random Forest
8
9 ◆ Support Vector Classification
10
11 ◆ Complement Naive Bayes
12
13 ◆ Multinomial Naive Bayes .
14
15 Initial text cleaning is done as necessary, stopwords, special
  characters were removed, text was tokenized, lemmatised and
  then combined together. Data was split using hold out method-
  test-train split using 80-20 ratio. Later they were vectorized
  using TF-IDF vectorizer, and then fed into each of the models.
  The final model with highest accuracy is indicated as the
  conclusion.
```

```
In [1]: 1 #importing libraries
        2
        3 import numpy as np
        4 import pandas as pd
        5 import re
        6 import matplotlib.pyplot as plt
        7 import seaborn as sns
```

```
In [2]: 1 df = pd.read_csv('Research_Article_train.csv') #reading the csv
```

In [3]: 1 df.head() *#displaying the first 5 elements of the dataframe*

Out [3]:

	ID	TITLE	ABSTRACT	Computer Science	Physics	Mathematics	Statistics	Quantitative Biology
0	1	Reconstructing Subject- Specific Effect Maps	Predictive models allow subject- specific inf...	1	0	0	0	0
1	2	Rotation Invariance Neural Network	Rotation invariance and translation invarian...	1	0	0	0	0
2	3	Spherical polyharmonics and Poisson kernels fo...	We introduce and develop the notion of spher...	0	0	1	0	0
3	4	A finite element approximation for the stochas...	The stochastic Landau-- Lifshitz-- Gilbert (LL...	0	0	1	0	0
4	5	Comparative study of Discrete Wavelet Transfor...	Fourier- transform infra-red (FTIR) spectra o...	1	0	0	1	0

```
In [4]: 1 df['text'] = df['TITLE']+" " + df['ABSTRACT']
2 df.drop(['ID', 'TITLE', 'ABSTRACT'], inplace = True, axis = 1)
3 df.head()
```

Out[4]:

	Computer Science	Physics	Mathematics	Statistics	Quantitative Biology	Quantitative Finance	text
0	1	0	0	0	0	0	Reconstructing Subject- Specific Effect Maps ...
1	1	0	0	0	0	0	Rotation Invariance Neural Network Rotation ...
2	0	0	1	0	0	0	Spherical polyharmonics and Poisson kernels fo...
3	0	0	1	0	0	0	A finite element approximation for the stochas...
4	1	0	0	1	0	0	Comparative study of Discrete Wavelet Transfor...

Pre-processing the data

```
In [5]: 1 from nltk.tokenize import TreebankWordTokenizer #tokenizing
2 tree_tokeniser=TreebankWordTokenizer()
```

```
In [6]: 1 from nltk.stem import WordNetLemmatizer #Word net Lemmatize
2 def lema_text(text):
3     lematized_text=[WordNetLemmatizer().lemmatize(i) for i in t
4     return lematized_text
```

```
In [7]: 1 from nltk.corpus import stopwords # importing stopwo
2 stop_words=set(stopwords.words('english'))
3
4
5 def cleaning_stopwords(text):
6     no_stopword_text = [w for w in text.split() if not w in sto
7     return ' '.join(no_stopword_text)
8
```

```
In [8]: 1 def pre_processing(df, text, cleaned_text):
2         df[cleaned_text] = df[text].str.lower() #making each word to
3         df[cleaned_text] = df[cleaned_text].apply(lambda strip: re.sub
4         df[cleaned_text] = df[cleaned_text].apply(lambda strip: re.sub
5         df[cleaned_text] = df[cleaned_text].apply(lambda s: ' '.join([
6         df[cleaned_text] = df[cleaned_text].apply(lambda s: tree_tokenize
7         df[cleaned_text] = df[cleaned_text].apply(lambda s: lemmatize_text
8         df[cleaned_text] = df[cleaned_text].apply(lambda s: ' '.join(
9         return df
```

```
In [9]: 1 df_new = pre_processing(df, 'text', 'cleaned_text')
2         df_new.head()
```

Out [9]:

	Computer Science	Physics	Mathematics	Statistics	Quantitative Biology	Quantitative Finance	text	
0	1	0	0	0	0	0	Reconstructing Subject- Specific Effect Maps ...	r si
1	1	0	0	0	0	0	Rotation Invariance Neural Network Rotation ...	n
2	0	0	1	0	0	0	Spherical polyharmonics and Poisson kernels fo...	p p
3	0	0	1	0	0	0	A finite element approximation for the stochas...	ε
4	1	0	0	1	0	0	Comparative study of Discrete Wavelet Transfor...	ξ

```
In [10]: 1 x=df_new.iloc[:,7] #selecting the input labels - x
          2 x
```

```
Out[10]: 0 reconstructing subjectspecific effect map pred...
          1 rotation invariance neural network rotation in...
          2 spherical polyharmonics poisson kernel polyhar...
          3 finite element approximation stochastic maxwel...
          4 comparative study discrete wavelet transforms ...

          ...
20967 contemporary machine learning guide practition...
20968 uniform diamond coating wcco hard alloy cuttin...
20969 analysing soccer game clustering conceptors pr...
20970 efficient simulation lefttail sum correlated l...
20971 optional stopping problem bayesians recently o...
Name: cleaned_text, Length: 20972, dtype: object
```

```
In [11]: 1 y=df_new.iloc[:,0:6] #selecting the output features - y
          2 y
```

```
Out[11]:
```

	Computer Science	Physics	Mathematics	Statistics	Quantitative Biology	Quantitative Finance
0	1	0	0	0	0	0
1	1	0	0	0	0	0
2	0	0	1	0	0	0
3	0	0	1	0	0	0
4	1	0	0	1	0	0
...
20967	1	1	0	0	0	0
20968	0	1	0	0	0	0
20969	1	0	0	0	0	0
20970	0	0	1	1	0	0
20971	0	0	1	1	0	0

20972 rows × 6 columns

```
In [21]: 1 from sklearn.model_selection import train_test_split
          2
          3 x_train,x_test,y_train,y_test = train_test_split(x,y, test_size
```

```
In [22]: 1 from sklearn.feature_extraction.text import TfidfVectorizer
2 tfidf_vectoriser=TfidfVectorizer(min_df=5, max_df=0.9, token_pa
3                                     max_features=10000,ngram_range
4 x_train_tvect=tfidf_vectoriser.fit_transform(x_train)
5 x_test_tvect=tfidf_vectoriser.transform(x_test)
```

1. Logistic regression for multi-label classification using a one-vs-rest

```
In [23]: 1 from sklearn.linear_model import LogisticRegression
2 from sklearn.multiclass import OneVsRestClassifier
3
4 # define model
5 model = LogisticRegression()
6
7 # define the ovr strategy
8 ovr = OneVsRestClassifier(model)
9
10 # fit model
11 ovr.fit(x_train_tvect, y_train)
12
13 # make predictions
14 yhat = ovr.predict(x_test_tvect)
15
16 from sklearn.metrics import accuracy_score
17
18 # View accuracy score
19 accuracy_m1 = accuracy_score(y_test, yhat)
20 print('The accuracy of Logitic Regression model with tfidf vect
```

The accuracy of Logitic Regression model with tfidf vectorizing is 0.6522050059594756

2. Decision tree classifier

```
In [15]: 1  from sklearn.tree import DecisionTreeClassifier
2
3  # calling the decision tree into a parameter named classifier
4  classifier = DecisionTreeClassifier()
5
6  #training the model by fitting the TFIDFvectorized data onto th
7  classifier.fit(x_train_tvect,y_train)
8
9
10 #passing on the TFIDF vectorized test data (x_test_tvect) to ev
11 yhat = classifier.predict(x_test_tvect)
12
13
14 from sklearn.metrics import accuracy_score
15
16 # View accuracy score
17 accuracy_m2 = accuracy_score(y_test, yhat)
18 print('The accuracy of Decision Tree model with tfidf vectorizi
```

The accuracy of Decision Tree model with tfidf vectorizing is 0.5022646007151371

3. Random Forest Classifier

```
In [16]: 1  from sklearn.ensemble import RandomForestClassifier
2
3  forest= RandomForestClassifier(n_estimators = 200)
4  model = forest.fit(x_train_tvect, y_train)
5
6  yhat= model.predict(x_test_tvect)
7  yhat
8
9  from sklearn.metrics import accuracy_score
10 # View accuracy score
11 accuracy_m3 = accuracy_score(y_test, yhat)
12 print('The accuracy of Random Forest model with tfidf vectorizi
```

The accuracy of Random Forest model with tfidf vectorizing is 0.5973778307508939

4. Multinomial Naive Bayes

```
In [17]: 1 from sklearn.naive_bayes import MultinomialNB
2
3 model_nb=MultinomialNB()
4 model_nb=OneVsRestClassifier(model_nb)
5 model_nb.fit(x_train_tvect,y_train)
6 yhat = model_nb.predict(x_test_tvect)
7
8 accuracy_m4 = accuracy_score(y_test, yhat)
9 print('The accuracy of Multinomial Naive Bayes model with tfidf
```

The accuracy of Multinomial Naive Bayes model with tfidf vectorizing is 0.6421930870083432

5. Complement Naive Bayes

```
In [18]: 1 from sklearn.naive_bayes import ComplementNB
2
3 model=ComplementNB()
4 model=OneVsRestClassifier(model_nb)
5 model.fit(x_train_tvect,y_train)
6 yhat = model_nb.predict(x_test_tvect)
7
8 accuracy_m5 = accuracy_score(y_test, yhat)
9 print('The accuracy of Complement Naive Bayes model with tfidf
```

The accuracy of Complement Naive Bayes model with tfidf vectorizing is 0.6421930870083432

6. Linear Support vector classification

```
In [19]: 1 from sklearn.svm import LinearSVC
2 from sklearn.multiclass import OneVsRestClassifier
3
4
5 model= LinearSVC()
6 ovr=OneVsRestClassifier(model)
7 ovr.fit(x_train_tvect,y_train)
8 yhat = ovr.predict(x_test_tvect)
9
10 accuracy_m6 = accuracy_score(y_test, yhat)
11 print('The accuracy of SVC model with tfidf vectorizing is ',ac
```

The accuracy of SVC model with tfidf vectorizing is 0.6424314660309892

Comparing the best among the 6 models

```
In [24]: 1 classification = {0:'Logistic Regression',1:'Decision Tree', 2:  
2           4: 'Gaussian Naive Bayes', 5:'Support Vector  
3 array1= [accuracy_m1,accuracy_m2,accuracy_m3, accuracy_m4, accu  
4 acc1 = np.argmax([accuracy_m1,accuracy_m2,accuracy_m3,accuracy_  
5 print('Best classifier according to TF IDF Vectorization is {a}
```

Best classifier according to TF IDF Vectorization is Logistic Regression with accuracy 0.6522

Conclusion

- ◆ From the results it can be concluded that Logistic Regression model with TF-IDF vectorization technique yields the best accuracy of approximately 65.22%. We further conduct the experiments based on this model and accuracy.
- ◆ In our comparison to other techniques for supervised classification such as SVMs or ensemble methods, logistic regression is rather fast and with a better accuracy.
- ◆ By applying a logarithmic transformation to the outcome variable, we can model a nonlinear association linearly

```
In [ ]: 1
```