

20 Newsgroups Text Classification

AML-2304 Natural Language Processing

Final Project Report

Prof. Harriet Huang

Group: Valar Morghulis

- Kishan Lakhtariya C0849341
- Harshkumar Patel C0848343
- Kanan Trivedi C0851492
- Mohammaddin Masbi C0850500

Abstract

Analyzed and processed text data of 20 newsgroups text datasets with different techniques to get meaningful words from the dataset and converted them into word vectors using two different word embedding strategies, TF-IDF and Word2Vec, and trained the vectorized data using three different machine learning models, from all the training we got maximum training accuracy of 96% and testing accuracy of 71% using Support Vector Machine Classifier with the vectorizer TF-IDF.

1. Introduction

- It is essential to classify the new topics for the chatbots and AI assistants to give recommendations based on the reader's interest in addressing this problem; we have chosen the news dataset with 20 different topics as labels and trained the model to predict the news text.

2. Dataset

- We have used the fetch_20newsgroups API to get the dataset which contains 18846 samples with 2 columns of data and a target.

```
1 # Fetching and Storing all data into news_dataset variable with shuffling
2 news_dataset = fetch_20newsgroups(subset='all',
3                                 shuffle=True,
4                                 remove=('headers', 'footers', 'quotes'))
```

- As the target is the integer value, we have mapped the target_labels with it for further analysis.

```
1 # Mapping labels with its integer values (Label Encoding)
2 label_map = dict((i, j) for i, j in enumerate(news_dataset.target_names))
3 label_map

{0: 'alt.atheism',
 1: 'comp.graphics',
 2: 'comp.os.ms-windows.misc',
 3: 'comp.sys.ibm.pc.hardware',
 4: 'comp.sys.mac.hardware',
 5: 'comp.windows.x',
 6: 'misc.forsale',
 7: 'rec.autos',
 8: 'rec.motorcycles',
 9: 'rec.sport.baseball',
10: 'rec.sport.hockey',
11: 'sci.crypt',
12: 'sci.electronics',
13: 'sci.med',
14: 'sci.space',
15: 'soc.religion.christian',
16: 'talk.politics.guns',
17: 'talk.politics.mideast',
18: 'talk.politics.misc',
19: 'talk.religion.misc'}
```

- The initial shape of the data is:

```
1 df.shape

(18846, 3)
```

```
1 # Converting train dataset into pandas dataframe and storing into train_df variable
2 df = pd.DataFrame({'data': news_dataset.data, 'target': news_dataset.target})
3 df["target_label"] = df["target"].map(label_map)
4 df.head()

          data  target  target_label
0 \n\nI am sure some bashers of Pens fans are pr...      10  rec.sport.hockey
1 My brother is in the market for a high-perform...       3  comp.sys.ibm.pc.hardware
2 \n\n\n\nFinally you said what you dream abou...      17  talk.politics.mideast
3 \nThink!\n\nIt's the SCSI card doing the DMA t...       3  comp.sys.ibm.pc.hardware
4  1) I have an old Jasmine drive which I cann...       4  comp.sys.mac.hardware
```

3. Exploratory Data Analysis (EDA):

- We have three data initially
 - i. Data: which contains the text of the news
 - ii. Target: which is the integer value of labels
 - iii. Target_labels: which is the string of the labels

		data	target	target_label
0	\n\nI am sure some bashers of Pens fans are pr...	10		rec.sport.hockey
1	My brother is in the market for a high-perform...	3		comp.sys.ibm.pc.hardware
2	\n\n\n\tFinally you said what you dream abou...	17		talk.politics.mideast
3	\nThink!\nIt's the SCSI card doing the DMA t...	3		comp.sys.ibm.pc.hardware
4	1) I have an old Jasmine drive which I cann...	4		comp.sys.mac.hardware
...
18841	DN> From: nyeda@cnsvax.uwec.edu (David Nye)\nD...	13		sci.med
18842	\nNot in isolated ground recepticles (usually ...	12		sci.electronics
18843	I just installed a DX2-66 CPU in a clone mothe...	3		comp.sys.ibm.pc.hardware
18844	\nWouldn't this require a hyper-sphere. In 3...	1		comp.graphics
18845	After a tip from Gary Crum (crum@fcom.cc.utah....	7		rec.autos
18846	rows x 3 columns			

- The total count and frequency of the labels are:

Category	Count
rec.sport.hockey	999
soc.religion.christian	997
rec.motorcycles	996
rec.sport.baseball	994
sci.crypt	991
rec.autos	990
sci.med	990
comp.windows.x	988
sci.space	987
comp.os.ms-windows.misc	985
sci.electronics	984
comp.sys.ibm.pc.hardware	982
misc.forsale	975
comp.graphics	973
comp.sys.mac.hardware	963
talk.politics.mideast	940
talk.politics.guns	910
alt.atheism	799
talk.politics.misc	775
talk.religion.misc	628

4. Preprocessing

- We performed pre-processing steps to convert text data into meaningful data to make the computer understandable the text data.
 - First, we plot the word cloud to know the frequency of the words, and we discover that many un-useful words have higher frequencies than essential words.

- Here we can see that the occurrence of the “n” is higher in the text data.
- So, we must perform data cleaning, including removing symbols and ASCII characters, stopping word removal, steaming and lemmatizing.
- We tokenized the data using NLTK’s word_tokenize function and stored all tokenized words in separate columns.

```

1 # Converting text paragraph into tokens using nltk.word_tokenize function and also removing special characters and un-usefull text
2 # And returning tokens list
3 def tokenize(data):
4     tokens = nltk.word_tokenize(data)
5     return [word for word in tokens if word.isalpha()]

1 # Applying tokenize function to data column of dataset and storing results into column name tokenized
2 df['tokenized'] = df.apply(lambda x: tokenize(x['data']), axis=1)
3 df.head()

          data      target      target_label      tokenized
0 \n\nI am sure some bashers of Pens fans are pr...    10      rec.sport.hockey [I, am, sure, some, bashers, of, Pens, fans, a...
1 My brother is in the market for a high-perform...     3      comp.sys.ibm.pc.hardware [My, brother, is, in, the, market, for, a, vid...
2 \n\n\n\nFinally you said what you dream abou...    17      talk.politics.mideast [Finally, you, said, what, you, dream, about, ...
3 \n\nThink!\n\nIt's the SCSI card doing the DMA t...     3      comp.sys.ibm.pc.hardware [Think, It, the, SCSI, card, doing, the, DMA, ...
4  1) I have an old Jasmine drive which I cann...     4      comp.sys.mac.hardware [I, have, an, old, Jasmine, drive, which, I, c...

```

- We applied the English stop word removal on the tokenized word using “stopwords.words()” of the NLTK function.

```

1 # Defining Function to remove english stopwords from tokenized data with the help of stopwords function of NLTK
2 def remove_stopword(data):
3
4     stopword = stopwords.words("english")
5     return [word for word in data if not word in stopword]

1 # Applying remove_stopword function to column tokenized of dataset and storing returned values into stopword_removed column
2 df['stopword_removed'] = df.apply(lambda x: remove_stopword(x['tokenized']),
3                                     axis=1)
4 df[['data', "stopword_removed"]].head()

          data      stopword_removed
0 \n\nI am sure some bashers of Pens fans are pr... [I, sure, bashers, Pens, fans, pretty, confuse...
1 My brother is in the market for a high-perform... [My, brother, market, video, card, supports, V...
2 \n\n\n\nFinally you said what you dream abou... [Finally, said, dream, Mediterranean, That, ne...
3 \n\nThink!\n\nIt's the SCSI card doing the DMA t... [Think, It, SCSI, card, DMA, transfers, NOT, d...
4  1) I have an old Jasmine drive which I cann... [I, old, Jasmine, drive, I, use, new, system, ...

```

- After removing the stop words, we convert the words to their root form, which is called stemming

```

1 # Converting words to it's root form using stemming technique
2 # Word stemming using PorterStemmer function of NLTK
3 # And return the stemmed word list
4
5
6 def stemming(data):
7
8     stemmer = PorterStemmer()
9     return [stemmer.stem(word) for word in data]

1 # Applying stemming function to dataset with column name stopword_removed
2 df['stemmed_words'] = df.apply(lambda x: stemming(x['stopword_removed']),
3                                 axis=1)
4 df[['data', 'stemmed_words']].head()

          data      stemmed_words
0 \n\nI am sure some bashers of Pens fans are pr... [i, sure, basher, pen, fan, pretti, confus, la...
1 My brother is in the market for a high-perform... [my, brother, market, video, card, support, ve...
2 \n\n\n\nFinally you said what you dream abou... [final, said, dream, mediterranean, that, new, ...
3 \n\nThink!\n\nIt's the SCSI card doing the DMA t... [think, it, scsi, card, dma, transfer, not, di...
4  1) I have an old Jasmine drive which I cann... [i, old, jasmin, drive, i, use, new, system, m...

```

- And also applied the lemmatizing to the stop words.

```

1 # Lemmatization also convert words to its root form
2 # But the only difference between lemmatization and stemming is that stemming just removes or stems
3 # lemmatization considers the context and converts the word to its meaningful base form, which is c
4
5
6 def lemmatization(data):
7     lemmatizer = WordNetLemmatizer()
8     return [lemmatizer.lemmatize(word) for word in data]

1 df['lemmatized_words'] = df.apply(lambda x: stemming(x['stopword_removed']),
2                                     axis=1)
3 df[['data', 'lemmatized_words']].head()

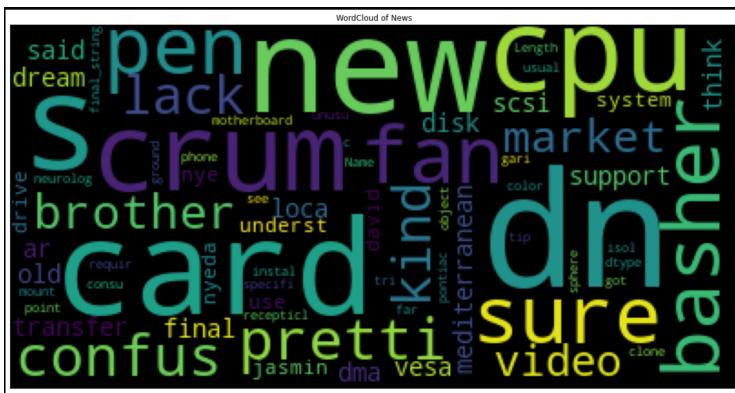
          data      lemmatized_words
0 \n\nI am sure some bashers of Pens fans are pr... [i, sure, basher, pen, fan, pretti, confus, la...
1 My brother is in the market for a high-perform... [my, brother, market, video, card, support, ve...
2 \n\n\n\nFinally you said what you dream abou... [final, said, dream, mediterranean, that, new, ...
3 \n\nThink!\n\nIt's the SCSI card doing the DMA t... [think, it, scsi, card, dma, transfer, not, di...
4  1) I have an old Jasmine drive which I cann... [i, old, jasmin, drive, i, use, new, system, m...

```

- The only difference between lemmatization and stemming is that stemming removes or stems the last few characters of a word, often leading to incorrect meanings and spelling.

- Lemmatization considers the context and converts the word to its meaningful base form, which is called Lemma.
 - The final step of the preprocessing is the joining of the lemmatized_words into a string.

- At the end of preprocessing, we again generated a word cloud on final_string, and we got a clean and clear view of the words.



5. Feature Extraction:

- Extracting features from the text data is called word vectorization or word embeddings. This is a very crucial part of Natural Language Processing because this feature vector will go into the model training. If this is not done correctly, the model will lead to misprediction.
 - Word Vectorization is the technique of converting words to numerical form to perform model training on text-based data.
 - For this problem, we have used two-word vectorizing techniques.

1. TF-IDF Vectorizer:

- a. Term Frequency-Inverse Document Frequency is a statistical measure that evaluates how relevant a word is to a document in a collection of documents.
 - b. This is done by multiplying two metrics: how many times a word appears in a document and the inverse document frequency of the word across a set of documents.
 - c. Here we have used the TfidfVectorizer class of the Sklearn library to convert the words to vectors.

```
1 # Creating instance of TfIdfVectorizer function of sklearn library
2 tf_idf = TfIdfVectorizer()
3 # Converting final_strings to its meaning full words vectors from dataframe
4 # Storing final traing data to X variable
5 X = tf_idf.fit_transform(df["final_string"])
6 # Storing Labels to Y variable
7 Y = df["target"]
8 # Printing Shape of data
9 X.shape
```

(18846, 57689)

```

1 # Displaying average non-zero values from dataset
2 X.nnz / float(X.shape[0])

```

e. 61.29640241961159

- f. Here we discovered that from all features, 61.29% of data is none zero, and 39% of features contain zeros because of the uneven length of news text.
- g. To train and evaluate the model, we split the dataset into train-test with a ratio of 80-20%.
- h. We have trained the model with three different classifiers.

i. Naive Bayes Classifier

- With Naive Bayes, we got a training accuracy of 81% and test accuracy of 69%

```

1 # Defining naive bayes classifier model with
2 # fitting train data into model
3
4 naive_bayes_classifier = MultinomialNB()
5 naive_bayes_classifier.fit(X_train, y_train)

```

● MultinomialNB
MultinomialNB()

```

1 # getting training accuracy of model
2 y_pred = naive_bayes_classifier.predict(X_train)
3 NB_train_acc = metrics.accuracy_score(y_train, y_pred)
4 print("Train Accuracy : ", NB_train_acc)
5 tf_idf_train_acc["NaiveBayes"] = NB_train_acc
6
7 # Performing prediction on test dataset
8 y_pred = naive_bayes_classifier.predict(X_test)
9 NB_test_acc = metrics.accuracy_score(y_test, y_pred)
10 print("Test Accuracy : ", NB_test_acc)
11 tf_idf_test_acc["NaiveBayes"] = NB_test_acc

```

Train Accuracy : 0.8114884584770496
Test Accuracy : 0.6896551724137931

ii. SGD Classifier

- With SGD, we got a training accuracy of 85% and test accuracy of 70%

```

1 # Creating model of SGD Classifier with SGDClassifier class of sklearn
2 # Fitting the training dataset into model
3 sgd = SGDClassifier(loss='hinge', penalty='l2', alpha=1e-3, random_state=42)
4 sgd.fit(X_train, y_train)

```

● SGDClassifier
SGDClassifier(alpha=0.001, random_state=42)

```

1 # getting training accuracy of model
2 y_pred = sgd.predict(X_train)
3 SGD_train_acc = metrics.accuracy_score(y_train, y_pred)
4 print("Train Accuracy : ", SGD_train_acc)
5 tf_idf_train_acc["SGD"] = SGD_train_acc
6
7 # Performing prediction on test dataset
8 y_pred = sgd.predict(X_test)
9 SGD_test_acc = metrics.accuracy_score(y_test, y_pred)
10 print("Test Accuracy : ", SGD_test_acc)
11 tf_idf_test_acc["SGD"] = SGD_test_acc

```

Train Accuracy : 0.857389227911913
Test Accuracy : 0.7026525198938992

iii. Support Vector Machine(SVM) Classifier

- With SVM, we got a training accuracy of 96% and test accuracy of 71%

```

1 svm = SVC()
2 svm.fit(X_train, y_train)

```

● SVC
SVC()

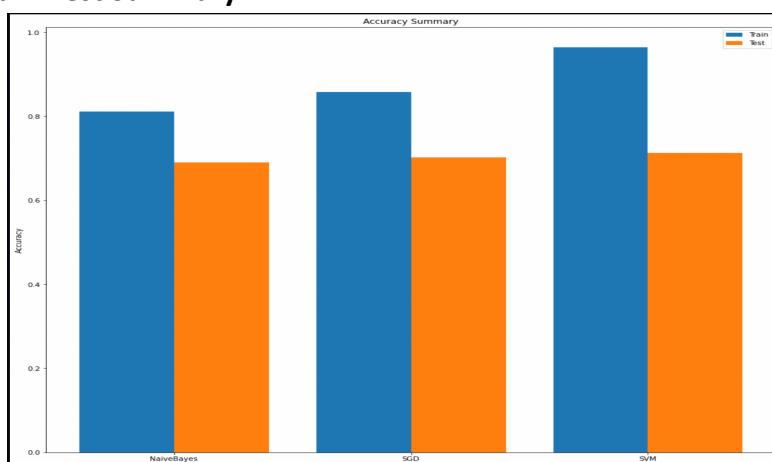
```

1 # getting training accuracy of model
2 y_pred = svm.predict(X_train)
3 SVM_train_acc = metrics.accuracy_score(y_train, y_pred)
4 print("Train Accuracy : ", SVM_train_acc)
5 tf_idf_train_acc["SVM"] = SVM_train_acc
6
7 # Performing prediction on test dataset
8 y_pred = svm.predict(X_test)
9 SVM_test_acc = metrics.accuracy_score(y_test, y_pred)
10 print("Test Accuracy : ", SVM_test_acc)
11 tf_idf_test_acc["SVM"] = SVM_test_acc

```

Train Accuracy : 0.964181480498806
Test Accuracy : 0.7124668435013263

iv. Train-Test Summary:



- From all three classification models, Support Vector Machine Classifier performed well with train accuracy of 96% and test accuracy of 71% other than Naïve Bayes with 81%-69% Train-Test and SGD with 85%-70% train-test accuracy.

2. Word2Vect Vectorizer:

- Word2vec is a technique for natural language processing to convert words into vectors. The word2vec algorithm uses a neural network model to learn word associations from a large corpus of text. Once trained, such a model can detect synonymous words or suggest different words for a partial sentence.

```

1 # Splitting the lemmatized_words list into train test dataset
2 X_train, X_test, y_train, y_test = train_test_split(df['lemmatized_words'],
3                                                 df['target'],
4                                                 test_size=0.2)

```

```

1 # Train the word2vec model with vector_size 800 and window size 600
2 w2v_model = gensim.models.Word2Vec(X_train,
3                                     vector_size=500,
4                                     window=100,
5                                     min_count=20)

```

- Here we have taken the lemmatized_words for training the network.
- We have trained the word2vec model using Gensim's Word2Vec function with vector_size=500 and window=100 with a min_count of 20.

```

1 # Generating aggregated sentence vectors based on the word vectors for each word in the sentence
2 words = set(w2v_model.wv.index_to_key)
3 X_train_vect = np.array([
4     [np.array([w2v_model.wv[i] for i in ls if i in words]) for ls in X_train])
5 X_test_vect = np.array([
6     [np.array([w2v_model.wv[i] for i in ls if i in words]) for ls in X_test]])

```

- Generating aggregated sentence vectors based on the word vectors for each word in the sentence

```

1 # Compute sentence vectors by averaging the word vectors for the words contained in the sentence
2 X_train_vect_avg = []
3 for v in X_train_vect:
4     if v.size:
5         X_train_vect_avg.append(v.mean(axis=0))
6     else:
7         X_train_vect_avg.append(np.zeros(500, dtype=float))
8
9 X_test_vect_avg = []
10 for v in X_test_vect:
11     if v.size:
12         X_test_vect_avg.append(v.mean(axis=0))
13     else:
14         X_test_vect_avg.append(np.zeros(500, dtype=float))

```

- Compute final sentence vectors by averaging the word vectors for the words contained in the sentence
- We have performed model training with three different classifiers.

i. Stochastic Gradient Descent(SGD) Classifier

- We have achieved a training accuracy of 59.9 % and a test accuracy of 59.3% with the SGD Classifier.

```

1 # getting training accuracy of model
2 y_pred = sgd.predict(X_train_vect_avg)
3 SGD_train_acc = metrics.accuracy_score(y_train, y_pred)
4 print("Train Accuracy : ", SGD_train_acc)
5 w2v_train_acc["SGD"] = SGD_train_acc
6
7 # Performing prediction on test dataset
8 y_pred = sgd.predict(X_test_vect_avg)
9 SGD_test_acc = metrics.accuracy_score(y_test, y_pred)
10 print("Test Accuracy : ", SGD_test_acc)
11 w2v_test_acc["SGD"] = SGD_test_acc

```

Train Accuracy : 0.599230565136641
Test Accuracy : 0.5952254641909814

ii. Random Forest Classifier

- With a random forest classifier, we got a training accuracy of 97% and test accuracy of 61%

- ```

1 # Train data with randomforest classifier
2 rf = RandomForestClassifier()
3 rf_model = rf.fit(X_train_vect_avg, y_train)

1 # getting training accuracy of model
2 y_pred = rf_model.predict(X_train_vect_avg)
3 rf_train_acc = metrics.accuracy_score(y_train, y_pred)
4 print("Train Accuracy : ", rf_train_acc)
5 w2v_train_acc["RandomForest"] = rf_train_acc
6
7 # Performing prediction on test dataset
8 y_pred = rf_model.predict(X_test_vect_avg)
9 rf_test_acc = metrics.accuracy_score(y_test, y_pred)
10 print("Test Accuracy : ", rf_test_acc)
11 w2v_test_acc["RandomForest"] = rf_test_acc

Train Accuracy : 0.9710798620323693
Test Accuracy : 0.6127320954907162

```

### iii. Support Vector Machine(SVM) Classifier

- Using the SVM classifier, we got a training accuracy of 96% and test accuracy of 71%

- ```

1 svm = SVC()
2 svm.fit(X_train_vect_avg, y_train)

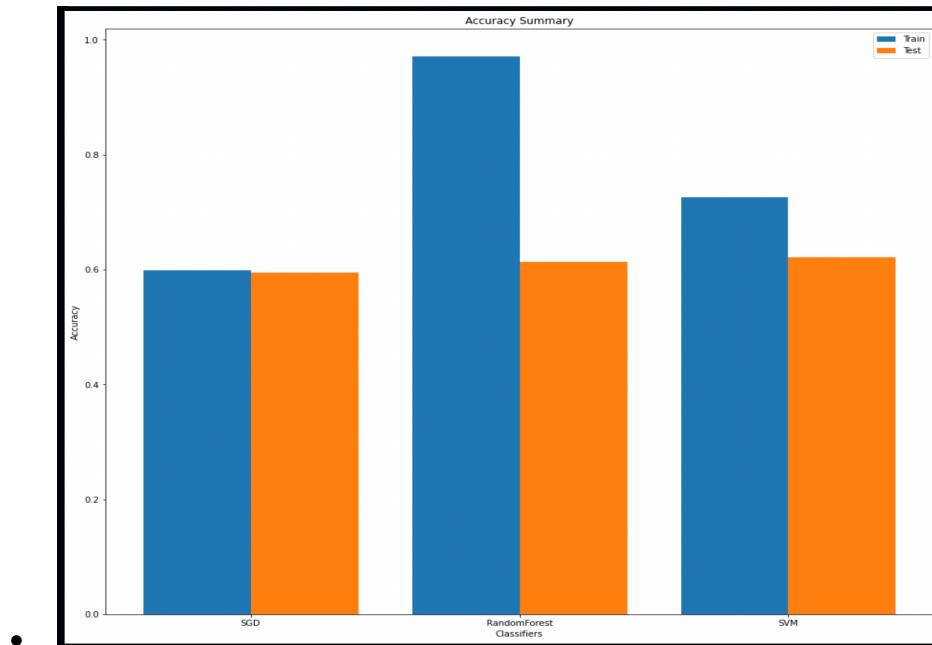
1 # getting training accuracy of model
2 y_pred = svm.predict(X_train_vect_avg)
3 SVM_train_acc = metrics.accuracy_score(y_train, y_pred)
4 print("Train Accuracy : ", SVM_train_acc)
5 w2v_train_acc["SVM"] = SVM_train_acc
6
7 # Performing prediction on test dataset
8 y_pred = svm.predict(X_test_vect_avg)
9 SVM_test_acc = metrics.accuracy_score(y_test, y_pred)
10 print("Test Accuracy : ", SVM_test_acc)
11 w2v_test_acc["SVM"] = SVM_test_acc

Train Accuracy : 0.7263863093658796
Test Accuracy : 0.6209549071618037

```

iv. Classifier Summary

- For the word2vec word embedding technique from all three classification models, the Support Vector Machine classifier performed well with the train-test, 73%-61% accuracy score,



Conclusion

- To conclude, we performed classification on 20 newsgroup datasets with 2 different word embedding techniques with all text pre-processing steps with feature selection to train the classification models.
- We observe that the TF-IDF vectorization technique performed better for this problem, with a 71% accuracy score higher than Word2Vec embedding using the SVM Classifier.
- Because the Word2Vec technique is most suitable for recommendation problems, but this is the problem of Multiclass text classification, so for that TI-IDF is best.