

▼ 🇬🇧 BBC NEWS CLASSIFICATION PROJECT

(Unsupervised Learning Project)

🏫 University of Colorado, Boulder - Unsupervised Algorithms in Machine Learning

By Mattison Hinline

Overview

This project looks at an unsupervised learning model technique called matrix factorization. All the required libraries to run the notebook are in the first coding cell. This notebook explores the training data while using common natural language processing techniques. Our goal is to classify different news articles into five different groups: business, tech, sport, entertainment or politics. We will train two models: (1) an unsupervised model using matrix factorization and then compare it with (2) a supervised model using KMeans clustering. Both models are submitted to have their testing accuracy.

1. Exploratory Data Analysis (EDA) 🧐

In this section we will explore the data given for the competition. It is important to look at what you have to work with before beginning any model building or model training. Normally, it is suggested to first split the data into training and testing sets before any EDA because you, as a researcher, do not want to be biased for the results which can lead to overfitting. For this reason, we will not visualize or explore the test data given, and use that data only for testing the models.

2. Model Building and Training 🏗️

Next, we will build multiple unsupervised models and train them on the test data. Once we have at least two strong models built, we will fine-tune those models and test on the test dataset. We will be using the train-test split on this dataset that is already given in the CSV files.

3. Model Comparisons 🤖

Once we have good, working models we will compare them for their performance. We will discuss which model is the best out of the models tested and why.

4. Conclusions 🧐

Finally, we will summary the project and results, including future suggestions for further analysis.

```
1 #import important libraries
2 import numpy as np
3 import pandas as pd
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 import os
7
```

```

8 #EDA and preprocessing
9 import re
10 import nltk.corpus
11 from nltk.corpus import stopwords
12 from nltk.tokenize import word_tokenize
13 from nltk.stem import WordNetLemmatizer
14 from string import digits
15
16 #modeling
17 from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
18 from sklearn.decomposition import NMF
19 from sklearn.metrics import accuracy_score
20 import sklearn.metrics as metrics
21 import itertools
22 from sklearn.cluster import KMeans
23 from sklearn.model_selection import train_test_split
24
25 #find the files names
26 for dirname, _, filenames in os.walk('/kaggle/input'):
27     for filename in filenames:
28         print(os.path.join(dirname, filename))

```

/kaggle/input/learn-ai-bbc/BBC News Train.csv
 /kaggle/input/learn-ai-bbc/BBC News Sample Solution.csv
 /kaggle/input/learn-ai-bbc/BBC News Test.csv

```

1 #label file paths
2 path_dir = '/kaggle/input/learn-ai-bbc/'
3 train_path = path_dir + 'BBC News Train.csv'
4 sample_solution_path = path_dir + 'BBC News Sample Solution.csv'
5 test_path = path_dir + 'BBC News Test.csv'

```

```

1 #import data
2 train = pd.read_csv(train_path)
3 sample_solution = pd.read_csv(sample_solution_path)
4 test = pd.read_csv(test_path)

```

▼ 1. Exploratory Analysis (EDA)

```

1 #look at what we are trying to submit
2 sample_solution

```

	ArticleId	Category
0	1018	sport
1	1319	tech
2	1138	business
3	459	entertainment
4	1020	politics
...
730	1923	sport

After importing the data, we first take a look at what our ultimate goal for the project is by looking at the submission dataframe (above). We need to create a dataframe comprising of classified articles (ArticleId) and which category they belong to (category). We also see that each ArticleID is unique, while the categories repeat themselves. This is good information to know before continuing. Now that we know our end goal structure, let's move on to the training data.

We can see in the dataframe below that we have three columns:

1. ArticleID: which is the identifying number for the article
2. Test: the article header and text
3. Category: category given to the article

```
1 #look at the training data
2 train
```

ArticleId

Text

Category

```
1 # gain more information from the dataframe
2 train.describe()
```

ArticleId	
count	1490.000000
mean	1119.696644
std	641.826283
min	2.000000
25%	565.250000
50%	1112.500000
75%	1680.750000
max	2224.000000

1490 rows x 3 columns

```
1 #check the type of data, null value counts and number of entries
2 train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1490 entries, 0 to 1489
Data columns (total 3 columns):
#   Column      Non-Null Count  Dtype
---  -
0   ArticleId   1490 non-null   int64
1   Text        1490 non-null   object
2   Category    1490 non-null   object
dtypes: int64(1), object(2)
memory usage: 35.0+ KB
```

Great. All this information look as we expected. We have two object columns and one integer column (for the ID's). It looks like we are not missing any rows. Since the data we are working with is text, we don't need to worry about numbers that are missing such as 9999, 0, etc. Before continuing, I want to make sure that we have no repeated articles in the data. From the code below, we can see that we have 1490 unique IDs and we know the dataframe has 1490 rows, we can assume that each article is unique and continue. I also wanted to see how many categories there are. We can see below that there are five categories in total: business, tech, politics, sport, entertainment.

```
1 # check for repeated articles
2 train['ArticleId'].nunique()
```

1490

```
1 train['Category'].unique()
```

```
array(['business', 'tech', 'politics', 'sport', 'entertainment'],  
      dtype=object)
```

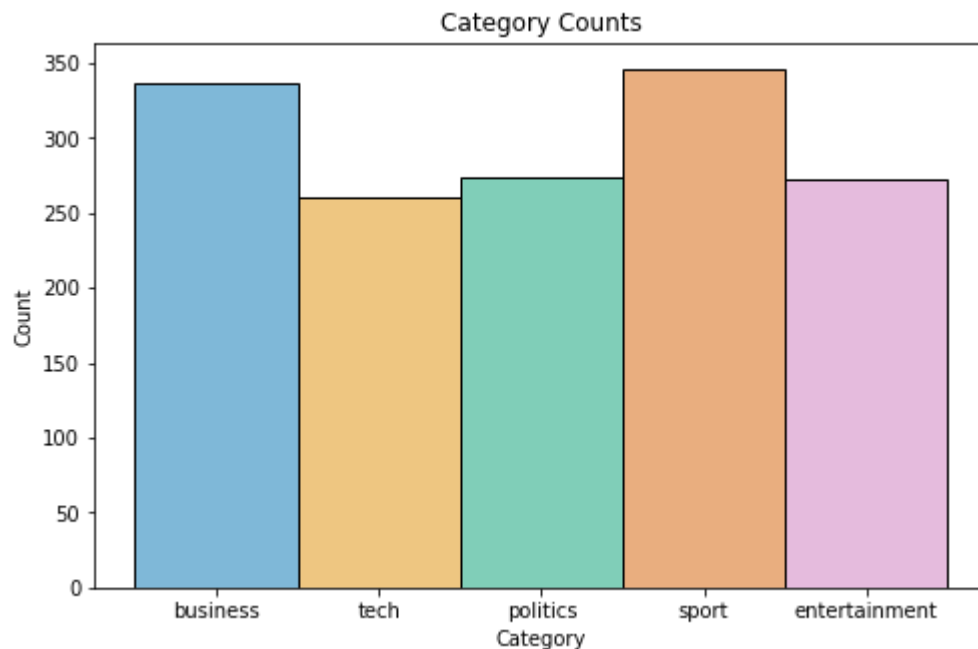
The first thing I want to look at is a couple of texts to get an idea of how they are saved in the dataframe. We will look at the first row. We can see that the text starts with the header, then has a fairly good amount of text afterwards which would be the article text. We can see that the data has already been preprocessed a bit because there are no uppercase characters. Since these are news articles, we can also assume that there are no spelling mistakes. Capitalization and spelling are two important factors when it comes to natural language processing. For this reason, we will be thankful that this has already been processed as such.

```
1 # first row  
2 train['Text'][0]
```

'worldcom ex-boss launches defence lawyers defending former worldcom chief bernie ebbers against a battery of fraud charges have called a company whistleblower as their first witness. cynthia cooper worldcom s ex-head of internal accounting alerted directors to irregular accounting practices at the us telecoms giant in 2002. her warnings led to the collapse of the firm following the discovery of an \$11bn (£5.7bn) accounting fraud. mr ebbers has pleaded not guilty to charges of fraud and conspiracy. prosecution lawyers have argued that mr ebbers orchestrated a series of accounting tricks at worldcom ordering employees to hide expenses and inflate revenues to meet wall street earnings estimates. but ms cooper who now runs her own consulting business told a jury in new york on wednesday that external auditors arthur andersen had approved worldcom s accounting in early 2001 and 2002. she said andersen had given a green light to the procedures and practices used by worldcom. mr ebber s lawyers have said he was unaware of the fraud arguing that auditors did not alert him to any problems. ms cooper also said that during shareholder meetings mr ebbers often passed over technical questions to the company s finance chief giving only brief answers himself. the prosecution s star witness former worldcom financial chief scott sullivan has said that mr ebbers ordered accounting adjustments at the firm telling him to hit our books . however ms cooper said mr sullivan had not mentioned anything uncomfortable about worldcom s accounting during a 2001 audit committee meeting. mr ebbers could face a jail sentence of 85 years if convicted of all the charges he is facing. worldcom emerged from bankruptcy protection in 2004 and is now known as mci. last week mci agreed to a buyout by verizon communications in a deal valued at \$6.75bn.'

Now let's visualize the data as much as we can before running some models. We can see that overall we have about even number of entries for each category. This is good because if one or two categories was severely underrepresented or, in contrast, overrepresentative in the data, then it may cause our model to be biased and/or perform poorly on some or all of the test data.

```
1 fig, ax = plt.subplots(figsize=(8, 5))  
2 sns.histplot(  
3     data = train,  
4     x = 'Category',  
5     hue = 'Category',  
6     palette = 'colorblind',  
7     legend = False,  
8     ).set(  
9         title = 'Category Counts');
```



While working with text data, it is important to make the text "readable" for the computer. To do this, we will take three steps:

1. remove punctuation
2. remove stop words (common English words such as 'to', 'the', 'of', etc

```

1 def clean_text(dataframe, text_col):
2     """
3     A helper function which takes a dataframe
4     and removes punctuation and stopwords.
5     """
6     #remove all punctuation
7     dataframe['no_punct'] = dataframe[text_col].apply(lambda row: re.sub(r'^\w\s+', "", row))
8
9     #remove numbers
10    dataframe['no_punct_num'] = dataframe['no_punct'].apply(lambda row: re.sub(r'[0-9]+', "", row))
11
12    #remove stopwords
13    stop_words = stopwords.words('english')
14    dataframe['no_stopwords'] = dataframe['no_punct_num'].apply(lambda x: ' '.join([word for word in x.split() if word not in stop_words]))
15
16    #remove extra spaces
17    dataframe['clean_text'] = dataframe['no_stopwords'].apply(lambda x: re.sub(' +', ' ', x))
18    return dataframe['clean_text']

1 #clean dataframe text column
2 clean_text(train, 'Text')

1 train['clean_text'][1]

```

'german business confidence slides german business confidence fell february knocking hopes speedy recovery europe largest economy munichbased research institute ifo said confidence index fell february january first decline three months study found outlook manufacturing retail sectors worsened

observers hoping confident business sector would signal economic activity picking surprised ifo index taken knock said dz bank economist bernd weidensteiner main reason probably domestic economy still weak particularly retail trade economy labour minister wolfgang clement called dip february ifo confidence figure mild decline said despite retreat index remained relatively high level expected modest economic upswing continue germany economy grew last year shrinking however economy contracted last three months mainly due reluctance consumers spend latest indications growth still proving elusive ifo president hanswerner sinn said improvement german domestic demand sluggish exports kept things going first half demand exports hit value euro hit record levels making german products less competitive overseas top unemployment rate stuck close manufacturing firms including daimlerchrysler siemens volkswagen negotiating unions cost cutting measures analysts said ifo figures germany continuing problems may delay interest rate rise european central bank eurozone interest rates comments senior officials recently focused threat inflation prompting fears interest rates may rise'

After cleaning the dataframe text cells, we will also tokenize and lemmatize the text. Tokenize entails splitting a string of words into a list of words. For example "cat sat on dog" would be converted to ['cat', 'sat', 'on', 'dog']. Tokenizing splits up each word which we can then use later on to train models easier. Next, we will lemmatize the text. We can choose to lemmatize or stem the words. For this project, I chose to lemmatize the words because it keeps a bit more information than would stemming. An example of lemmatizing would be to take the words 'running', 'horses', and 'adjustable', and we lemmatize them to be 'run', 'horse', 'adjust'. This keeps the words general meaning but allows the model to learn better. Additionally, we will ensure that all words are in lowercase form. These cleaning steps mentioned before and here are important because, for example, the computer could look at two sentences such as "Running big reddish dogs." versus "Run big red dog!" and these would be considered different even though they are quite similar. After cleaning, both sentences would be converted to ['run', 'big', 'red', 'dog'] and therefore these two "articles" probably would be classified together, which is the goal of our model.

```

1 # tokenize text function
2 wordnet_lemmatizer = WordNetLemmatizer()
3 def lemmatizer(text):
4     """
5     A helper function to lemmatize an entire sentence/string
6     """
7     lem = [wordnet_lemmatizer.lemmatize(word.lower()) for word in text]
8     return lem
9
10 def tokenize_lemmatize(dataframe, text_col):
11     """
12     A helper function to tokenize then lemmatize the string.
13     Also, add column which counts the number of words in that string.
14     """
15     dataframe['tokenized'] = dataframe.apply(lambda row: nltk.word_tokenize(row[text_col]), axis=1)
16     dataframe['lemmatized'] = dataframe['tokenized'].apply(lambda string: lemmatizer(string))
17     dataframe['num_words'] = dataframe['lemmatized'].apply(lambda lst: len(lst))
18     return

```

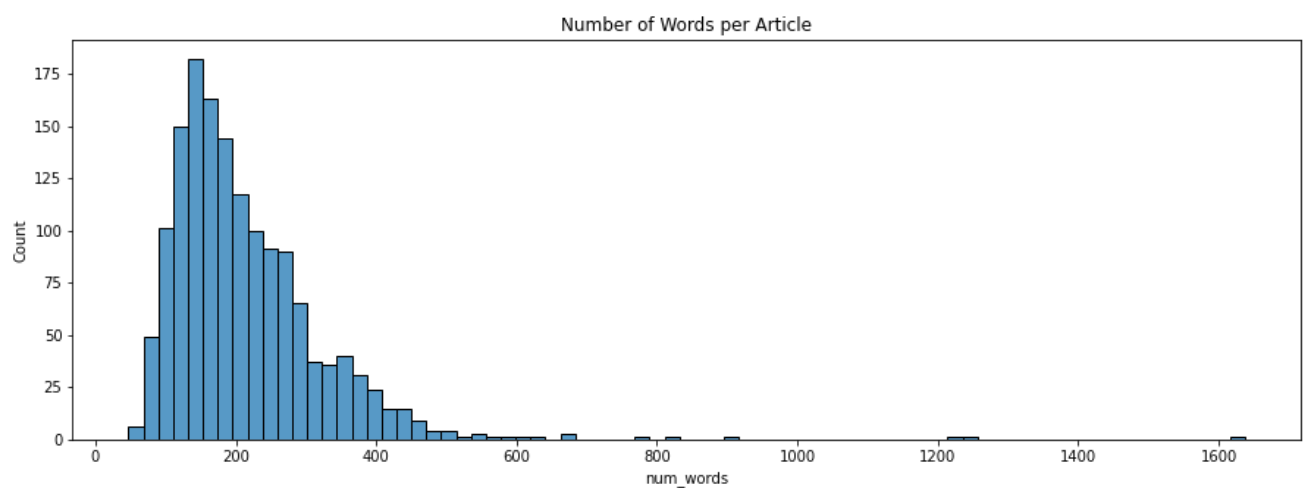
```

1 tokenize_lemmatize(train, 'clean_text')

```

After cleaning, we can see (below) the number of words per article. We see most articles are around 200 words. However, we also see we have some severe outliers that reach up to more than 750 words! We actually will remove these outliers as they might actually impact our model later on, in addition to creating more features (words) to have to calculate within the model.

```
1 # number of tokens (words) per article
2 fig, ax = plt.subplots(figsize=(15, 5))
3 sns.histplot(
4     data = train,
5     x = 'num_words',
6     palette = 'colorblind',
7     ).set(
8         title = 'Number of Words per Article');
```



```
1 #remove outlier articles (longer than 750 words)
2 train = train[train['num_words'] < 750]
3 len(train)

1484
```

Let's also look at the number of words per category (boxplot below). Per category, we also see quite a few outliers. We will leave these this time. We also see that the mean of each category is similar, with tech and politics having more words, and variance, than the rest of the topics.

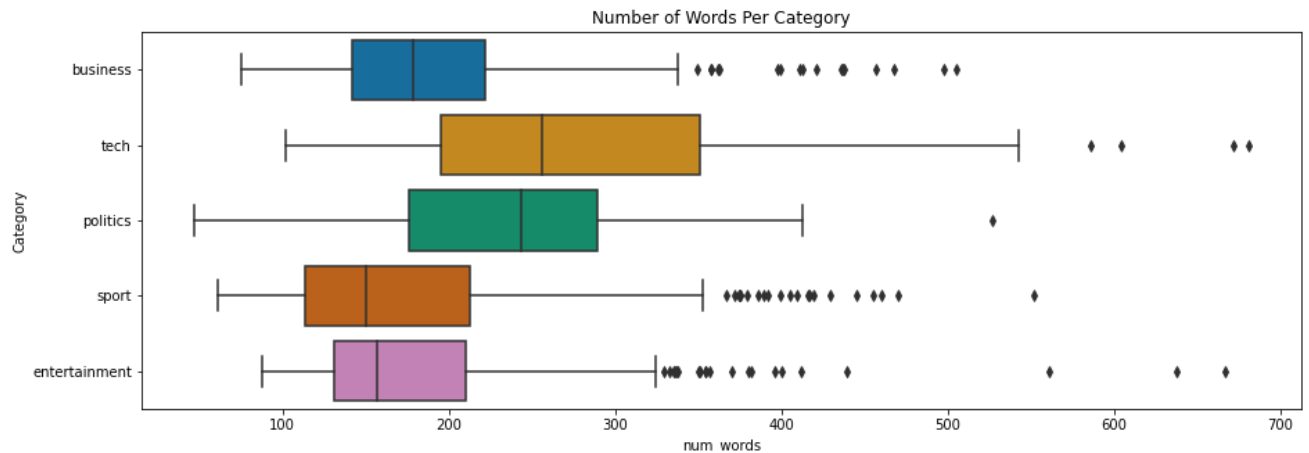
```
1 # words per category
2 fig, ax = plt.subplots(figsize=(15, 5))
3 sns.boxplot(
4     data = train,
5     x = 'num_words',
6     y = 'Category',
```



```

7 palette = 'colorblind'
8 ).set(
9     title = 'Number of Words Per Category';

```



▼ 2. Model Building and Training

```

1 train_df = train.copy()

```

```

1 def predict(w_matrix):
2     sortedW = np.argsort(w_matrix)
3     n_predictions, maxVal = sortedW.shape
4     predictions = [[sortedW[i][maxVal - 1]] for i in range(n_predictions)]
5     topics = np.empty(n_predictions, dtype = np.int64)
6     for i in range(n_predictions):
7         topics[i] = predictions[i][0]
8     return topics

```

```

1 def label_permute(ytdf,yp,n=5):
2     """
3     ytdf: labels dataframe object
4     yp: clustering label prediction output
5     Returns permuted label order and accuracy.
6     Example output: (3, 4, 1, 2, 0), 0.74
7     """
8     perms = list(itertools.permutations([0, 1, 2, 3, 4])) #create permutation list
9     best_labels = []
10    best_acc = 0
11    current = {}
12    labels = ['business', 'tech', 'politics', 'sport', 'entertainment']
13    for perm in perms:

```

```

14     for i in range(n):
15         current[labels[i]] = perm[i]
16         if len(current) == 5:
17             conditions = [
18                 (ytdf['Category'] == current['business']),
19                 (ytdf['Category'] == current['tech']),
20                 (ytdf['Category'] == current['politics']),
21                 (ytdf['Category'] == current['sport']),
22                 (ytdf['Category'] == current['entertainment'])]
23             ytdf['test'] = ytdf['Category'].map(current)
24             current_accuracy = accuracy_score(ytdf['test'], yp)
25             if current_accuracy > best_acc:
26                 best_acc = current_accuracy
27                 best_labels = perm
28                 ytdf['best'] = ytdf['test']
29     return best_labels, best_acc

```

```

1  #create vectorizer
2  tfidvec = TfidfVectorizer(min_df = 2,
3                           max_df = 0.95,
4                           norm = 'l2',
5                           stop_words = 'english')
6  tfidvec_train = tfidvec.fit_transform(train_df['clean_text'])
7
8  #create model
9  nmf_model = NMF(n_components=5,
10                 init='nndsvda',
11                 solver = 'mu',
12                 beta_loss = 'kullback-leibler',
13                 l1_ratio = 0.5,
14                 random_state = 101)
15 nmf_model.fit(tfidvec_train)
16
17 #view results
18 yhat_train = predict(nmf_model.transform(tfidvec_train))
19 label_order, accuracy = label_permute(train_df, yhat_train )
20 print('accuracy=', accuracy)

```

accuracy= 0.9609164420485176

For this model, using matrix factorization, I found the best combination of parameters as above which resulted in the highest accuracy. I also tried other combinations, changing the TfidfVectorizer and/or NMF model parameters. Particularly, I played around with min_df and max_df in the TfidfVectorizer

- using max_df values of 0.85, 0.90, 0.95
- using min_df values of 0, 1, and 2
- using beta_loss of 'frobenius' and 'kullback-leibler'
- using solver of 'mu' and 'cd'

Although all models performed with higher than 85%, the combination above (in the code) was the best model.

```

1 #show best labels for the trained model
2 label_dict = {4:'business', 2:'tech', 1:'politics', 0:'sport', 3:'entertainment'}
3 for i in range(5):
4     print(f'{label_order[i]}: {label_dict[label_order[i]]}')

```

```

4: business
2: tech
1: politics
0: sport
3: entertainment

```

```

1 #first clean testing data as we did with the training data
2 clean_text(test, 'Text')
3 tfidvec_test = tfidvec.transform(test['clean_text'])
4 yhat_test = predict(nmf_model.transform(tfidvec_test))

```

```

1 #create a submission dataframe
2 test_predictions = pd.DataFrame(columns=['ArticleId', 'Category', 'yhat'])
3 test_predictions['ArticleId'] = test['ArticleId']
4 test_predictions['yhat'] = yhat_test
5 test_predictions['Category'] = test_predictions['yhat'].apply(lambda i: label_dict[i])
6
7 #delete columns unneeded for submission
8 test_predictions = test_predictions.drop('yhat', 1)
9 print(test_predictions.head(15))

```

	ArticleId	Category
0	1018	sport
1	1319	tech
2	1138	sport
3	459	business
4	1020	sport
5	51	sport
6	2025	politics
7	1479	politics
8	27	entertainment
9	397	business
10	1644	business
11	263	tech
12	765	politics
13	2134	tech
14	297	entertainment

/opt/conda/lib/python3.7/site-packages/ipykernel_launcher.py:8: FutureWarning: In a future version c

```

1 #save and submit test dataframe
2 try:
3     test_predictions.to_csv('submission.csv', index=False)
4 except:
5     pass
6
7 # #public and private score was 0.96326

```

To get the accuracy on the test set, I went ahead and submitted the results from the unsupervised model. Our model got a test accuracy score of 0.96326, or 96.3%, which is pretty good! Let's see if we can do even better using a supervised model.

▼ 3. Model Comparisons

For this project, we are asked to use unsupervised learning to classify text articles using a matrix factorization model. Traditionally, supervised models would perform better with this type of data if we have pre-labeled text, which we do. Therefore, we will compare the unsupervised learning model above to a supervised model below. Since there was no given direction for which supervised model to use specifically, for fun let's try a KMeans clustering model. For good measure, we will re-import the data again since we'll be working with a new model.

```
1 #import data
2 train = pd.read_csv(train_path)
3 test = pd.read_csv(test_path)

4
5
6
7
8 #clean data
9 clean_text(train, 'Text')
10
11
12
13 #split data into X and y
14 y_train = train['Category'].values
15 X_train = train['clean_text'].values
16
17
18 #create new vectorizer for supervised learning model
19 tfidfvec_supervised = TfidfVectorizer(min_df = 2,
20                                     max_df = 0.95,
21                                     norm = 'l2',
22                                     stop_words = 'english')
23 tfSuper_train = tfidfvec_supervised.fit_transform(X_train)
24
25
26 #create KMeans Model and train
27 kmeans = KMeans(n_clusters = 5,
28                 init = 'k-means++',
29                 algorithm = 'full',
30                 random_state = 101)
31 yhat_train_super = kmeans.fit_predict(tfSuper_train)
32
33
34 #get accuracy
35 y_train_df = pd.DataFrame(y_train, columns=['Category'])
36 label_order, accuracy = label_permute(y_train_df, yhat_train_super)
37 print('accuracy=', accuracy)
38 print(label_order, '\n')
39
40
41 #show label order
42 label_dict = {3:'business', 1:'tech', 4:'politics', 2:'sport', 0:'entertainment'}
```

```

30 for i in range(5):
31     print(f'{label_order[i]}: {label_dict[label_order[i]]}')

accuracy= 0.9369127516778524
(3, 1, 0, 2, 4)

3: business
1: tech
0: entertainment
2: sport
4: politics

```

Now we will test the model on the test set. Again, we need to repeat the same steps to clean the data like we did in the training set.

```

1 #clean data
2 clean_text(test, 'Text')
3
4 #split data
5 X_test = test['clean_text'].values
6
7 #create vectorizer (do not fit it!)
8 tfidfvec_supervised.transform(X_test)
9 yhat_test = kmeans.predict(tfidfvec_supervised.transform(X_test))
10
11 #create a submission dataframe
12 test_predictions = pd.DataFrame(columns=['ArticleId', 'Category', 'yhat'])
13 test_predictions['ArticleId'] = test['ArticleId']
14 test_predictions['yhat'] = yhat_test
15 test_predictions['Category'] = test_predictions['yhat'].apply(lambda i: label_dict[i])
16
17 #delete columns unneeded for submission
18 test_predictions = test_predictions.drop('yhat', 1)
19 print(test_predictions.head(2))

   ArticleId  Category
0        1018    sport
1        1319     tech
/opt/conda/lib/python3.7/site-packages/ipykernel_launcher.py:18: FutureWarning: In a future version

1 #save and submit test dataframe
2 # try:
3 #     test_predictions.to_csv('submission.csv', index=False)
4 # except:
5 #     pass
6
7 # #public and private score was 0.62993

```

It looks like we got a training accuracy of 93.69% but only a testing accuracy of 62.99%. Wow! This suggests that our model is probably overfitting to the training data and does a poor job predicting on new data. If we wanted to train a more powerful supervised learning model, we could use techniques such as ensemble methods and/or cross-validation (Kfold), or different

models such as decision tree, random forest, SVM, etc. However, since this class focuses on unsupervised learning techniques, and that is our main focus in this project, we will simply compare our KMeans model to the matrix factorization model.

▼ 4. Conclusions

To summarize this project, we first cleaned the training data in common NLP preprocessing ways and explored the data. Then we created a matrix factorization model and got a testing accuracy of 96.3%. We got this score by find-tuning some parameters and using the training accuracy as a guide. The unsupervised model did quite well coompared to the supervised learning model. I suspect the supervised model overfit to the training data, which our unsupervised model did not. We made sure to preprocess the data in the same way for all training and testing runs.

Future Project Enhancement Options: Here are a few things that could be performed in future NLP projects which may impact the results in this study by improving or decreasing the performance. It is worth trying out. This project removed captialized letters, all numbers, and punctuation. In addition, we did not deal with misspelled or 'non-real' words. These things can impact the outcomes of a model and are worth modifying and trying out in different ways to see if better results can be obtained.

End 🇵🇸
