

```
In [2]: import pandas
import numpy as np
import matplotlib.pyplot as plt

#seaborn is an open-source Python library built on top of matplotlib. It is used for data visualization and exploratory data
#analysis. Seaborn works easily with dataframes and the Pandas library. The graphs created can also be customized easily.
import seaborn as sns
from sklearn.feature_extraction.text import CountVecorizer
#used for bag of words and it extract feature from text document
from sklearn.feature_extraction.text import TfidfTransformer
#tfidf = term frequency is used for
from sklearn import feature_extraction, linear_model, model_selection, preprocessing
from sklearn.metrics import accuracy_score
from sklearn.model_selection import train_test_split
from sklearn.pipeline import Pipeline
#cross validation when setting different parameters

In [3]: fake = pandas.read_csv(r"C:\Users\lisha\Downloads\Fake.csv")
true = pandas.read_csv(r"C:\Users\lisha\Downloads\True.csv")

In [4]: fake.shape

Out[4]: (23481, 4)

In [5]: true.shape

Out[5]: (21437, 4)

In [6]: fake.head()

Out[6]:
```

	title	text	subject	date
0	Donald Trump Sends Out Embarrassing New Year...	Donald Trump just couldn't wish all Americans ...	News	December 31, 2017
1	Drunk Bragging Trump Staffer Started Russian...	House Intelligence Committee Chairman Devin Nu...	News	December 31, 2017
2	Sheriff David Clarke Becomes An Internet Joke...	On Friday, it was revealed that former Milwaukee...	News	December 30, 2017
3	Trump Is So Obsessed He Even Has Obama's Name...	On Christmas day, Donald Trump announced that...	News	December 29, 2017
4	Pope Francis Just Called Out Donald Trump Dur...	Pope Francis used his annual Christmas Day mes...	News	December 25, 2017

```
In [7]: true.head()

Out[7]:
```

	title	text	subject	date
0	As U.S. budget fight looms, Republicans flip t...	WASHINGTON (Reuters) - The head of a conservat...	politicsNews	December 31, 2017
1	U.S. military to accept transgender recruits o...	WASHINGTON (Reuters) - Transgender people wil...	politicsNews	December 29, 2017
2	Senior U.S. Republican senator: Let Mr. Muel...	WASHINGTON (Reuters) - The special counsel inv...	politicsNews	December 31, 2017
3	FBI Russia probe helped by Australian diplom...	WASHINGTON (Reuters) - Trump campaign adviser ...	politicsNews	December 30, 2017
4	Trump wants Postal Service to charge much mor...	SEATTLE/WASHINGTON (Reuters) - President Donal...	politicsNews	December 29, 2017

```
In [8]: # add flag to track fake and real
fake['target'] = 'fake'
true['target'] = 'true'

In [9]: fake.head()

Out[9]:
```

	title	text	subject	date	target
0	Donald Trump Sends Out Embarrassing New Year...	Donald Trump just couldn't wish all Americans ...	News	December 31, 2017	fake
1	Drunk Bragging Trump Staffer Started Russian...	House Intelligence Committee Chairman Devin Nu...	News	December 31, 2017	fake
2	Sheriff David Clarke Becomes An Internet Joke...	On Friday, it was revealed that former Milwaukee...	News	December 30, 2017	fake
3	Trump Is So Obsessed He Even Has Obama's Name...	On Christmas day, Donald Trump announced that...	News	December 29, 2017	fake
4	Pope Francis Just Called Out Donald Trump Dur...	Pope Francis used his annual Christmas Day mes...	News	December 25, 2017	fake

```
In [10]: true.head()

Out[10]:
```

	title	text	subject	date	target
0	As U.S. budget fight looms, Republicans flip t...	WASHINGTON (Reuters) - The head of a conservat...	politicsNews	December 31, 2017	true
1	U.S. military to accept transgender recruits o...	WASHINGTON (Reuters) - Transgender people wil...	politicsNews	December 29, 2017	true
2	Senior U.S. Republican senator: Let Mr. Muel...	WASHINGTON (Reuters) - The special counsel inv...	politicsNews	December 31, 2017	true
3	FBI Russia probe helped by Australian diplom...	WASHINGTON (Reuters) - Trump campaign adviser ...	politicsNews	December 30, 2017	true
4	Trump wants Postal Service to charge much mor...	SEATTLE/WASHINGTON (Reuters) - President Donal...	politicsNews	December 29, 2017	true

```
In [11]: # combine datasets
data = pandas.concat([fake, true]).reset_index(drop = True)
data.shape

Out[11]: (44898, 5)

In [12]: data.head()

Out[12]:
```

	title	text	subject	date	target
0	Donald Trump Sends Out Embarrassing New Year...	Donald Trump just couldn't wish all Americans ...	News	December 31, 2017	fake
1	Drunk Bragging Trump Staffer Started Russian...	House Intelligence Committee Chairman Devin Nu...	News	December 31, 2017	fake
2	Sheriff David Clarke Becomes An Internet Joke...	On Friday, it was revealed that former Milwaukee...	News	December 30, 2017	fake
3	Trump Is So Obsessed He Even Has Obama's Name...	On Christmas day, Donald Trump announced that...	News	December 29, 2017	fake
4	Pope Francis Just Called Out Donald Trump Dur...	Pope Francis used his annual Christmas Day mes...	News	December 25, 2017	fake

```
In [13]: data.tail()

Out[13]:
```

	title	text	subject	date	target
44893	Fully committed NATO backs new U.S. approach...	BRUSSELS (Reuters) - NATO allies on Tuesday we...	worldnews	August 22, 2017	true
44894	LexisNexis withdrew two products from Chinese...	LONDON (Reuters) - LexisNexis, a provider of l...	worldnews	August 22, 2017	true
44895	Minsk cultural hub becomes haven from authori...	MINSK (Reuters) - In the shadow of disused Sov...	worldnews	August 22, 2017	true
44896	Vatican upbeat on possibility of Pope Francis...	MOSCOW (Reuters) - Vatican Secretary of State...	worldnews	August 22, 2017	true
44897	Indonesia to buy \$1.14 billion worth of Russia...	JAKARTA (Reuters) - Indonesia will buy 11 Sukh...	worldnews	August 22, 2017	true

```
In [14]: # shuffle the data
from sklearn.utils import shuffle
data = shuffle(data)
data = data.reset_index(drop=True)

In [15]: data.head()

Out[15]:
```

	title	text	subject	date	target
0	Illinois governor's office warns of crippling ...	CHICAGO (Reuters) - Potential action this week...	politicsNews	August 23, 2016	true
1	U.S. ethics office tweets sarcasm at Trump on ...	WASHINGTON (Reuters) - President-elect Donald ...	politicsNews	November 30, 2016	true
2	UK PM May is focused on tackling extension, sp...	LONDON (Reuters) - British Prime Minister Ther...	worldnews	November 30, 2017	true
3	Scaramucci's SkyBridge sells itself, investmen...	BOSTON/BEIJING (Reuters) - The hedge fund inve...	politicsNews	January 17, 2017	true
4	Republican Rep. Brooks: 30-40 Republican 'no' ...	WASHINGTON (Reuters) - Republican U.S. Represe...	politicsNews	March 23, 2017	true

```
In [16]: data.info()

Out[16]:
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 44898 entries, 0 to 44897
Data columns (total 5 columns):
#    Column  Non-Null Count  Dtype
#  --  --
0     title    44898 non-null     object
1     text     44898 non-null     object
2     subject  44898 non-null     object
3     date     44898 non-null     object
4     target   44898 non-null     object
dtypes: object(5)
memory usage: 5.1+ MB
```

```
In [17]: # removing date
data.drop(['date'],axis=1,inplace=True)
data.head()

Out[17]:
```

	title	text	subject	target
0	Illinois governor's office warns of crippling ...	CHICAGO (Reuters) - Potential action this week...	politicsNews	true
1	U.S. ethics office tweets sarcasm at Trump on ...	WASHINGTON (Reuters) - President-elect Donald ...	politicsNews	true
2	UK PM May is focused on tackling extension, sp...	LONDON (Reuters) - British Prime Minister Ther...	worldnews	true
3	Scaramucci's SkyBridge sells itself, investmen...	BOSTON/BEIJING (Reuters) - The hedge fund inve...	politicsNews	true
4	Republican Rep. Brooks: 30-40 Republican 'no' ...	WASHINGTON (Reuters) - Republican U.S. Represe...	politicsNews	true

```
In [18]: data.drop(['title'],axis=1,inplace=True)
data.head()

Out[18]:
```

	text	subject	target
0	CHICAGO (Reuters) - Potential action this week...	politicsNews	true
1	WASHINGTON (Reuters) - President-elect Donald...	politicsNews	true
2	LONDON (Reuters) - British Prime Minister Ther...	worldnews	true
3	BOSTON/BEIJING (Reuters) - The hedge fund inve...	politicsNews	true
4	WASHINGTON (Reuters) - Republican U.S. Represe...	politicsNews	true

```
In [19]: #convert lower case
data['text'] = data['text'].apply(lambda x: x.lower())
data.head()

Out[19]:
```

	text	subject	target
0	chicago (reuters) - potential action this week...	politicsNews	true
1	washington (reuters) - president-elect donald...	politicsNews	true
2	london (reuters) - british prime minister ther...	worldnews	true
3	bostonbeijing (reuters) - the hedge fund inve...	politicsNews	true
4	washington (reuters) - republican u.s. represe...	politicsNews	true

```
In [20]: # remove punctuation
import string
def punctuation_remove(text):
    all_list = [char for char in text if char not in string.punctuation]
    clean_str = ''.join(all_list)
    return clean_str
data['text'] = data['text'].apply(punctuation_remove)

In [21]: data.head()

Out[21]:
```

	text	subject	target
0	chicago reuters potential action this week by...	politicsNews	true
1	washington reuters presidentelect donald trum...	politicsNews	true
2	london reuters british prime minister theresa...	worldnews	true
3	bostonbeijing reuters hedge fund investme...	politicsNews	true
4	washington reuters republican us representat...	politicsNews	true

```
In [22]: # Removing stopwords
#nltk is a standard python library with prebuilt functions and utilities for the ease of use and implementation.
#it is one of the most used libraries for natural language processing and computational linguistics.
import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords
stop = stopwords.words('english')
data['text'] = data['text'].apply(lambda x: ' '.join([word for word in x.split() if word not in (stop)]))

[nltk_data] downloading package stopwords to
[nltk_data]   c:\users\lisha\appdata\local\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
```

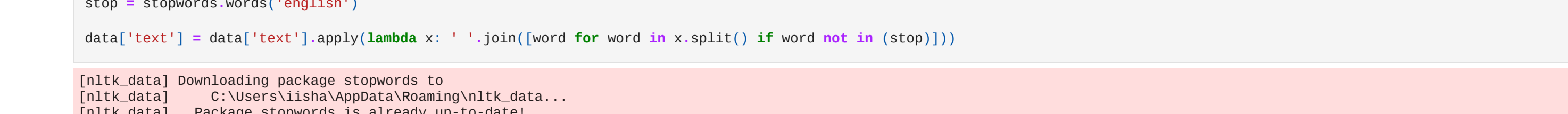
```
In [23]: data.head()

Out[23]:
```

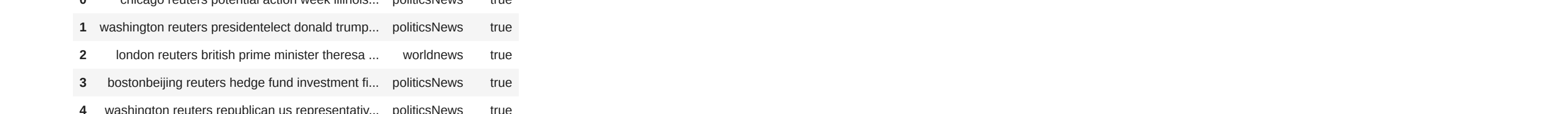
	text	subject	target
0	chicago reuters potential action week illino...	politicsNews	true
1	washington reuters presidentelect donald trum...	politicsNews	true
2	london reuters british prime minister theresa...	worldnews	true
3	bostonbeijing reuters hedge fund investme...	politicsNews	true
4	washington reuters republican us representat...	politicsNews	true

Basic data exploration

```
In [24]: # How many articles per subject?
print(data.groupby(['subject'])['text'].count())
data.groupby(['subject'])['text'].count().plot(kind="bar")
plt.show()
```



```
In [25]: # How many fake and real articles?
print(data.groupby(['target'])['text'].count())
data.groupby(['target'])['text'].count().plot(kind="bar")
plt.show()
```



```
In [26]: !pip install wordcloud
Requirement already satisfied: wordcloud in c:\users\lisha\anaconda3\lib\site-packages (1.8.0)
Requirement already satisfied: matplotlib in c:\users\lisha\anaconda3\lib\site-packages (from wordcloud) (3.3.4)
Requirement already satisfied: pillow in c:\users\lisha\anaconda3\lib\site-packages (from wordcloud) (8.2.0)
Requirement already satisfied: kiwisolver<=1.0.1 in c:\users\lisha\anaconda3\lib\site-packages (from matplotlib->wordcloud) (1.3.1)
Requirement already satisfied: cycler<=0.10 in c:\users\lisha\anaconda3\lib\site-packages (from matplotlib->wordcloud) (0.10.0)
Requirement already satisfied: pyparsing<=2.0.4, >=2.1.2, <=3.6, >=2.0.3 in c:\users\lisha\anaconda3\lib\site-packages (from matplotlib->wordcloud) (2.4.7)
Requirement already satisfied: python-dateutil<=2.1 in c:\users\lisha\anaconda3\lib\site-packages (from matplotlib->wordcloud) (2.8.1)
Requirement already satisfied: six in c:\users\lisha\anaconda3\lib\site-packages (from cycler<=0.10->matplotlib->wordcloud) (1.12.0)
```

```
In [27]: # word cloud for fake news
from wordcloud import WordCloud
fake_data = data[data['target'] == "fake"]
all_words = ' '.join([text for text in fake_data.text])
wordcloud = WordCloud(width=800, height=500,
                      max_font_size=110,
                      collocations=False).generate(all_words)
plt.figure(figsize=(10,7))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```

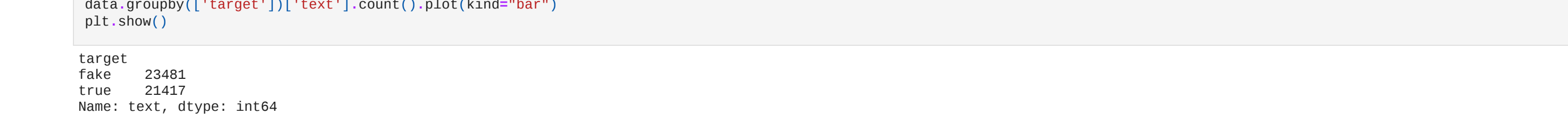


```
In [28]: # word cloud for real news
from wordcloud import WordCloud
real_data = data[data['target'] == "true"]
all_words = ' '.join([text for text in real_data.text])
wordcloud = WordCloud(width=800, height=500,
                      max_font_size=110,
                      collocations=False).generate(all_words)
plt.figure(figsize=(10,7))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```



```
In [29]: # Most frequent words counter
from nltk import tokenize
token_space = tokenize.WhitespaceTokenizer()
def counter(text, column_text, quantity):
    all_words = ' '.join([text for text in text[column_text]])
    token_phrase = token_space.tokenize(all_words)
    frequency = nltk.FreqDist(token_phrase)
    df_frequency = pandas.DataFrame({"word": list(frequency.keys()),
                                     "frequency": list(frequency.values())})
    df_frequency = df_frequency.nlargest(columns = "frequency", n = quantity)
    plt.figure(figsize=(12,8))
    ax = sns.barplot(data=df_frequency, x = "word", y = "frequency", color = "blue")
    ax.set(ylabel = "count")
    plt.xticks(rotation='vertical')
    plt.show()
```

```
In [30]: # Most frequent words in fake news
counter(data[data['target'] == "fake"], "text", 20)
```



```
In [31]: # Most frequent words in real news
counter(data[data['target'] == "true"], "text", 20)
```



```
In [32]: # Function to plot the confusion matrix
from sklearn import metrics
import itertools
def plot_confusion_matrix(cm, classes,
                          normalize=False,
                          title='Confusion matrix',
                          cmap=plt.cm.Blues):
    plt.imshow(cm, interpolation='nearest', cmap=cmap)
    plt.title(title)
    plt.colorbar()
    tick_marks = np.arange(len(classes))
    plt.xticks(tick_marks, classes, rotation=45)
    plt.yticks(tick_marks, classes)
    if normalize:
        cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
        print("Normalized confusion matrix")
    else:
        print('Confusion matrix, without normalization')
    thresh = cm.max() / 2.
    for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
        plt.text(j, i, cm[i, j],
                 horizontalalignment="center",
                 color="white" if cm[i, j] > thresh else "black")
    plt.tight_layout()
    plt.ylabel('True label')
    plt.xlabel('Predicted label')
```

Split data

```
In [33]: # Split the data
X_train,X_test,y_train,y_test = train_test_split(data['text'], data.target, test_size=0.2, random_state=42)

In [34]: X_train.head()

Out[34]:
```

```
36335    washington reuters us state parties approve...
32384    berlin reuters german departement cited progr...
24419    astana reuters russian foreign minister serget...
24740    republicans even moderately literate would kno...
27839    washington reuters company promoting plan unit...
Name: text, dtype: object
```

```
In [35]: y_train.head()

Out[35]:
```

```
36335    true
32384    true
24419    true
24740    fake
27839    true
Name: target, dtype: object
```

Decision Tree Classifier

```
In [36]: from sklearn.tree import DecisionTreeClassifier
# Vectorizing and applying TF-IDF
pipe = Pipeline([('vect', CountVecorizer()),
                  ('tridf', TfidfTransformer()),
                  ('model', DecisionTreeClassifier(criterion='entropy',
                                                    max_depth=28,
                                                    splitters='best',
                                                    random_state=42))])
# Fitting the model
model = pipe.fit(X_train, y_train)
# Accuracy
prediction = model.predict(X_test)
print("Accuracy: %f" % metrics.f1_score(y_test, prediction)*100,2)))
accuracy: 99.58%
```

```
In [ ]:
```