

```
[15]: 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 [16]: fake = pandas.read_csv(r"C:\Users\lisha\Downloads\Fake.csv")
true = pandas.read_csv(r"C:\Users\lisha\Downloads\True.csv")
```

```
In [17]: fake.shape
```

```
Out[17]: (23481, 4)
```

```
In [18]: true.shape
```

```
Out[18]: (21417, 4)
```

```
In [19]: fake.head()
```

	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 Milwauk...	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 [20]: true.head()
```

	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. Muek...	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 [21]: # add flag to track fake and real
fake['target'] = 'f'
true['target'] = 't'
```

```
In [22]: fake.head()
```

	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 Milwauk...	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 [23]: true.head()
```

	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. Muek...	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 [24]: # combine datasets
data = pandas.concat([fake, true]).reset_index(drop = True)
```

```
Out[24]: data.shape
```

```
In [25]: data.head()
```

	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 Milwauk...	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 [26]: data.tail()
```

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

```
In [28]: data.head()
```

	title	text	subject	date	target
0	IRRATIONAL GEORGETOWN PROFESSOR Has Month-Long...	A Georgetown University associate professor ha...	left-news	Dec 20, 2016	fake
1	U.S. lawmakers want documents on Russia electi...	WASHINGTON (Reuters) - The House of Representa...	politicsNews	January 25, 2017	true
2	NSA chief, intelligence director warn comment...	WASHINGTON (Reuters) - Two top U.S. intelligen...	politicsNews	June 7, 2017	true
3	Obama's Party With #BlackLivesMatter Organizer...	The Obama s are classless we know that. The Ob...	left-news	Dec 5, 2015	fake
4	OBAMA'S EPA GESTAPO TO SKIP Hearing On CO Mine...	Priorities priorities Arizona Sen. John McCain...	left-news	Apr 7, 2016	fake

```
In [29]: data.info()
```

```
<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 [30]: # removing date
data.drop(['date'],axis=1,inplace=True)
data.head()
```

	title	text	subject	target
0	IRRATIONAL GEORGETOWN PROFESSOR Has Month-Long...	A Georgetown University associate professor ha...	left-news	fake
1	U.S. lawmakers want documents on Russia electi...	WASHINGTON (Reuters) - The House of Representa...	politicsNews	true
2	NSA chief, intelligence director warn comment...	WASHINGTON (Reuters) - Two top U.S. intelligen...	politicsNews	true
3	Obama's Party With #BlackLivesMatter Organizer...	The Obama s are classless we know that. The Ob...	left-news	fake
4	OBAMA'S EPA GESTAPO TO SKIP Hearing On CO Mine...	Priorities priorities Arizona Sen. John McCain...	left-news	fake

```
In [31]: data.drop(['title'],axis=1,inplace=True)
data.head()
```

	text	subject	target
0	A Georgetown University associate professor ha...	left-news	fake
1	WASHINGTON (Reuters) - The House of Representa...	politicsNews	true
2	WASHINGTON (Reuters) - Two top U.S. intelligen...	politicsNews	true
3	The Obama s are classless we know that. The Ob...	left-news	fake
4	Priorities priorities Arizona Sen. John McCain...	left-news	fake

```
In [32]: #convert lower case
data['text'] = data['text'].apply(lambda x: x.lower())
data.head()
```

	text	subject	target
0	a georgetown university associate professor ha...	left-news	fake
1	washington (reuters) - the house of representa...	politicsNews	true
2	washington (reuters) - two top u.s. intelligen...	politicsNews	true
3	the obama s are classless we know that. the ob...	left-news	fake
4	priorities priorities arizona sen. john mccain...	left-news	fake

```
In [33]: # remove punctuation
import string
def punctuation_removal(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_removal)
```

```
In [34]: data.head()
```

	text	subject	target
0	a georgetown university associate professor ha...	left-news	fake
1	washington reuters the house of representati...	politicsNews	true
2	washington reuters two top us intelligence of...	politicsNews	true
3	the obama s are classless we know that the oba...	left-news	fake
4	priorities priorities arizona sen john mccain ...	left-news	fake

```
In [35]: # Removing 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)]))
```

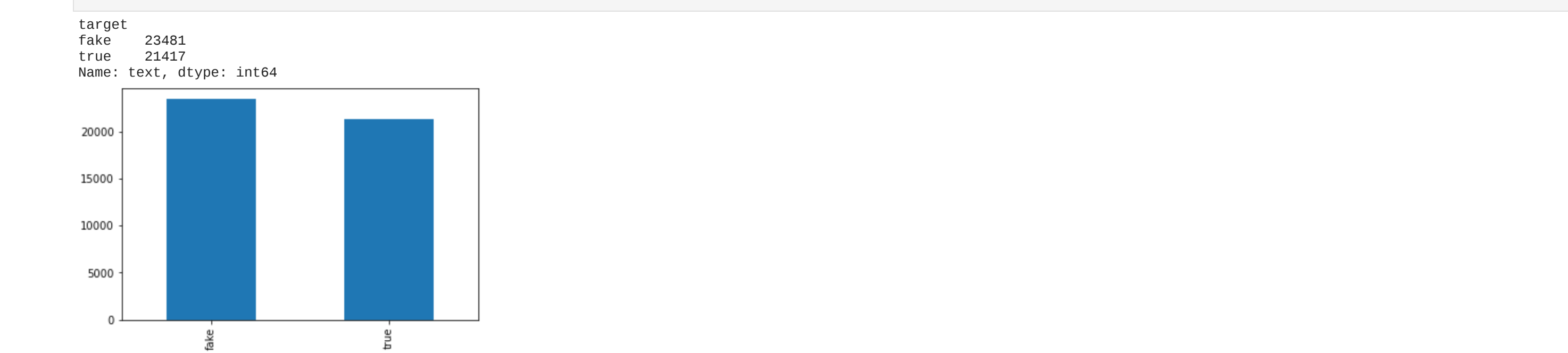
```
[nltk.data] downloading package stopwords to
[nltk.data]   c:\users\lisha\appdata\local\nltk_data...
[nltk.data]   Package stopwords is already up-to-date!
```

```
In [36]: data.head()
```

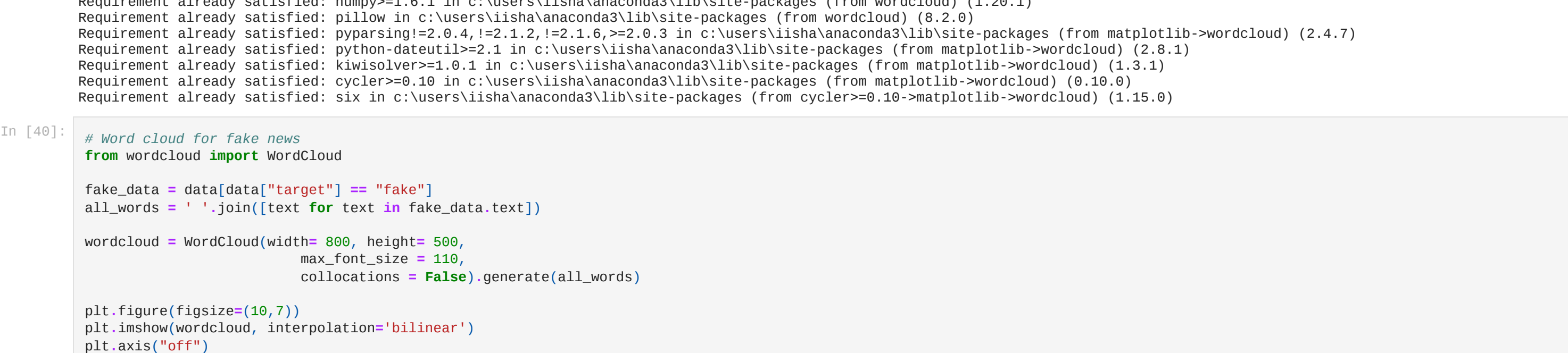
	text	subject	target
0	georgetown university associate professor mont...	left-news	fake
1	washington reuters house representatives intel...	politicsNews	true
2	washington reuters two top us intelligence off...	politicsNews	true
3	obama classless know obama guest list holiday...	left-news	fake
4	priorities priorities arizona sen john mccain ...	left-news	fake

## Basic data exploration

```
In [37]: # How many articles per subject?
print(data.groupby(['subject'])['text'].count())
data.groupby(['subject'])['text'].count().plot(kind="bar")
plt.show()
```

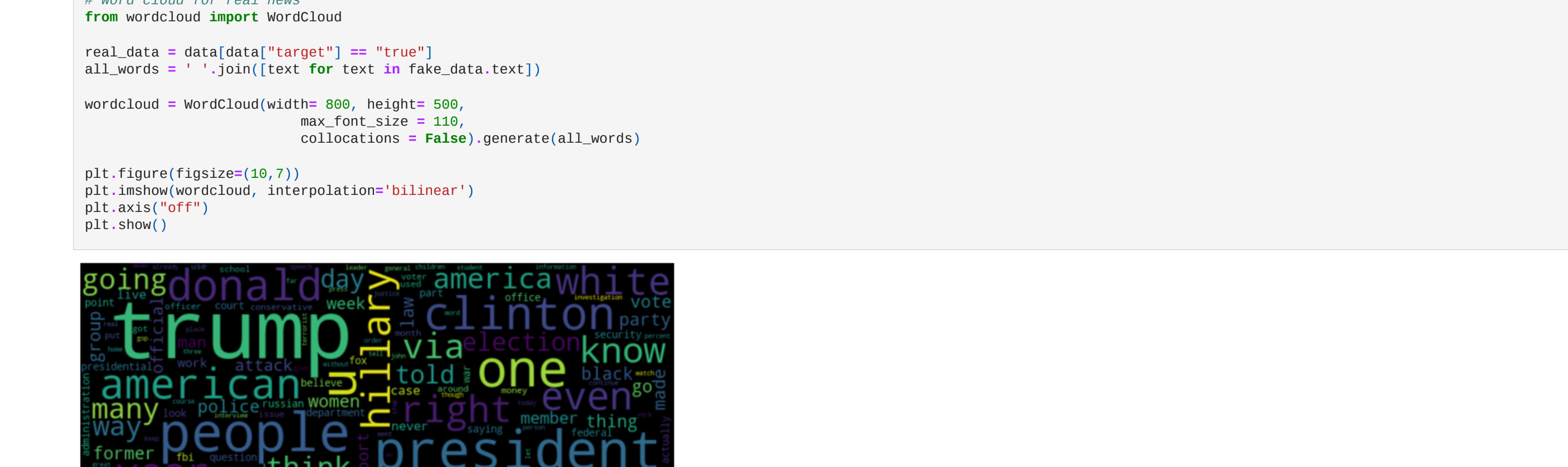


```
In [38]: # How many fake and real articles?
print(data.groupby(['target'])['text'].count())
data.groupby(['target'])['text'].count().plot(kind="bar")
plt.show()
```



```
In [39]: !pip install wordcloud
Requirement already satisfied: wordcloud in c:\users\lisha\anaconda3\lib\site-packages (1.8.1)
Requirement already satisfied: matplotlib in c:\users\lisha\anaconda3\lib\site-packages (from wordcloud) (3.3.4)
Requirement already satisfied: numpy>=1.6.3 in c:\users\lisha\anaconda3\lib\site-packages (from wordcloud) (1.20.1)
Requirement already satisfied: pillow in c:\users\lisha\anaconda3\lib\site-packages (from wordcloud) (8.2.0)
Requirement already satisfied: pypparsing>=2.0.4, !=2.1.2, !=2.1.6, >=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: Keras>=0.1.0 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: six in c:\users\lisha\anaconda3\lib\site-packages (from cycler>=0.10->matplotlib-wordcloud) (1.15.0)
```

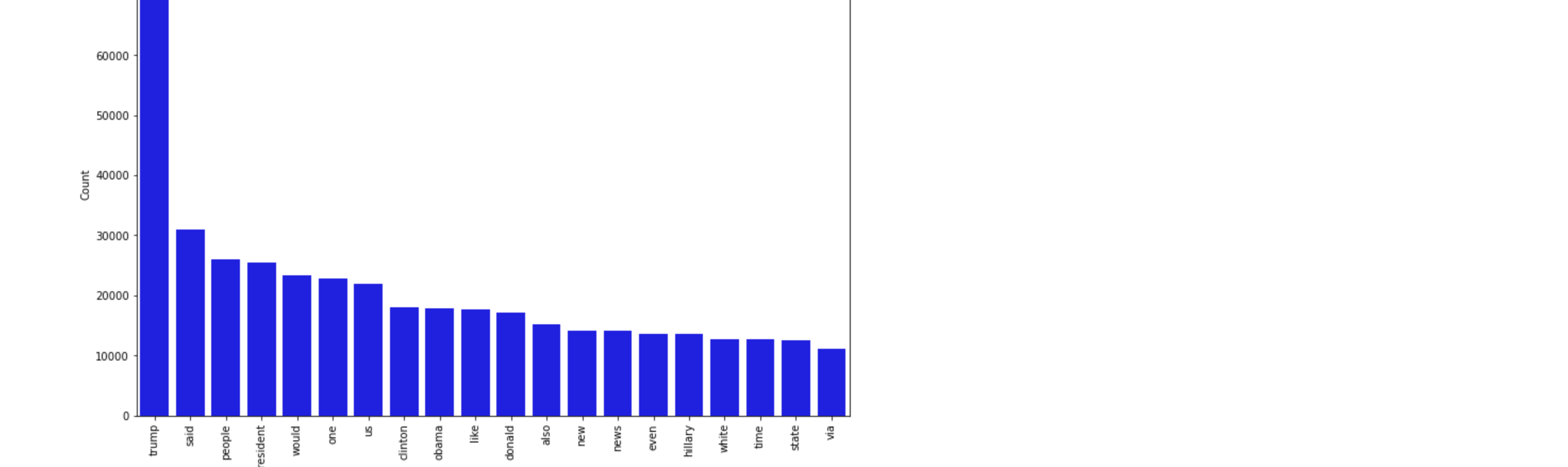
```
In [40]: # word cloud for fake news
from wordcloud import WordCloud
fake_data = data[data['target'] == "fake"]
all_words = ' '.join([text for text in fake_data.text])
wordcloud = WordCloud(width=800, height=500,
                      max_font_size=110,
                      collocations=False).generate(all_words)
plt.figure(figsize=(10,7))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis('off')
plt.show()
```



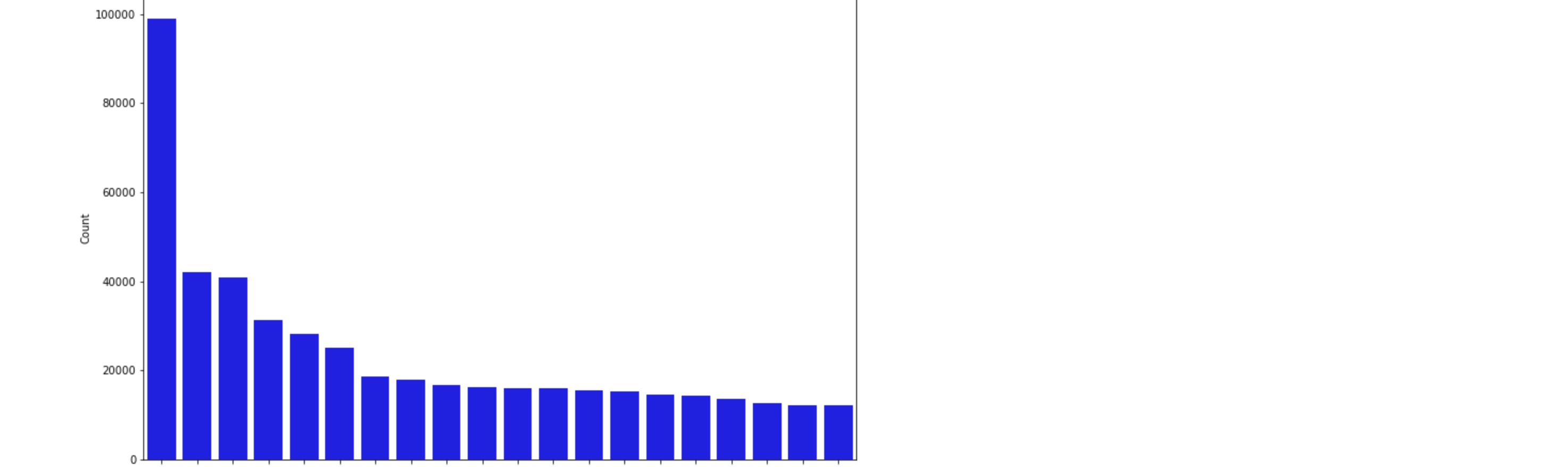
```
In [41]: # word cloud for real news
from wordcloud import WordCloud
real_data = data[data['target'] == "true"]
all_words = ' '.join([text for text in real_data.text])
wordcloud = WordCloud(width=800, height=500,
                      max_font_size=110,
                      collocations=False).generate(all_words)
plt.figure(figsize=(10,7))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis('off')
plt.show()
```



```
In [42]: # 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=(10,8))
    ax = sns.barplot(x=df_frequency, x = "word", y = "frequency", color = "blue")
    ax.set(ylabel = "count")
    plt.xticks(rotation='vertical')
    plt.show()
```



```
In [44]: # Most frequent words in real news
counter(data[data['target'] == "true"], "text", 20)
```



```
In [45]: # 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 [46]: # 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 [47]: X_train.head()
```

```
Out[47]: 36335    21st century wire says censorship running ramp...
32384    gothenburg sweden reuters deadlock divorce tal...
24419    washington reuters republican us senators john...
24740    new york reuters us republican presidential ca...
27839    obama reuters indonesia foreign minister fly b...
```

```
In [48]: y_train.head()
```

```
Out[48]: 36335    fake
32384    true
24419    true
24740    true
27839    true
Name: target, dtype: object
```

## Decision Tree Classifier

```
In [49]: from sklearn.tree import DecisionTreeClassifier
```

```
# Vectorizing and applying TF-IDF
pipe = Pipeline([('vect', CountVecorizer()),
                  ('tridf', TfidfTransformer()),
                  ('model', DecisionTreeClassifier(criterion='entropy',
                                                  max_depth=20,
                                                  splitter='best',
                                                  random_state=42))])
```

```
# Fitting the model
model = pipe.fit(X_train, y_train)
# Accuracy
prediction = model.predict(X_test)
print("accuracy: {:.1%}".format(round(accuracy_score(y_test, prediction)*100,2)))
```

```
In [ ]:
```