



```
from google.colab import files
uploaded = files.upload()

Choose Files News.csv
• News.csv(text/csv) - 116525400 bytes, last modified: 2/26/2024 - 100% done
Saving News.csv to News.csv
```

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

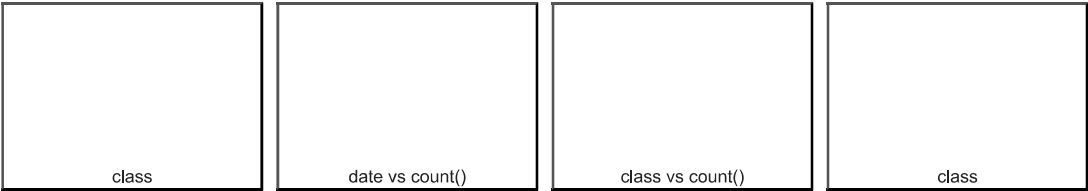
```
data = pd.read_csv('News.csv',index_col=0)
data.head()
```

| | title | text | subject | date | class | |
|---|--|---|---------|-------------------|-------|---|
| 0 | Donald Trump Sends Out Embarrassing New Year'... | Donald Trump just couldn't wish all Americans ... | News | December 31, 2017 | 0 |  |
| 1 | Drunk Bragging Trump Staffer Started Russian ... | House Intelligence Committee Chairman Devin Nu... | News | December 31, 2017 | 0 |  |
| 2 | Sheriff David Clarke Becomes An Internet Joke... | On Friday, it was revealed that former Milwauk... | News | December 30, 2017 | 0 | |

Next steps:

Generate code with data

 View recommended plots



```
data.shape

(44919, 5)

data = data.drop(["title", "subject","date"], axis = 1)

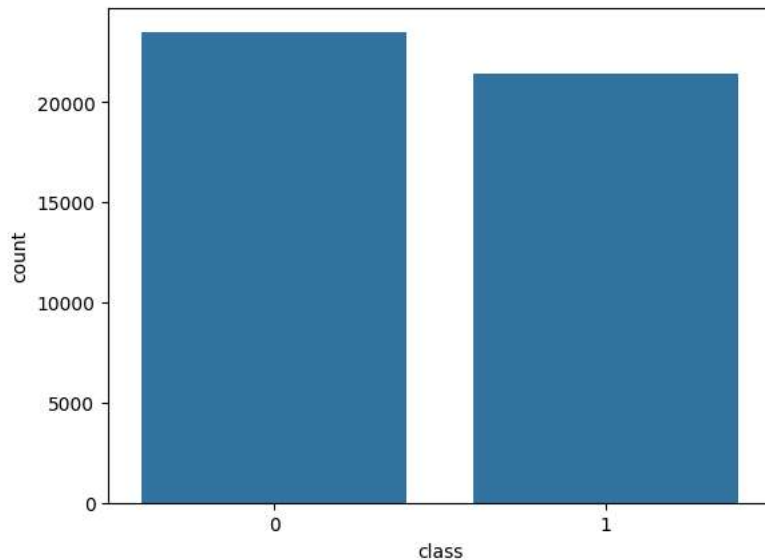
data.isnull().sum()

text      0
class     0
dtype: int64

data = data.sample(frac=1)
data.reset_index(inplace=True)
data.drop(["index"], axis=1, inplace=True)

sns.countplot(data=data,
              x='class',
              order=data['class'].value_counts().index)
```

<Axes: xlabel='class', ylabel='count'>



```

from tqdm import tqdm
import re
import nltk
nltk.download('punkt')
nltk.download('stopwords')
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.stem.porter import PorterStemmer
from wordcloud import WordCloud

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Unzipping tokenizers/punkt.zip.
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.

def preprocess_text(text_data):
    preprocessed_text = []

    for sentence in tqdm(text_data):
        sentence = re.sub(r'[^\w\s]', '', sentence)
        preprocessed_text.append(' '.join(token.lower()
                                           for token in str(sentence).split()
                                           if token not in stopwords.words('english')))

    return preprocessed_text

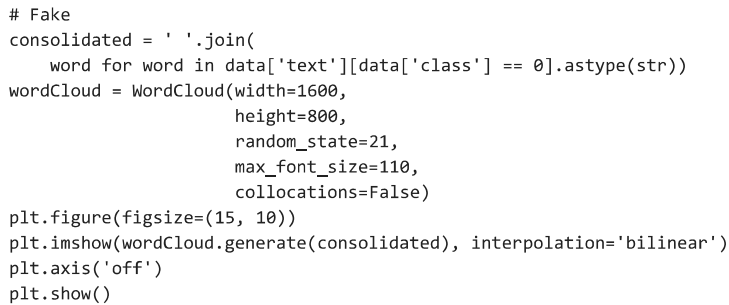
preprocessed_review = preprocess_text(data['text'].values)
data['text'] = preprocessed_review

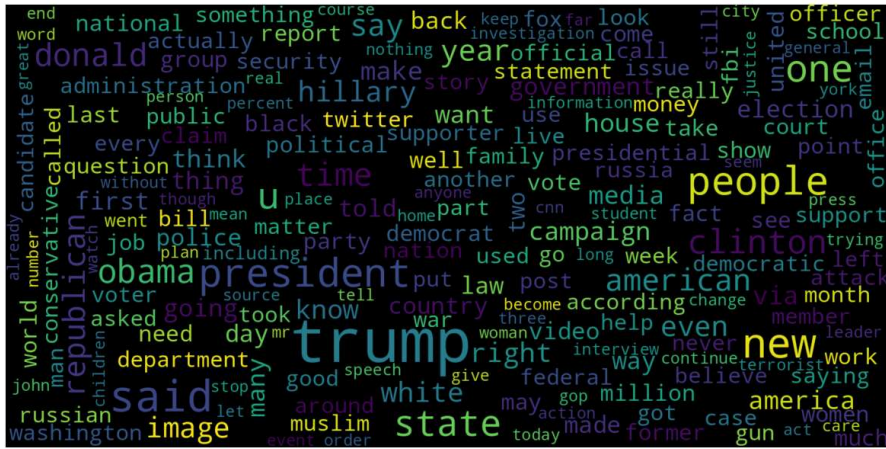
100%|██████████| 44919/44919 [38:49<00:00, 19.28it/s]

consolidated = ' '.join(
    word for word in data['text'][data['class'] == 1].astype(str))
wordCloud = WordCloud(width=1600,
                      height=800,
                      random_state=21,
                      max_font_size=110,
                      collocations=False)

plt.figure(figsize=(15, 10))
plt.imshow(wordCloud.generate(consolidated), interpolation='bilinear')
plt.axis('off')
plt.show()

```





```
from sklearn.feature_extraction.text import CountVectorizer

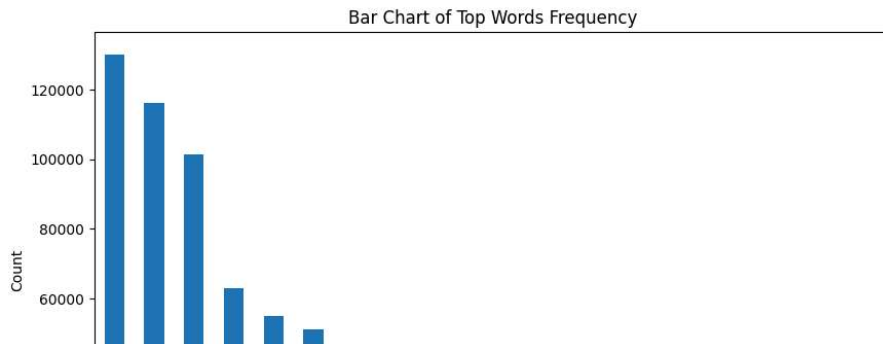
def get_top_n_words(corpus, n=None):
    vec = CountVectorizer().fit(corpus)
    bag_of_words = vec.transform(corpus)
    sum_words = bag_of_words.sum(axis=0)
    words_freq = [(word, sum_words[0, idx])
                   for word, idx in vec.vocabulary_.items()]
    words_freq = sorted(words_freq, key=lambda x: x[1],
                        reverse=True)

    return words_freq[:n]

common_words = get_top_n_words(data['text'], 20)
df1 = pd.DataFrame(common_words, columns=['Review', 'count'])

df1.groupby('Review').sum()['count'].sort_values(ascending=False).plot(
    kind='bar',
    figsize=(10, 6),
    xlabel="Top Words",
    ylabel="Count",
    title="Bar Chart of Top Words Frequency"
)
```

```
<Axes: title={'center': 'Bar Chart of Top Words Frequency'}, xlabel='Top Words',
ylabel='Count'>
```



```
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.linear_model import LogisticRegression
```

```
x_train, x_test, y_train, y_test = train_test_split(data['text'],
                                                    data['class'],
                                                    test_size=0.25)
```

```
from sklearn.feature_extraction.text import TfidfVectorizer
```

```
vectorization = TfidfVectorizer()
x_train = vectorization.fit_transform(x_train)
x_test = vectorization.transform(x_test)
```

```
from sklearn.linear_model import LogisticRegression
```

```
model = LogisticRegression()
model.fit(x_train, y_train)
```

```
# testing the model
print(accuracy_score(y_train, model.predict(x_train)))
print(accuracy_score(y_test, model.predict(x_test)))
```

```
0.9936774614859449
0.9873552983081033
```

```
from sklearn.tree import DecisionTreeClassifier
```

```
model = DecisionTreeClassifier()
model.fit(x_train, y_train)
```

```
# testing the model
print(accuracy_score(y_train, model.predict(x_train)))
print(accuracy_score(y_test, model.predict(x_test)))
```

```
1.0
0.9957257346393589
```

```
from sklearn import metrics
cm = metrics.confusion_matrix(y_test, model.predict(x_test))
```

```
cm_display = metrics.ConfusionMatrixDisplay(confusion_matrix=cm,
                                              display_labels=[False, True])
```