

```
In [1]: 1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 from matplotlib import style
5 style.use('ggplot')
6 import re
7 from nltk.tokenize import word_tokenize
8 from nltk.stem import PorterStemmer
9 from nltk.corpus import stopwords
10 stop_words = set(stopwords.words('english'))
11 from sklearn.feature_extraction.text import TfidfVectorizer
12 from sklearn.model_selection import train_test_split
13
14
```

```
In [2]: 1 df = pd.read_csv('IMDB Dataset.csv')
2 df.head()
3
```

```
Out[2]:
```

	review	sentiment
0	One of the other reviewers has mentioned that ...	positive
1	A wonderful little production.   The...	positive
2	I thought this was a wonderful way to spend ti...	positive
3	Basically there's a family where a little boy ...	negative
4	Petter Mattei's "Love in the Time of Money" is...	positive

```
In [3]: 1 df.shape
```

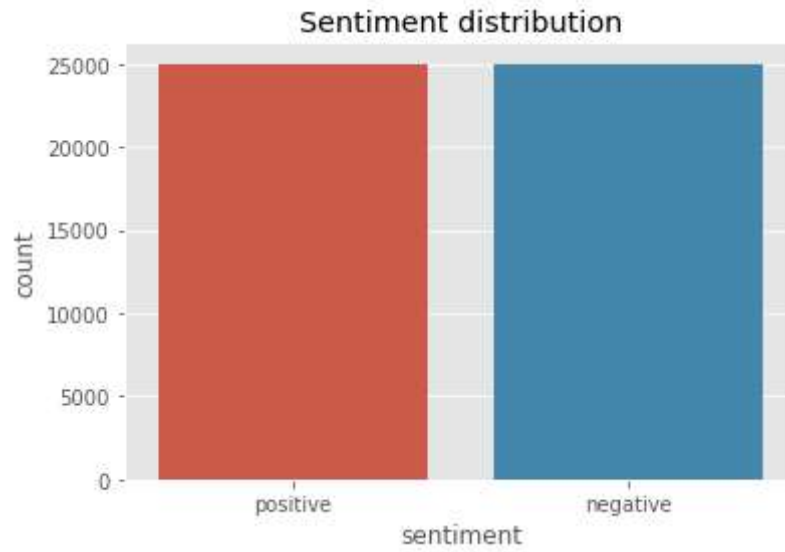
```
Out[3]: (50000, 2)
```

```
In [4]: 1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 50000 entries, 0 to 49999
Data columns (total 2 columns):
#   Column      Non-Null Count  Dtype
---  -
0   review      50000 non-null  object
1   sentiment   50000 non-null  object
dtypes: object(2)
memory usage: 781.4+ KB
```

```
In [5]: 1 sns.countplot(x='sentiment', data=df)
        2 plt.title("Sentiment distribution")
```

```
Out[5]: Text(0.5, 1.0, 'Sentiment distribution')
```



```
In [6]: 1 for i in range(5):
        2     print("Review: ", [i])
        3     print(df['review'].iloc[i], "\n")
        4     print("Sentiment: ", df['sentiment'].iloc[i], "\n\n")
```

Review: [0]

One of the other reviewers has mentioned that after watching just 1 Oz episode you'll be hooked. They are right, as this is exactly what happened with me.<br /><br />The first thing that struck me about Oz was its brutality and unflinching scenes of violence, which set in right from the word GO. Trust me, this is not a show for the faint hearted or timid. This show pulls no punches with regards to drugs, sex or violence. Its is hardcore, in the classic use of the word.<br /><br />It is called OZ as that is the nickname given to the Oswald Maximum Security State Penitentiary. It focuses mainly on Emerald City, an experimental section of the prison where all the cells have glass fronts and face inwards, so privacy is not high on the agenda. Em City is home to many..Aryans, Muslims, gangstas, Latinos, Christians, Italians, Irish and more....so scuffles, death stares, dodgy dealings and shady agreements are never far away.<br /><br />I would say the main appeal of the show is due to the fact that it goes where other shows wouldn't dare. Forget pretty pictures painted for mainstream audiences, forget charm, forget romance...OZ doesn't mess around. The first episode I ever saw struck me as so nasty it was surreal, I couldn't say I was ready for it, but as I watched more, I developed a taste for Oz, and got accustomed to the high levels of graphic violence. Not just violence, but injustice (crooked guards who'll be sold out for a nickel, inmates who'll kill on order and get away with it, well mannered, middle class inmates being turned into prison bitches due to their lack of street skills or prison experience) Watching Oz, you may become comfortable with what is uncomfortable viewing....thats if you can get in touch with your darker side.

Sentiment: positive

Review: [1]

A wonderful little production. <br /><br />The filming technique is very unassuming- very old-time-BBC fashion and gives a comforting, and sometimes discomforting, sense of realism to the entire piece. <br /><br />The actors are extremely well chosen- Michael Sheen not only "has got all the polari" but he has all the voices down pat too! You can truly see the seamless editing guided by the references to Williams' diary entries, not only is it well worth the watching but it is a terrifically written and performed piece. A masterful production about one of the great master's of comedy and his life. <br /><br />The realism really comes home with the little things: the fantasy of the guard which, rather than use the traditional 'dream' techniques remains solid then disappears. It plays on our knowledge and our senses, particularly with the scenes concerning Orton and Halliwell and the sets (particularly of their flat with Halliwell's murals decorating every surface) are terribly well done.

Sentiment: positive

Review: [2]

I thought this was a wonderful way to spend time on a too hot summer weekend, sitting in the air conditioned theater and watching a light-hearted comedy. The plot is simplistic, but the dialogue is witty and the characters are likable (even the well bread suspected serial killer). While some may be disappointed when they realize this is not Match Point 2: Risk Addiction, I thought it was proo

f that Woody Allen is still fully in control of the style many of us have grown to love.<br /><br />This was the most I'd laughed at one of Woody's comedies in years (dare I say a decade?). While I've never been impressed with Scarlet Johanson, in this she managed to tone down her "sexy" image and jumped right into a average, but spirited young woman.<br /><br />This may not be the crown jewel of his career, but it was wittier than "Devil Wears Prada" and more interesting than "Superman" a great comedy to go see with friends.

Sentiment: positive

Review: [3]

Basically there's a family where a little boy (Jake) thinks there's a zombie in his closet & his parents are fighting all the time.<br /><br />This movie is slower than a soap opera... and suddenly, Jake decides to become Rambo and kill the zombie.<br /><br />OK, first of all when you're going to make a film you must Decide if its a thriller or a drama! As a drama the movie is watchable. Parents are divorcing & arguing like in real life. And then we have Jake with his closet which totally ruins all the film! I expected to see a BOOGEYMAN similar movie, and instead i watched a drama with some meaningless thriller spots.<br /><br />3 out of 10 just for the well playing parents & descent dialogs. As for the shots with Jake: just ignore them.

Sentiment: negative

Review: [4]

Petter Mattei's "Love in the Time of Money" is a visually stunning film to watch. Mr. Mattei offers us a vivid portrait about human relations. This is a movie that seems to be telling us what money, power and success do to people in the different situations we encounter. <br /><br />This being a variation on the Arthur Schnitzler's play about the same theme, the director transfers the action to the present time New York where all these different characters meet and connect. Each one is connected in one way, or another to the next person, but no one seems to know the previous point of contact. Stylishly, the film has a sophisticated luxurious look. We are taken to see how these people live and the world they live in their own habitat.<br /><br />The only thing one gets out of all these souls in the picture is the different stages of loneliness each one inhabits. A big city is not exactly the best place in which human relations find since fulfillment, as one discerns is the case with most of the people we encounter.<br /><br />The acting is good under Mr. Mattei's direction. Steve Buscemi, Rosario Dawson, Carol Kane, Michael Imperioli, Adrian Grenier, and the rest of the talented cast, make these characters come alive.<br /><br />We wish Mr. Mattei good luck and await anxiously for his next work.

Sentiment: positive

In [7]:

```
1 def no_of_words(text):  
2     words= text.split()  
3     word_count = len(words)  
4     return word_count
```

```
In [8]: 1 df['word count'] = df['review'].apply(no_of_words)
```

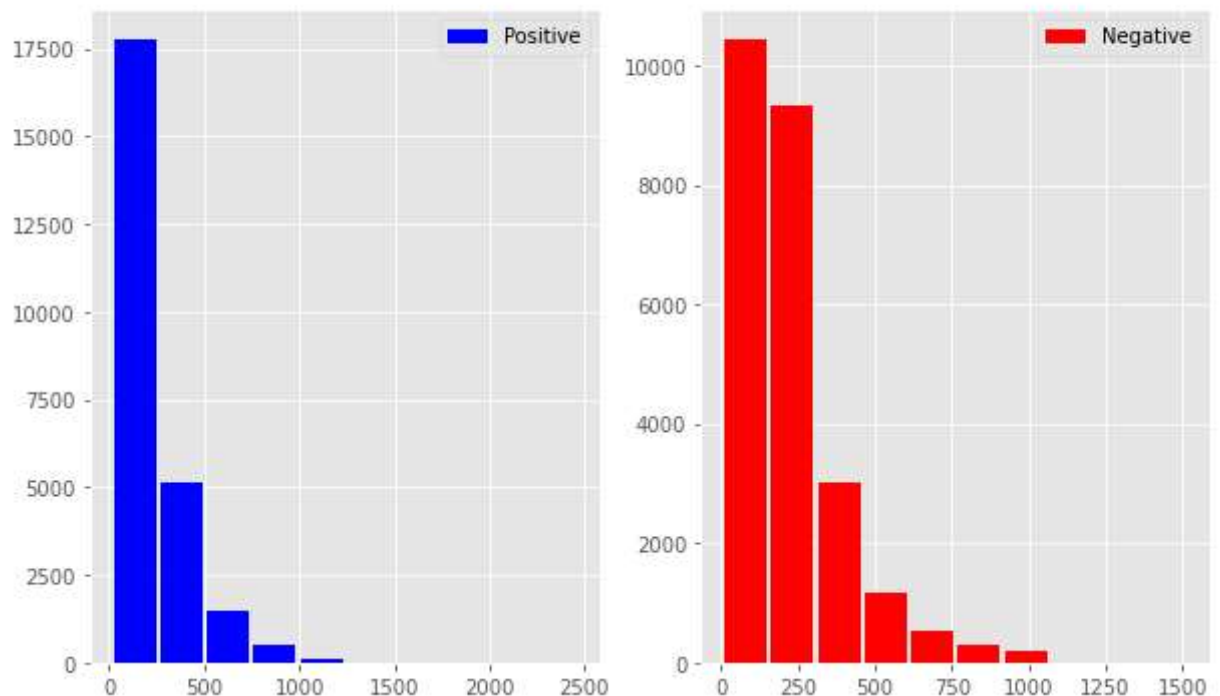
```
In [9]: 1 df.head()
```

Out[9]:

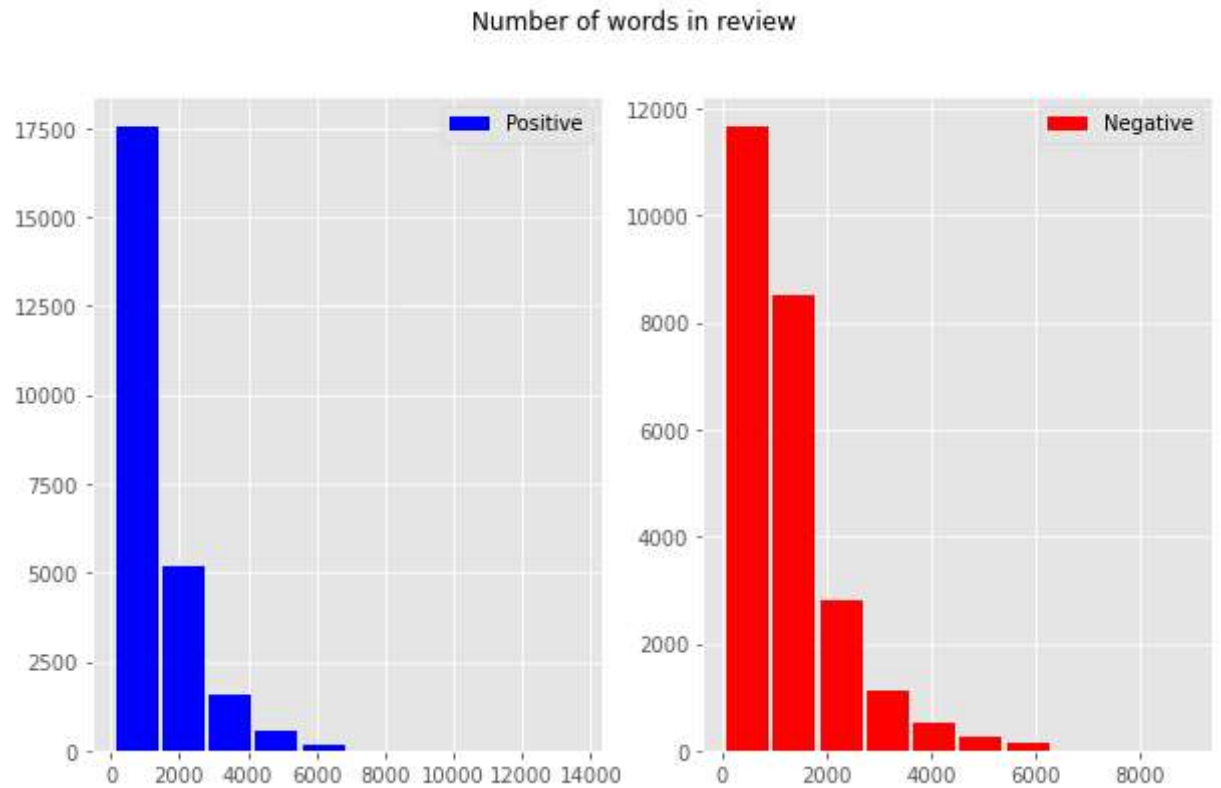
	review	sentiment	word count
0	One of the other reviewers has mentioned that ...	positive	307
1	A wonderful little production.   The...	positive	162
2	I thought this was a wonderful way to spend ti...	positive	166
3	Basically there's a family where a little boy ...	negative	138
4	Petter Mattei's "Love in the Time of Money" is...	positive	230

```
In [10]: 1 fig, ax = plt.subplots(1,2, figsize=(10,6))
2 ax[0].hist(df[df['sentiment'] == 'positive']['word count'], label='Positive')
3 ax[0].legend(loc='upper right');
4 ax[1].hist(df[df['sentiment'] == 'negative']['word count'], label='Negative')
5 ax[1].legend(loc='upper right');
6 fig.suptitle("Number of words in review")
7 plt.show()
```

Number of words in review



```
In [11]: 1 fig, ax = plt.subplots(1,2, figsize=(10,6))
2 ax[0].hist(df[df['sentiment'] == 'positive']['review'].str.len(), label='Pos
3 ax[0].legend(loc='upper right');
4 ax[1].hist(df[df['sentiment'] == 'negative']['review'].str.len(), label='Neg
5 ax[1].legend(loc='upper right');
6 fig.suptitle("Number of words in review")
7 plt.show()
```



```
In [12]: 1 df.sentiment.replace("positive", 1, inplace=True)
2 df.sentiment.replace("negative", 2, inplace=True)
```

In [13]:

```
1 df.head()
```

Out[13]:

	review	sentiment	word count
0	One of the other reviewers has mentioned that ...	1	307
1	A wonderful little production.   The...	1	162
2	I thought this was a wonderful way to spend ti...	1	166
3	Basically there's a family where a little boy ...	2	138
4	Petter Mattei's "Love in the Time of Money" is...	1	230

In [14]:

```
1 def data_processing(text):
2     text= text.lower()
3     text = re.sub('<br />', '', text)
4     text = re.sub(r"https\S+|www\S+|http\S+", '', text, flags = re.MULTILINE)
5     text = re.sub(r'\@w+|\#', '', text)
6     text = re.sub(r'^\w\s', '', text)
7     text_tokens = word_tokenize(text)
8     filtered_text = [w for w in text_tokens if not w in stop_words]
9     return " ".join(filtered_text)
```

In [16]:

```
1 duplicated_count = df.duplicated().sum()
2 print("Number of duplicate entries: ", duplicated_count)
```

Number of duplicate entries: 418

In [17]:

```
1 df = df.drop_duplicates('review')
```

In [18]:

```
1 stemmer = PorterStemmer()
2 def stemming(data):
3     text = [stemmer.stem(word) for word in data]
4     return data
```

In [19]:

```
1 df.review = df['review'].apply(lambda x: stemming(x))
```

In [20]:

```
1 df['word count'] = df['review'].apply(no_of_words)
2 df.head()
```

Out[20]:

	review	sentiment	word count
0	One of the other reviewers has mentioned that ...	1	307
1	A wonderful little production.   The...	1	162
2	I thought this was a wonderful way to spend ti...	1	166
3	Basically there's a family where a little boy ...	2	138
4	Petter Mattei's "Love in the Time of Money" is...	1	230

```
In [21]: 1 pos_reviews = df[df.sentiment == 1]
          2 pos_reviews.head()
```

Out[21]:

	review	sentiment	word count
0	One of the other reviewers has mentioned that ...	1	307
1	A wonderful little production.   The...	1	162
2	I thought this was a wonderful way to spend ti...	1	166
4	Petter Mattei's "Love in the Time of Money" is...	1	230
5	Probably my all-time favorite movie, a story o...	1	119

```
In [23]: 1 from collections import Counter
          2 count = Counter()
          3 for text in pos_reviews['review'].values:
          4     for word in text.split():
          5         count[word] +=1
          6 count.most_common(15)
```

Out[23]:

```
[('the', 290932),
 ('and', 165372),
 ('a', 155251),
 ('of', 148673),
 ('to', 127921),
 ('is', 107829),
 ('in', 90145),
 ('that', 62191),
 ('I', 61914),
 ('it', 53162),
 ('this', 51403),
 ('/><br', 48800),
 ('as', 46419),
 ('with', 43153),
 ('was', 41934)]
```

```
In [24]: 1 pos_words = pd.DataFrame(count.most_common(15))
          2 pos_words.columns = ['word', 'count']
          3 pos_words.head()
```

Out[24]:

	word	count
0	the	290932
1	and	165372
2	a	155251
3	of	148673
4	to	127921



```
In [26]: 1 neg_reviews = df[df.sentiment == 2]
        2 neg_reviews.head()
```

Out[26]:

	review	sentiment	word count
3	Basically there's a family where a little boy ...	2	138
7	This show was an amazing, fresh & innovative i...	2	174
8	Encouraged by the positive comments about this...	2	130
10	Phil the Alien is one of those quirky films wh...	2	96
11	I saw this movie when I was about 12 when it c...	2	180

```
In [28]: 1 count = Counter()
        2 for text in neg_reviews['review'].values:
        3     for word in text.split():
        4         count[word] +=1
        5 count.most_common(15)
```

Out[28]:

```
[('the', 273542),
 ('a', 149568),
 ('and', 134388),
 ('of', 132924),
 ('to', 131974),
 ('is', 93846),
 ('in', 78593),
 ('I', 69449),
 ('that', 63687),
 ('this', 61414),
 ('it', 53956),
 ('/><br', 51411),
 ('was', 49969),
 ('for', 39373),
 ('with', 38797)]
```

```
In [29]: 1 neg_words = pd.DataFrame(count.most_common(15))
        2 neg_words.columns = ['word', 'count']
        3 neg_words.head()
```

Out[29]:

	word	count
0	the	273542
1	a	149568
2	and	134388
3	of	132924
4	to	131974

```
In [31]: 1 X = df['review']
        2 Y = df['sentiment']
```

```
In [32]: 1 vect = TfidfVectorizer()
2 X = vect.fit_transform(df['review'])
```

```
In [33]: 1 x_train, x_test, y_train, y_test = train_test_split(X, Y, test_size=0.3, ran
```

```
In [34]: 1 print("Size of x_train: ", (x_train.shape))
2 print("Size of y_train: ", (y_train.shape))
3 print("Size of x_test: ", (x_test.shape))
4 print("Size of y_test: ", (y_test.shape))
```

```
Size of x_train: (34707, 101895)
Size of y_train: (34707,)
Size of x_test: (14875, 101895)
Size of y_test: (14875,)
```

```
In [44]: 1 from sklearn.naive_bayes import MultinomialNB
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.metrics import accuracy_score, classification_report, confusion
4 import warnings
5 warnings.filterwarnings('ignore')
```

```
In [47]: 1 mnb = MultinomialNB()
2 mnb.fit(x_train, y_train)
3 mnb_pred = mnb.predict(x_test)
4 mnb_acc = accuracy_score(mnb_pred, y_test)
5 print("Test accuracy: {:.2f}%".format(mnb_acc*100))
```

Test accuracy: 86.11%

```
In [48]: 1 print(confusion_matrix(y_test, mnb_pred))
2 print("\n")
3 print(classification_report(y_test, mnb_pred))
```

```
[[6276 1195]
 [ 871 6533]]
```

	precision	recall	f1-score	support
1	0.88	0.84	0.86	7471
2	0.85	0.88	0.86	7404
accuracy			0.86	14875
macro avg	0.86	0.86	0.86	14875
weighted avg	0.86	0.86	0.86	14875

```
In [49]: 1 logreg = LogisticRegression()
2 logreg.fit(x_train, y_train)
3 logreg_pred = logreg.predict(x_test)
4 logreg_acc = accuracy_score(logreg_pred, y_test)
5 print("Test accuracy: {:.2f}%".format(logreg_acc*100))
```

Test accuracy: 89.60%

```
In [50]: 1 print(confusion_matrix(y_test, logreg_pred))
          2 print("\n")
          3 print(classification_report(y_test, logreg_pred))
```

```
[[6771  700]
 [ 847 6557]]
```

	precision	recall	f1-score	support
1	0.89	0.91	0.90	7471
2	0.90	0.89	0.89	7404
accuracy			0.90	14875
macro avg	0.90	0.90	0.90	14875
weighted avg	0.90	0.90	0.90	14875