

DEPENDENCIES

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import plotly.express as px
import nltk
from sklearn.feature_extraction.text import CountVectorizer
from wordcloud import WordCloud, STOPWORDS
from nltk.stem import WordNetLemmatizer
from nltk.tokenize import word_tokenize
import re, string, unicodedata
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score, f1_score
from sklearn.model_selection import train_test_split
from string import punctuation
from nltk import pos_tag
from nltk.corpus import wordnet
import re
import warnings
warnings.filterwarnings('ignore')
import matplotlib.pyplot as plt
```

DATA UPLOADING

```
df=pd.read_csv('/content/drive/MyDrive/Sentiment Analysis/Data/IMDB-Dataset.csv', encoding='latin-1')
```

DATA CLEANING

```
from nltk.corpus import stopwords
stop_words = stopwords.words('english')
new_stopwords = ["would", "shall", "could", "might"]
stop_words.extend(new_stopwords)
stop_words.remove("not")
stop_words=set(stop_words)
print(stop_words)
```

```
{'then', 'hers', 'y', 'over', 'himself', 'during', 'haven't', 'aren't', 'couldn't', 'doing', 'further', 'no', 'would', 'it', 'up',
```

PREPROCESSING DATA BY REMOVING STOPWORDS, URLS, SPECIAL CHARACTERS AND EXPANDING CONTRACTIONS

```
#Removing special character
def remove_special_character(content):
    return re.sub('\W+', ' ', content )#re.sub('[^&@#!]*\]', '', content)

# Removing URL's
def remove_url(content):
    return re.sub(r'http\S+', '', content)

#Removing the stopwords from text
def remove_stopwords(content):
    clean_data = []
    for i in content.split():
        if i.strip().lower() not in stop_words and i.strip().lower().isalpha():
            clean_data.append(i.strip().lower())
    return " ".join(clean_data)

# Expansion of english contractions
def contraction_expansion(content):
    content = re.sub(r"won't", "would not", content)
    content = re.sub(r"can't", "can not", content)
    content = re.sub(r"don't", "do not", content)
    content = re.sub(r"shouldn't", "should not", content)
```

```
content = re.sub(r"needn't", "need not", content)
content = re.sub(r"hasn't", "has not", content)
content = re.sub(r"haven't", "have not", content)
content = re.sub(r"weren't", "were not", content)
content = re.sub(r"mightn't", "might not", content)
content = re.sub(r"didn't", "did not", content)
content = re.sub(r"n't", " not", content)
'''content = re.sub(r'\re', " are", content)
content = re.sub(r'\s', " is", content)
content = re.sub(r'\d', " would", content)
content = re.sub(r'\ll', " will", content)
content = re.sub(r'\t', " not", content)
content = re.sub(r'\ve', " have", content)
content = re.sub(r'\m', " am", content)'''
return content
```

#Data preprocessing using above methods

```
def data_cleaning(content):
    content = contraction_expansion(content)
    content = remove_special_character(content)
    content = remove_url(content)

    content = remove_stopwords(content)
    return content
```

APPLYING THE DATA CLEANING ON OUR DATASET

```
pd.options.display.max_colwidth = 1000
df['Reviews_clean']=df['Reviews'].apply(data_cleaning)
df.head(5)
```

Ratings		Reviews	Movies	Resenhas	Reviews_clean
0	1		Disaster Movie	<p>* IsenÃ§Ã£o de responsabilidade: eu sÃ³ assisti esse filme como um acordo condicional. E eu vejo filmes de graÃ§a. Eu nÃ£o seria pego morto dando meu dinheiro suado a esses idiotas. Bem, para explicar a profundidade desse 'filme', eu poderia escrever minha crÃ¡tica mais curta de todos os tempos. NÃ£o vÃ¡ este filme. Ã de longe o filme mais estÃ³pido, lamenta, preguiÃ§oso e inacreditavelmente UNFUNNY que eu jÃ¡ vi. Ã um desastre total. Mas como o meu Ã³dio por este filme e por outros, se estende muito alÃ©m de uma experiÃªncia, acho que vou continuar um pouco. NÃ£o conheÃ§o nenhuma das pessoas do filme alÃ©m de Carmen Electra, Vanessa Minnillo, e Kim Kardashian, mas isso nÃ£o importa. Eles sÃ£o todos horrÃveis, embora eu ache que esse seja o ponto. A ediÃ§Ã£o Ã horrÃvel e, possivelmente, erros de continuidade flagrantes tornam essa porcariÃa ainda mais horrÃvel do que eu pensava. Agora eu sei que esses filmes nÃ£o devem ser sÃ©rios, mas vamos lÃ¡, Ã© o cinema 101 que se alguÃ©m f...</p>	<p>disclaimer watched movie conditional agreement see films free not caught dead giving hard earned money idiots well explain depth film write shortest review ever not see movie far stupidest lamest lazy unbelievably unfunny movie ever seen total disaster since hatred movie others like extends far beyond one viewing think go bit not know people movie besides carmen electra vanessa minnillo kim kardashian not matter horrible though think point editing flat horrible possibly blatant continuity errors make crapfast even crappier thought know films not supposed serious come film making someone gets minor facial cut next shot someone gets cut sword blood least cut though since narnia films get away give disaster movie pass jokes thoughtless mindless physical gags obviously take popular movies last year late well including best picture nominees know saddest thing stupid movies not care much money make many cameos sorry ass excuses films taking away jobs actors writers directors truly deserv...</p>
				<p>Estou escrevendo isso na esperanÃ§a de que isso seja colocado sobre a revisÃ£o anterior deste "filme". Como alguÃ©m pode achar divertido esse desleixo estÃ¡ completamente alÃ©m de mim. Antes de mais nada, um filme de parÃ³dia intitulado "Filme de desastre" deveria ser, de fato, uma parÃ³dia de filmes de desastre. Agora eu jÃ¡ vi 1 (sim, conte-os, 1) filme de desastre sendo falsificado, sendo "Twister". Como Juno, Homem de Ferro, Batman, O Hulk, Alvin e os Esquilos, Amy Winehouse ou Hancock se registram como filmes de Desastre? Selzterwater e Failburg mostraram mais uma vez que</p>	<p>writing hopes gets put previous review film anyone find slop entertaining completely beyond first spoof film entitled disaster movie indeed spoof disaster films seen yes count disaster film spoofed twister juno iron man batman hulk alvin chipmunks amy winehouse hancock register disaster films selzterwater failburg shown lack sort writing skill humor unfortunately tortured date movie epic movie know exactly expect two plot jokes bad references cheaply remade scenes films someone informed satire copy paste one film another though not sav</p>

DATA SUMMARY

two...no plot, no jokes just bad references

sido torturado com Date Movie e Epic

write not believe people still pay see

```
#Checking empty values
df.isna().sum()
```

```
Ratings      0
Reviews      0
Movies       25
Resenhas      0
Reviews_clean 0
dtype: int64
```

turd sandwich but instead I'm iust going to

cocÃ' mas em vez disso vou fazer

```
#Basic descriptive statistics
df['Ratings'].describe()
```

```
count    150000.000000
mean         5.500000
std         2.872291
min         1.000000
25%         3.000000
50%         5.500000
75%         8.000000
max        10.000000
Name: Ratings, dtype: float64
```

me up or pre pubescent annoying iine

que naverÃ uma ma de buggers irritantes

hitch movies continue make form

```
#Basic statistics on cleaned reviewed column
df['Reviews_clean'].describe()
```

```
count    150000
unique    149755
top      story soundtrack dialog graphic reasonableness entertainment overall
freq              10
Name: Reviews_clean, dtype: object
```

```
#Checking unique reviews and movies
print('Unique reviews:%s' % df.Reviews_clean.nunique())
```

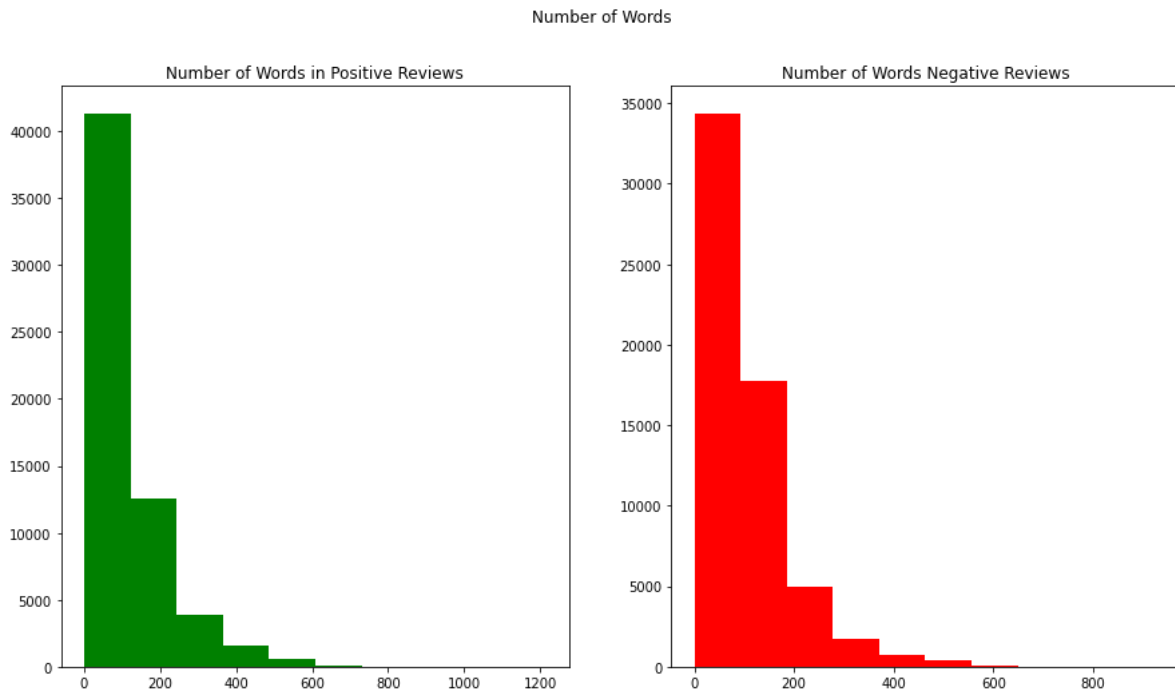
EXPLORATORY DATA ANALYSIS

```
1    15000
2    15000
4    15000
3    15000
5    15000
6    15000
8    15000
7    15000
10   15000
9    15000
Name: Ratings, dtype: int64
```

[illegible]

4/8

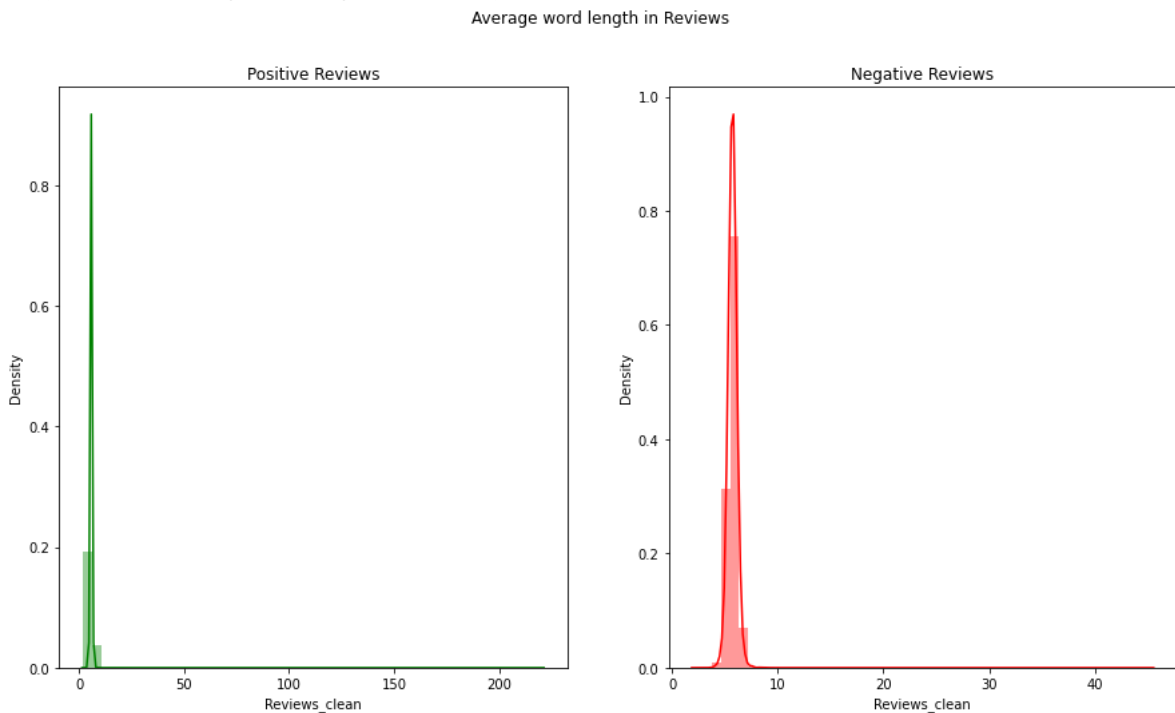

```
figure.suptitle('Number of Words')
plt.show()
```



```
#Visulaize average of words in reviews
```

```
figure,(pos_ax,neg_ax)=plt.subplots(1,2,figsize=(15,8))
pos_word=df[df['Ratings']>=7]['Reviews_clean'].str.split().apply(lambda x : [len(i) for i in x])
sns.distplot(pos_word.map(lambda x: np.mean(x)),ax=pos_ax,color='green')
pos_ax.set_title('Positive Reviews')
neg_word=df[df['Ratings']<=4]['Reviews_clean'].str.split().apply(lambda x : [len(i) for i in x])
sns.distplot(neg_word.map(lambda x: np.mean(x)),ax=neg_ax,color='red')
neg_ax.set_title('Negative Reviews')
figure.suptitle('Average word length in Reviews')
```

```
Text(0.5, 0.98, 'Average word length in Reviews')
```

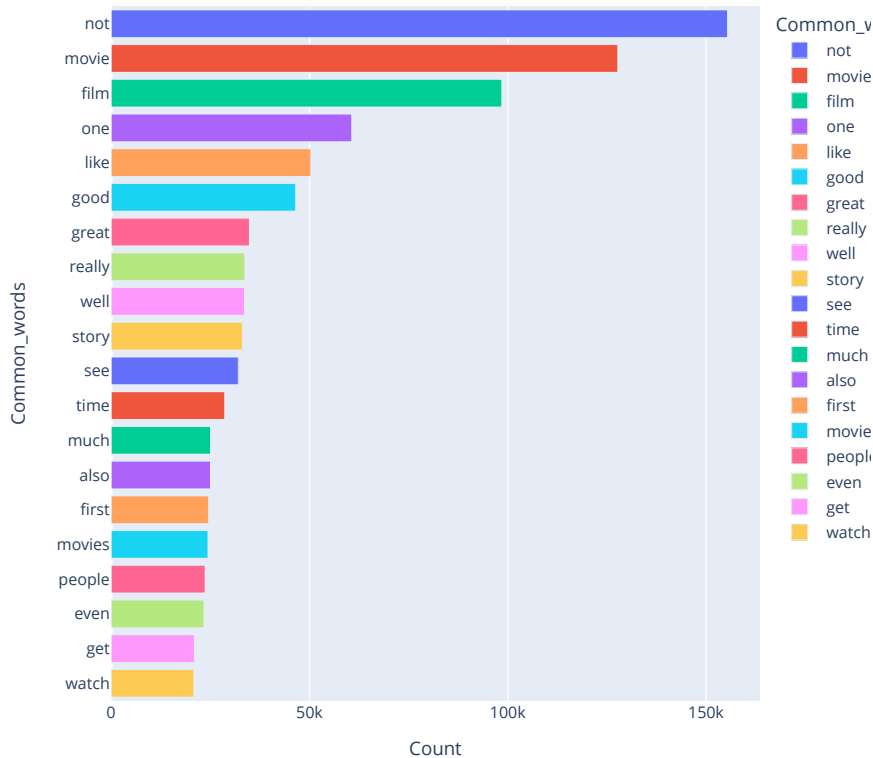


```
#Get important feature by using CountVectorizer
def get_top_text_ngrams(corpus, n, g):
    vec = CountVectorizer(ngram_range=(g, g)).fit(corpus)
    bag_of_words = vec.transform(corpus)
    sum_words = bag_of_words.sum(axis=0)
```

```
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]
```

```
most_common_uni = get_top_text_ngrams(df.Reviews_clean[df['Ratings']>=7],20,1)
most_common_uni = dict(most_common_uni)
temp = pd.DataFrame(columns = ["Common_words" , 'Count'])
temp["Common_words"] = list(most_common_uni.keys())
temp["Count"] = list(most_common_uni.values())
fig = px.bar(temp, x="Count", y="Common_words", title='Common Words in Positive Reviews', orientation='h',
             width=700, height=700,color='Common_words')
fig.show()
```

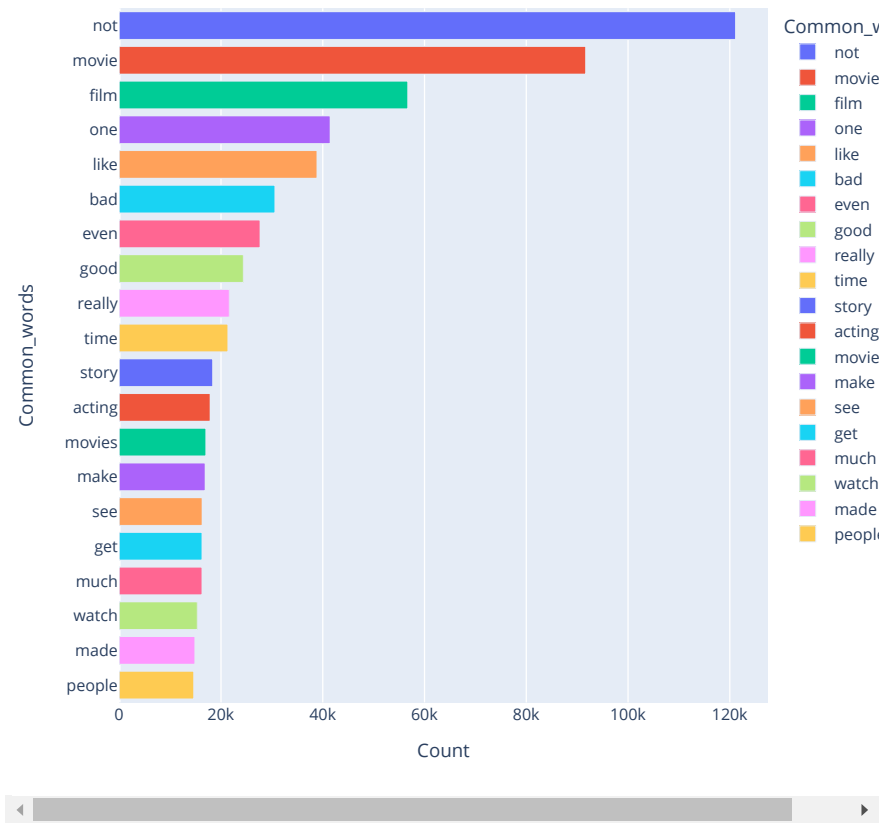
Common Words in Positive Reviews



```
most_common_uni = get_top_text_ngrams(df.Reviews_clean[df['Ratings']<4],20,1)
most_common_uni = dict(most_common_uni)
temp = pd.DataFrame(columns = ["Common_words" , 'Count'])
temp["Common_words"] = list(most_common_uni.keys())
temp["Count"] = list(most_common_uni.values())
fig = px.bar(temp, x="Count", y="Common_words", title='Common bigram in Negative Reviews', orientation='h',
             width=700, height=700,color='Common_words')
fig.show()
```



Common bigram in Negative Reviews



[Colab paid products](#) - [Cancel contracts here](#)

✓ 10s completed at 8:49 PM

