

STEP1:Finding a Dataset

```
In [1]: import pandas as pd
import numpy as np
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import BernoulliNB
import nltk
nltk.download('stopwords')

data = pd.read_csv("IMDB Dataset.csv")
print(data.head())
```

```
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\dell\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

	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

STEP2:Data Preparation, Tokenization, Stopwords Removal and Stemming

```
In [2]: '''Here we will:
1.remove links and all the special characters from the review column
2.tokenize and remove the stopwords from the review column
3.stem the words in the review column'''
import nltk
import re
nltk.download('stopwords')
stemmer = nltk.SnowballStemmer("english")
from nltk.corpus import stopwords
import string
stopword=set(stopwords.words('english'))

def clean(text):
    text = str(text).lower()
    text = re.sub('\[.*?\]', '', text)
    text = re.sub('https?://\S+|www\.\S+', '', text)
    text = re.sub('<.*?>+', '', text)
    text = re.sub('[%s]' % re.escape(string.punctuation), '', text)
    text = re.sub('\n', '', text)
    text = re.sub('\w*\d\w*', '', text)
    text = [word for word in text.split(' ') if word not in stopword]
    text=" ".join(text)
    text = [stemmer.stem(word) for word in text.split(' ')]
    text=" ".join(text)
    return text
data["review"] = data["review"].apply(clean)
```

```
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\dell\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

Step 3: Text Vectorization

```
In [3]: x = np.array(data["review"])
        y = np.array(data["sentiment"])

        cv = CountVectorizer()
        X = cv.fit_transform(x)
        X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                            test_size=0.20,
                                                            random_state=42)
```

Step 4: Text Classification

```
In [5]: from sklearn.linear_model import PassiveAggressiveClassifier
        model = PassiveAggressiveClassifier()
        model.fit(X_train, y_train)
        user = input("Enter a Text: ")
        data = cv.transform([user]).toarray()
        output = model.predict(data)
        print(output)
```

```
Enter a Text: good movie
['positive']
```

```
In [ ]:
```