## **Imports**

```
import pandas as pd
import numpy as np
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word tokenize
from nltk.stem import PorterStemmer
from sklearn.feature extraction.text import TfidfVectorizer
from sklearn.model selection import train test split
from sklearn.naive bayes import MultinomialNB
from sklearn.metrics import accuracy score, precision score,
recall score, f1 score, classification report
import matplotlib.pyplot as plt
import seaborn as sns
nltk.download('stopwords')
[nltk data] Downloading package stopwords to /root/nltk data...
[nltk data] Unzipping corpora/stopwords.zip.
True
nltk.download('punkt')
[nltk data] Downloading package punkt to /root/nltk data...
[nltk data] Unzipping tokenizers/punkt.zip.
True
```

### Load the Data

```
import pandas as pd
df = pd.read csv('/content/drive/MyDrive/SideBoys/sem
7/AIHC/drugsComTrain raw.csv')
df.head()
   uniqueID
                             drugName
                                                          condition \
0
                            Valsartan Left Ventricular Dysfunction
     206461
1
      95260
                           Guanfacine
                                                               ADHD
2
     92703
                               Lybrel
                                                      Birth Control
3
                           Ortho Evra
                                                      Birth Control
     138000
      35696 Buprenorphine / naloxone
                                                  Opiate Dependence
                                              review rating
date \
0 "It has no side effect, I take it in combinati...
                                                           9 20-May-
```

```
12
   "My son is halfway through his fourth week of ...
                                                             8 27-Apr-
1
10
   "I used to take another oral contraceptive, wh...
2
                                                                14-Dec-
09
3
   "This is my first time using any form of birth...
                                                                 3-Nov-
15
   "Suboxone has completely turned my life around...
                                                             9 27-Nov-
16
   usefulCount
0
            27
1
           192
2
            17
3
            10
4
            37
```

## Preprocessing

```
# Remove duplicates
df = df.drop duplicates()
# Handle missing values if any
df = df.dropna()
# Text cleaning
stop words = set(stopwords.words("english"))
stemmer = PorterStemmer()
def clean text(text):
    text = text.lower()
    words = word tokenize(text)
    words = [word for word in words if word.isalnum() and word not in
stop wordsl
    words = [stemmer.stem(word) for word in words]
    return " ".join(words)
df["cleaned review"] = df["review"].apply(clean text)
df.head()
   uniqueID
                             drugName
                                                           condition \
0
     206461
                            Valsartan
                                       Left Ventricular Dysfunction
1
      95260
                           Guanfacine
                                                                ADHD
2
                                                       Birth Control
      92703
                                Lybrel
3
     138000
                           Ortho Evra
                                                       Birth Control
      35696
             Buprenorphine / naloxone
                                                   Opiate Dependence
                                               review rating
```

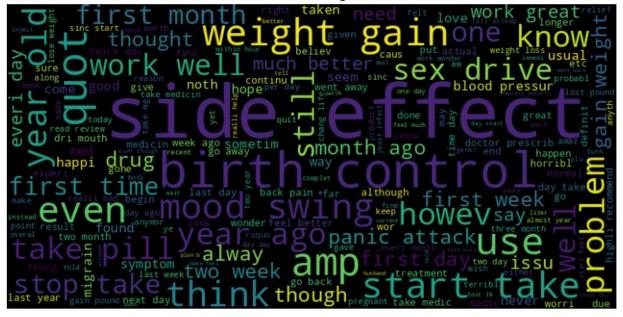
```
date \
   "It has no side effect, I take it in combinati...
                                                            9 2012-05-
20
   "My son is halfway through his fourth week of ...
                                                            8 2010-04-
1
27
2
   "I used to take another oral contraceptive, wh...
                                                            5 2009-12-
14
3
   "This is my first time using any form of birth...
                                                            8 2015-11-
03
   "Suboxone has completely turned my life around...
                                                            9 2016-11-
27
   usefulCount
                                                    cleaned review
sentiment
                     side effect take combin bystol 5 mg fish oil
            27
Neutral
           192
                son halfway fourth week intuniv becam concern ...
Positive
                use take anoth oral contracept 21 pill cycl li...
            17
Positive
            10 first time use form birth control 039 glad wen...
Positive
            37
                suboxon complet turn life around feel healthie...
Positive
```

## Exploratory Data Analysis (EDA)

```
# Generate summary statistics
summary stats = df.describe()
# Visualize drug ratings
plt.figure(figsize=(10, 6))
sns.histplot(df["rating"], bins=10, kde=True)
plt.title("Distribution of Drug Ratings")
plt.xlabel("Rating")
plt.ylabel("Count")
plt.show()
# Generate word frequency plot
from wordcloud import WordCloud
wordcloud = WordCloud(width=800, height=400).generate("
".join(df["cleaned review"]))
plt.figure(figsize=(10, 6))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.title("Word Cloud of Drug Reviews")
plt.show()
```

# 50000 - 50000 - 20000

Word Cloud of Drug Reviews



# Sentiment Analysis

from textblob import TextBlob

```
def get_sentiment(review):
    analysis = TextBlob(review)
    if analysis.sentiment.polarity > 0:
        return "Positive"
    elif analysis.sentiment.polarity < 0:
        return "Negative"
    else:
        return "Neutral"

df["sentiment"] = df["cleaned_review"].apply(get_sentiment)

# Calculate sentiment trends over time
df["date"] = pd.to_datetime(df["date"])
sentiment_over_time = df.groupby(df["date"].dt.year)
["sentiment"].value_counts().unstack().fillna(0)</pre>
```

## Text Mining and Feature Extraction

```
tfidf_vectorizer = TfidfVectorizer(max_features=1000)
X = tfidf_vectorizer.fit_transform(df["cleaned_review"])
```

## Machine Learning Models

```
X train, X test, y train, y test = train test split(X,
df["sentiment"], test_size=0.2, random state=42)
# Train a Multinomial Naive Bayes classifier
clf = MultinomialNB()
clf.fit(X train, y train)
y_pred = clf.predict(X test)
# Evaluate the model
accuracy = accuracy score(y test, y pred)
precision = precision_score(y_test, y_pred, average='weighted')
recall = recall_score(y_test, y_pred, average='weighted')
f1 = f1 score(y test, y pred, average='weighted')
print("Model Performance:")
print(f"Accuracy: {accuracy}")
print(f"Precision: {precision}")
print(f"Recall: {recall}")
print(f"F1 Score: {f1}")
print("\nClassification Report:")
print(classification report(y test, y pred))
```

Model Performance:

Accuracy: 0.6525872817955112 Precision: 0.7533290929215211 Recall: 0.6525872817955112 F1 Score: 0.5420630300069711

Classification Report:

CCU	133111CGC10	ii itcporci			
		precision	recall	f1-score	support
	Negative	0.97	0.11	0.19	8555
	Neutral	0.85	0.02	0.03	3546
	Positive	0.64	1.00	0.78	19979
	accuracy			0.65	32080
	macro avg	0.82	0.37	0.34	32080
wei	ghted avg	0.75	0.65	0.54	32080
	_				

## Data Visualization

```
# Visualize sentiment trends over time
sentiment_over_time.plot(kind="bar", stacked=True)
plt.title("Sentiment Trends Over Time")
plt.xlabel("Year")
plt.ylabel("Count")
plt.legend(title="Sentiment")
plt.show()
```

