

Sentiment Analysis on Amazon Fine Food Review

Introduction:

display my analy

```
# Imports
import pandas as pd
import numpy as np
import seaborn as sns
import nltk
from nltk import word_tokenize, sent_tokenize
from nltk.corpus import stopwords

%matplotlib inline
import matplotlib.pyplot as plt
```

Impo
cf.g

```
plt.figure(figsize=(10, 5))
plt.imshow(wordcloud)
plt.axis("off");
from wordcloud import WordCloud, STOPWORDS
from sklearn.feature_extraction.text import CountVectorizer
```

```

def __iter__(self): return 0

print("All imports installed...!")

```

All imports installed...!

In [2]:

```

amazon = pd.read_csv('Reviews.csv')
amazon.head()

```

Out[2]:

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1		1	5	1303862400 Good Quality Dog Food
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0		0	1	1346976000 Not as Advertised
2	3	B000LQOCHO	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1		1	4	1219017600 "Delight" says it all
3	4	B000UA0QIQ	A395BORC6FGVXV	Karl	3		3	2	1307923200 Cough Medicine
4	5	B006K2ZZ7K	A1UQRSCLF8GW1T	Michael D. Bigham "M. Wassir"	0		0	5	1350777600 Great taffy

Methodology:

To prepare for this analysis, I visualized the product scores from the dataset in a histogram using the plotly library.

In [3]:

```

# Visualizing Product Scores - Histogram

fig = px.histogram(amazon, x="Score")
fig.update_layout(title_text = "Product Score")
fig.show()

```

Product Score

A histogram titled "Product Score" showing the distribution of scores. The x-axis represents the score, and the y-axis represents the frequency, with labels at 250k, 300k, and 350k. A single blue bar is present at the highest score bin, indicating a frequency of approximately 350k.

3 4 B000UA00

After building the sentiment column, I also created word clouds to display the most frequently used words for both positive and negative product reviews, respectfully. In addition, I made a product sentiment histogram to show the distribution of reviews with sentiment across the dataset.

```
In [6]: # Positive Word Cloud  
  
positive = amazon[amazon["Sentiment"] == 1]  
  
text = " ".join(review for review in positive.Text)  
  
text = text.replace('\n', '')  
  
stopwords = set(STOPWORDS)  
stopwords.update(["br", "href"])
```

```
wordcloud.postive = WordCloud(width = 3000, height = 2000, random_state = 1, stopwords=stopwords, background_color="black")
plot_cloud(wordcloud.postive)
```

```
# Negative Word Cloud

negative = amazon[amazon["Sentiment"] == -1]

text = " ".join(review for review in negative.Text)

text = text.replace("\n", "")

stopwords = set(STOPWORDS)
stopwords.update(["good", "great", "br", "href"])
```

Product Sentiment

A histogram titled "Product Sentiment" showing the distribution of sentiment rates. The x-axis is labeled "Sentiment_Rate" and has two categories: "Positive" and "Negative". The y-axis is labeled "count" and ranges from 0 to 450k with increments of 50k. The "Positive" bar is orange and reaches a height of approximately 440k. The "Negative" bar is orange and reaches a height of approximately 80k.

Sentiment_Rate	count
Positive	~440k
Negative	~80k

Finally, I created a text classification model to train and establish the accuracy of my data. I start by pre-processing the textual data using NLTK to remove special characters, lowercasing text, and stopwords. Then, I test the accuracy of the sentiment model by performing the Multi Nominal Naive Bayes Classification function using the scikit-learn library.

```
In [9]: amazon.Summary = amazon['Summary'].str.replace('^\w\s+', '')  
amazon.head()  
  
<ipython-input-9-a813c68aaaa4>:1: FutureWarning:  
The default value of regex will change from True to False in a future version.  
  
Out[9]:
```

	Id	ProductId		UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian		1		1	5	1303862400	Good Quality Dog Food
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa		0		0	1	1346976000	Not as Advertised
2	3	B000LQOCHO	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"		1		1	4	1219017600	Delight says it all
3	4	B000UA0QIQ	A395BORC6FGVXV	Karl		3		3	2	1307923200	Cough Medicine
4	5	B006K2ZZ7K	A1UQRSCLF8GW1T	Michael D. Bigham "M."		0		0	5	1350777600	Great taffy

```
In [10]: sentiment_df = amazon[["Summary", "Sentiment"]]
sentiment_df.head()

Out[10]:
```

	Summary	Sentiment
0	Good Quality Dog Food	1
1	Not as Advertised	-1
2	Delight says it all	1
3	Cough Medicine	-1
4	Great taffy	1

Data Pre-Processing

```
In [11]: df = sentiment_df

df["Summary"] = df["Summary"].astype(str)

# Change to lowercasing for all text reviews in 'Summary'

df["Summary"] = df["Summary"].apply(lambda x: " ".join(x.lower() for x in x.split()))

stop = set(stopwords)

df["Summary"] = df["Summary"].apply(lambda x: " ".join(x for x in x.split() if x not in stop))
```