

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()

```

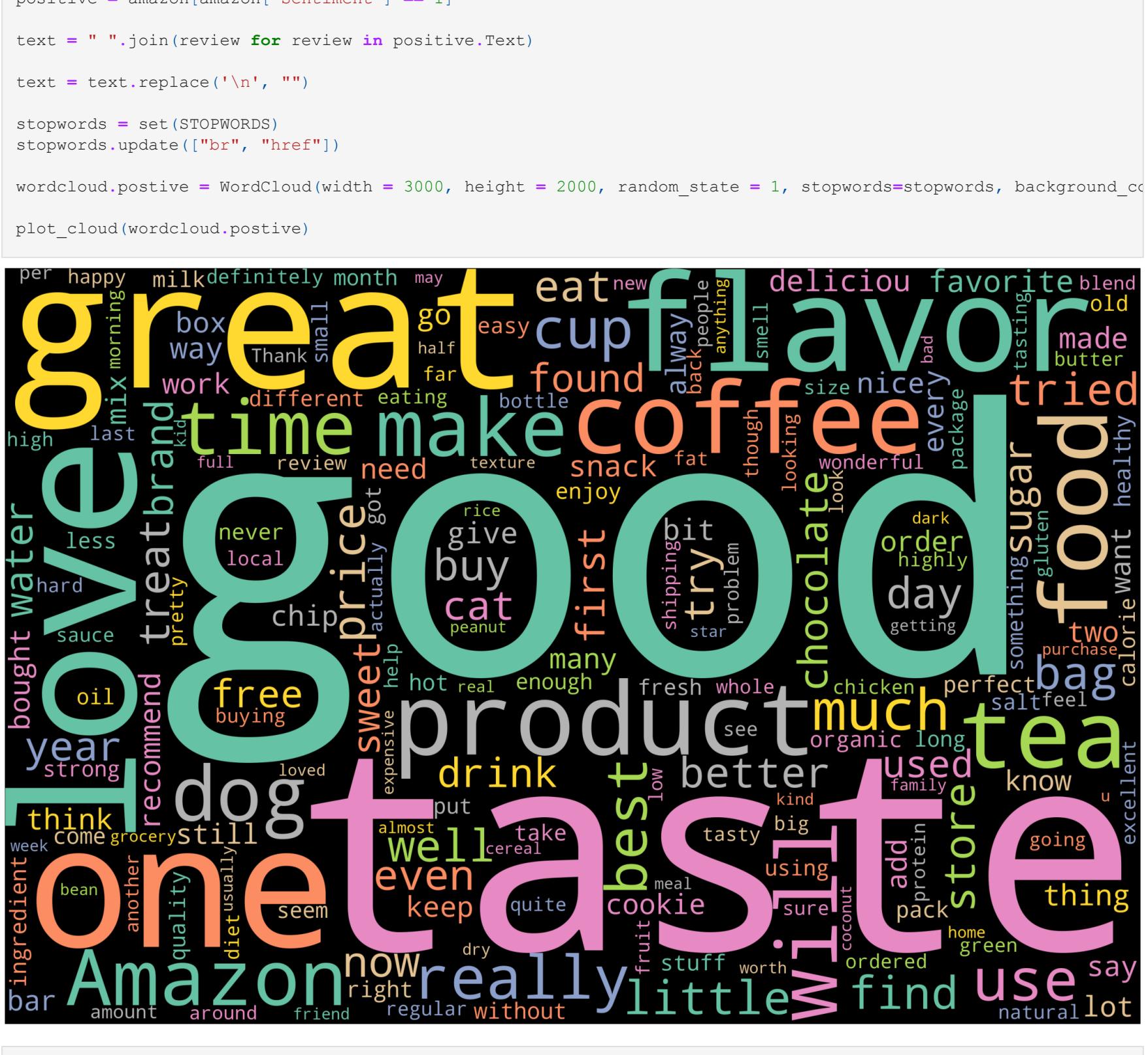
Product Score

The figure is a histogram titled "Product Score". The x-axis is labeled "Score" and has tick marks at 0, 1, 2, 3, 4, and 5. The y-axis is labeled "Count" and has tick marks at 200k, 250k, 300k, and 350k. The histogram shows a very tall bar at score 5, reaching a count of approximately 350,000. There are much shorter bars for scores 0, 1, 2, 3, and 4, indicating that most products have a score of 5.

2 - 3 B000EQQC

3 4 B000UA0QIQ A395BORC6FGVXV Karl 3 3 2

Wassir"



```
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"])

wordcloud.negative = WordCloud(width = 3000, height = 2000, random_state = 1, stopwords=stopwords, backgroundcolor="white")
```



A horizontal bar chart with a single data point at 0. The bar is orange and labeled "Positive".

From the orange histogram, we can see that the product sentiment is more positive than negative.

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.

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	B008T3GRG4	ATD8/F6ZCVE5NK	all pa	0	0	1	1346976000	Advertised
2	3	B000LQOCHO	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	1	4	1219017600	Delight says it all

