

5.1. A description of the dataset used.

The dataset that I am using are the reviews of user feedback on a product on amazon. There are almost 3,000 reviews for one product and over 30,000 reviews in this dataset. As the image below shows how the header column splits each category. The dataset that I will be looking into is the review.text category.

[illegible]

5.2. Details of the preprocessing steps.

To facilitate the preprocessing steps, we'll start by importing the necessary tools. We'll use Pandas to read the CSV file, spaCy for natural language processing, and TextBlob for sentiment analysis. Since dealing with a large DataFrame could be cumbersome, we'll streamline the process by initially focusing on the 'reviews.text' column. This approach helps us efficiently manage the data we need to analyse. After importing the tools, the next step involves data cleansing. We'll use the '.dropna()' method to remove any missing information or values from the dataset, minimising the potential for errors during analysis.

Following data cleansing, we'll create a list specifically for storing the reviews. This list will undergo preprocessing steps such as removing stopwords and tokenization, which aids in sentiment analysis and polarity identification.

Additionally, we'll implement a function that enables users to input a review for analysis. This function will output the sentiment score along with a classification indicating whether the sentiment is positive, negative, or neutral.

```

2 import pandas as pd
3 import numpy
4 from textblob import TextBlob
5 nlp = spacy.load("en_core_web_sm")
6
7 # Path to the CSV file
8 filePath = "amazon_product_reviews.csv"
9 df = pd.read_csv(filePath)
10
11 # Select the "review_text" column
12 reviews_data = df["review_text"]
13
14 # Remove the missing values or reviews
15 clean_data = reviews_data.dropna()
16
17 # Create a list
18 review_list = []
19
20 # Remove stopwords
21 for review in clean_data[0:20]: # change "clean_data" to "clean_data[XX:XX]" for min max values
22     nlp_review = nlp(review)
23     filtered_text = " ".join(token.text for token in nlp_review if not token.is_stop)
24     review_list.append(filtered_text)
25
26 print(review_list)

```

5.3. Evaluation of results.

In conclusion of my results, my model works. It is able to produce positive, negative and neutral results with their corresponding polarity. I have tested it with the first 20 results and as the picture showed how it would look. Looking at these results the polarity seems to sit in the region of 0.3-0.5.

[illegible]

5.4. Insights into the model's strengths and limitations.

With my model I decided to check what would happen if I did not remove any stop words as you could see in the pictures below. I would be able to conclude that the strength of my model with no stop words helps the accuracy of the model. However, the limitation would be with the small library with spacy 'en_core_web_sm' rather than having 'en_core_web_md'. With a larger library this might be able to produce much more accurate results to my model.

```

41 for list in review_list:
42     sentiment_label, sentiment_score = analyze_sentiment(list)
43     print("Review Sentiment:", sentiment_label, sentiment_score)
44
45 clean_data_limit = clean_data[:200]
46 for list in clean_data_limit:
47     sentiment_label, sentiment_score = analyze_sentiment(list)
48     print("actual Review Sentiment:", sentiment_label, sentiment_score)
49
50

```

[illegible]

```
actual Review Sentiment: Positive 0.325
actual Review Sentiment: Positive 0.8
actual Review Sentiment: Positive 0.6
actual Review Sentiment: Positive 0.7548833333333334
actual Review Sentiment: Positive 0.5555555555555556
actual Review Sentiment: Positive 0.375
actual Review Sentiment: Positive 0.5249999999999999
actual Review Sentiment: Positive 0.5444444444444444
actual Review Sentiment: Positive 0.4806666666666667
actual Review Sentiment: Positive 0.2166666666666667
actual Review Sentiment: Negative -0.4299999999999999
actual Review Sentiment: Positive 0.5949666666666667
actual Review Sentiment: Positive 0.7876666666666667
actual Review Sentiment: Positive 0.25
actual Review Sentiment: Positive 0.4166666666666667
actual Review Sentiment: Positive 0.7666666666666667
actual Review Sentiment: Negative -0.1963750000000000
actual Review Sentiment: Positive 0.5283333333333334
```

(These images are from the code from the left, producing results to the right.)

```
import pandas as pd
import spacy

from textblob import TextBlob

nlp = spacy.load('en_core_web_md')

# Path to the CSV file
filepath = "amazon_product_reviews.csv"
df = pd.read_csv(filepath)

# Select the relevant columns
review_data = df[['review_text']]

# Remove the missing values or review
clean_data = review_data.dropna()

# Iterate a list
review_list = []

# Remove duplicates
for review in clean_data[0:100]:
    # change "clean_data" to "clean_data[0:100]" for min size value
    nlp_review = nlp(review)
    # filter out the words that are not taken for taken in nlp_review if not taken, it_stop
    review_list.append(filter(review))

print(review_list)

def analyse_sentiment(review):
    # create a TextBlob object
    blob = TextBlob(review)
    # get the polarity score (where -1 is very negative, 0 is neutral, and 1 is very positive)
    polarity = blob.sentiment.polarity
    # determine sentiment label based on the polarity score
    if polarity > 0:
        sentiment_label = "Positive"
    elif polarity < 0:
        sentiment_label = "Negative"
    else:
        sentiment_label = "Neutral"
    # return both the sentiment label and the polarity score
    return sentiment_label, polarity

for list in review_list:
    sentiment_label, sentiment_score = analyse_sentiment(list)
    print("Review Sentiment: ", sentiment_label, sentiment_score)

clean_data_list = clean_data[0:100]

list in clean_data_list:
    sentiment_label, sentiment_score = analyse_sentiment(list)
    print("Final Review Sentiment: ", sentiment_label, sentiment_score)
```

[illegible]

(Image 1 and 2, are the code and terminal results to help explain)