

## Summary Report:

### 5.1. A description of the dataset used.

The dataset is stored in a DataFrame named `df`, which contains consumer reviews from Amazon. It has multiple columns, including `'reviews.text'`, which contains the text of the reviews. Missing values in the `'reviews.text'` column were removed using the `dropna()` function. Text cleaning techniques were used to preprocess the reviews, including removing stopwords, converting text to lowercase, and stripping white spaces. Sentiment analysis was performed on the cleaned reviews using `TextBlob` to analyse polarity. The function `analyse_polarity()` created was defined to calculate the polarity score of each review, and the results were stored in a new column named `'sentiment'` in the DataFrame. The similarity between two review products was calculated using `spaCy`'s `similarity()` function. Two reviews (from positions 40 and 10 in the DataFrame) were selected, and processed using `spaCy`, and their similarity score was calculated. The polarity and sentiment analysis were printed for a specific review (at position 300) to demonstrate the process.

### 5.2. Details of the preprocessing steps.

The NLP preprocessing pipeline is initiated within the `'remove_stopwords()'` function using `spaCy`'s language model (`nlp`). This pipeline includes tokenisation, part-of-speech tagging, dependency parsing, and other linguistic analyses. During the NLP preprocessing, `spaCy` tokenises the input text into individual tokens which are represented as a `spaCy` Token object. After tokenisation, the function iterates over each token in the processed doc object and filters out tokens that are stopwords using the `.is_stop` attribute such as "the", "is", "and", etc., that are typically removed during text preprocessing. The remaining non-stopword tokens are then joined together into a single string using `' '.join(filtered_tokens)`. This string represents the processed text without stopwords. Finally, the function converts the processed text to lowercase using `.lower()` and strips leading and trailing white spaces using `.strip()`. This ensures that the text is clean and properly formatted.

### 5.3. Evaluation of results.

The attribute `sentiment` is used to analyse and determine whether the review expresses positive, negative or neutral sentiment. The model is correctly working as demonstrated by the example below:

1. When index 0 is added, the results are Polarity score: -0.016666666666666663  
Sentiment: Negative. A polarity score of less than 0 is indicative of negative sentiment.
2. When index 3890 is analysed, the results are polarity score: 0.41111111111111115  
Sentiment: Positive. A polarity score greater than 0 is a positive sentiment
3. When index 4999 is analysed, the results are polarity score: 0.0  
Sentiment: Neutral. A polarity score of 0 is a neutral sentiment.

To conclude, the model on sample product reviews accurately and effectively identifies the sentiment of the review text, as evident in the example above.

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

There are various strengths to the model which include efficient data loading and preprocessing from a CSV file using Pandas which performs basic data cleaning operations such as removing missing values and stopwords. Additionally, spaCy is known for its efficiency, speed in processing large volumes of text data and ease of use with a range of pre-trained models. It is optimised for performance and can handle complex NLP tasks with minimal computational resources. Moreover, the model integrates both spaCy and TextBlob, two popular NLP libraries in Python, for text-processing tasks. SpaCy is used for tokenisation and stopwords removal, while TextBlob is used for sentiment analysis. Sentiment analysis and similarity comparison calculate the polarity score for each review. It also compares the similarity between two reviews using spaCy's similarity function.

However, there are a few limitations to the model. The lack of error-handling mechanisms may lead to issues if unexpected errors occur during data loading or processing. Adding error handling code would improve the robustness and reliability of the model and therefore the script. While TextBlob is a convenient tool for sentiment analysis, its accuracy may vary depending on the dataset and context. Using additional sentiment analysis models or custom-trained classifiers could provide more robust results, especially for domain-specific data. Likewise, there is room for improvement in the preprocessing steps as the model is relatively basic, focusing mainly on stopwords removal, lowercase conversion, and whitespace stripping. More advanced preprocessing techniques such as lemmatisation, stemming, or handling special characters may improve the quality of text data further. Lastly, the model exhibits a limited explanation of similarity calculation or what it represents. More explanation or documentation would enhance understanding for users.