

Data Collection and Preprocessing Phase

Date	15 March 2024
Team ID	SWTID1727420425
Project Title	Analysis of amazon cell phone reviews using nlp technique
Maximum Marks	6 Marks

Data Exploration and Pre processing Report

The Amazon cell phone review analysis involves two datasets: item metadata and customer reviews. Initial exploration revealed insights into review counts, ratings distribution, and product coverage. Data pre processing included handling missing values, cleaning review text, and converting review dates for trend analysis. Sentiment labels (positive, neutral, negative) were derived from ratings, and features like review length were added. This prepared data is now clean and ready for advanced NLP tasks like sentiment analysis and topic modeling.

Section	Description
Data Overview	The project utilizes two datasets for analyzing Amazon cell phone reviews: one containing product metadata (e.g., name, brand, price) and another with customer reviews (e.g., ratings, review text, review date). The review dataset provides insights into customer feedback, while the item dataset complements it with product details. Together, these datasets enable the application of NLP techniques to understand customer sentiments, trends, and preferences, forming the foundation for comprehensive analysis.
Resizing	In the context of this project, resizing involves standardizing the size and structure of review text for efficient NLP processing. Long reviews were truncated or summarized to a manageable length, while very short reviews were either padded or excluded, depending on their relevance. This step ensures uniformity across the dataset, improving the performance and accuracy of NLP models by focusing on meaningful and consistent input data.

Normalization	Normalization in this project focuses on transforming the review text into a consistent format for NLP analysis. This includes converting all text to lowercase, removing special characters, punctuation, and numbers, and standardizing spellings. Additionally, lemmatization was applied to reduce words to their base forms, ensuring uniformity while preserving meaning. These steps help eliminate noise and improve the accuracy of NLP techniques such as sentiment analysis and topic modeling.
Edge Detection	Edge detection is not directly applicable to textual data analysis in this project, as it is a technique primarily used in image processing to identify object boundaries. However, in the context of NLP, analogous techniques might involve identifying "edges" of significant meaning, such as detecting key phrases, transitions in sentiment, or abrupt shifts in topic within reviews, which can be explored through advanced methods like keyword extraction or topic segmentation.
Color Space Conversion	Color space conversion is not directly relevant to text-based NLP tasks in this project, as it pertains to image processing and the transformation of pixel data between color models (e.g., RGB to grayscale). For this project, the focus remains on textual data processing, where equivalent steps involve text cleaning and transformation rather than handling visual data.
Data Preprocessing Code Screenshots	
Loading Data	

	<pre>+ Code + Text</pre> <ul style="list-style-type: none"> Data Collection <pre>[] #load both items and review data items = pd.read_csv('/content/20191226-items.csv') reviews = pd.read_csv('/content/20191226-reviews.csv')</pre> <ul style="list-style-type: none"> Checking the structure of dataset <pre>[] # Check column names in each file print("Columns in reviews:", reviews.columns) print("Columns in items:", items.columns)</pre> <p>Columns in reviews: Index(['asin', 'name', 'rating', 'date', 'verified', 'title', 'body', 'helpfulVotes'], dtype='object')</p> <p>Columns in items: Index(['asin', 'brand', 'title', 'url', 'image', 'rating', 'reviewUrl', 'totalReviews', 'price', 'originalPrice'], dtype='object')</p>
Merging the Data	<ul style="list-style-type: none"> Merging the data <pre>[] # Merge the datasets on a common column (e.g., 'product_id') merged_data = pd.merge(reviews, items, on='asin', how='left')</pre> <pre>[] # Inspect the merged dataset print("Merged Data Shape:", merged_data.shape) print("Columns in Merged Data:", merged_data.columns)</pre> <p>Merged Data Shape: (67986, 17) Columns in Merged Data: Index(['asin', 'name', 'rating_x', 'date', 'verified', 'title_x', 'body', 'helpfulVotes', 'brand', 'title_y', 'url', 'image', 'rating_y', 'reviewUrl', 'totalReviews', 'price', 'originalPrice'], dtype='object')</p>
Handling missing data	<ul style="list-style-type: none"> handling missing data <pre>[] features = ['rating_x', 'verified', 'title_x', 'body', 'brand', 'price', 'originalPrice'] target = 'helpfulVotes'</pre> <pre>[] # Keep only necessary columns and drop rows with missing values filtered_data = merged_data[features + [target]].dropna()</pre> <pre>[] print("Filtered Data Shape:", filtered_data.shape)</pre> <p>Filtered Data Shape: (27069, 8)</p>
Data Exploration	<ul style="list-style-type: none"> Data exploration 

