

Amazon Product Reviews Sentiment Analysis using NLP

Authors:

- Wambui Githinji
- Lynette Mwiti

- Felix Njoroge
- Wilfred Lekishorumongi
- Monica Mwangi
- Joan Maina

Problem Statement

Reviews are critical to businesses as they offer insights into customer satisfaction, preferences and areas of improvement.

Businesses need to understand and interpret these reviews in order to cut through the competition. Lots of reviews are generated daily and manually analyzing them is impractical.

Objectives

Use Sentiment analysis to help the businesses get actionable insights from the feedback received from customers.

The approach taken with the analysis seeks to

- Determine the sentiment of the reviews (positive or negative) to understand overall customer satisfaction and feedback.
- Utilize sentiment analysis to help our stakeholders understand customer preferences across various products.
- Conduct exploratory data analysis to understand the distribution of sentiments over time, across brands and products.
- Leverage customer reviews to identify areas for improvement in products based on user experience.
- Build a classifier model to help predict reviews as positive or negative

Data Sources

Data for this project was obtained from Kaggle [repository]
(<https://www.kaggle.com/datasets/datafiniti/consumer-reviews-of-amazon-products?resource=download>)

The data represents:

- **Brand:** The brand name of the product being reviewed.
- **Categories:** Categories or tags that classify the product (e.g., electronics, home, books).
- **Keys:** Keywords or identifiers associated with the product.
- **Manufacturer:** The company or entity that manufactures the product.
- **Reviews.date:** The date when the review was posted.
- **Reviews.dateAdded:** Additional date-related information, possibly indicating when the review was added to the dataset.
- **Reviews.dateSeen:** Dates indicating when the review was observed or recorded (possibly by a data aggregator or platform).
- **Reviews.didPurchase:** Boolean (true/false) indicating whether the reviewer claims to have purchased the product.
- **Reviews.doRecommend:** Boolean (true/false) indicating whether the reviewer recommends the product.
- **Reviews.id:** Unique identifier for each review.
- **Reviews.numHelpful:** Number of users who found the review helpful.
- **Reviews.rating:** Rating given by the reviewer (typically on a scale such as 1 to 5 stars).
- **Reviews.sourceURLs:** URLs pointing to the source of the review.
- **Reviews.text:** The main body of the review text.
- **Reviews.title:** The title or headline of the review.
- **Reviews.userCity:** City location of the reviewer.
- **Reviews.userProvince:** Province or state location of the reviewer.
- **Reviews.username:** Username or identifier of the reviewer.

These are the variables this analysis will focus on to derive insights.

Methodology

The process can be divided into these many parts.

Data preparation

- **Text Cleaning:** Remove or handle punctuation, special characters, numbers, and stopwords
- **Tokenization:** Split text into words or subwords.
- **Text Normalization:** Convert text to lowercase, perform stemming or lemmatization.
- **Padding/Truncation: bold text** Ensure all text sequences are of the same length.
- **Train-Test Split:** Divide your data into training, validation, and test sets

EDA

Visualisations and insights. For each characteristic we will be:

- Creating visualisations
- Drawing conclusions
- Providing recommendations

Feature Engineering

In the feature engineering section, we process and transform the textual data for further analysis and modeling.

The methods used are;

- Sentiment Analysis
- Visualization with Word Clouds
- Text Vectorization to convert textual data into numerical form using TF-IDF and Count Vectorization.
- Word Embedding using Word2Vec and FastText
- Extraction of bigrams and trigrams

Model selection and building

The models used include a simple RNN and LSTM model

Hyperparameter tuning

Optimize hyperparameters for better performance

Model evaluation

Evaluate performance using the accuracy score

Analyze results

Look at the AUC/ROC curves and other evaluation tools

Data preparation

Importing Libraries

```

#Basic libraries
import pandas as pd
import numpy as np

#NLTK libraries
import nltk
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
from nltk.tokenize import word_tokenize

nltk.download('punkt')
nltk.download('stopwords')
nltk.download('wordnet')
import re
import string
!pip install wordcloud
from wordcloud import WordCloud, STOPWORDS
from nltk.stem.porter import PorterStemmer

from sklearn.feature_extraction.text import TfidfVectorizer

# Machine Learning libraries
import sklearn
from sklearn.svm import SVC
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import MinMaxScaler
from sklearn.ensemble import ExtraTreesClassifier
from sklearn.pipeline import make_pipeline
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.naive_bayes import BernoulliNB
from sklearn.neighbors import KNeighborsClassifier
from sklearn.multiclass import OneVsRestClassifier
from sklearn.svm import SVC
from sklearn.pipeline import Pipeline
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import label_binarize
from sklearn import svm, datasets
from sklearn import preprocessing

!pip install tensorflow
!pip install keras
!pip install numpy pandas scikit-learn

from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Embedding, SimpleRNN, Dense
from tensorflow.keras.preprocessing.text import Tokenizer

```

```
from tensorflow.keras.preprocessing.sequence import pad_sequences
```

#Metrics libraries

```
from sklearn import metrics
from sklearn.metrics import classification_report
from sklearn.model_selection import cross_val_score
from sklearn.metrics import roc_auc_score
from sklearn.metrics import roc_curve, auc
```

#Visualization libraries

```
import matplotlib.pyplot as plt
from matplotlib import rcParams
import seaborn as sns
from plotly import tools
import plotly.graph_objs as go
from plotly.offline import iplot
%matplotlib inline
```

#Ignore warnings

```
import warnings
warnings.filterwarnings('ignore')
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
```

```
[nltk_data] Package punkt is already up-to-date!
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
```

```
[nltk_data] Package stopwords is already up-to-date!
```

```
[nltk_data] Downloading package wordnet to /root/nltk_data...
```

```
[nltk_data] Package wordnet is already up-to-date!
```

```
Requirement already satisfied: wordcloud in
/usr/local/lib/python3.10/dist-packages (1.9.3)
```

```
Requirement already satisfied: numpy>=1.6.1 in
/usr/local/lib/python3.10/dist-packages (from wordcloud) (1.25.2)
```

```
Requirement already satisfied: pillow in
/usr/local/lib/python3.10/dist-packages (from wordcloud) (10.3.0)
```

```
Requirement already satisfied: matplotlib in
/usr/local/lib/python3.10/dist-packages (from wordcloud) (3.7.1)
```

```
Requirement already satisfied: contourpy>=1.0.1 in
/usr/local/lib/python3.10/dist-packages (from matplotlib->wordcloud)
(1.2.1)
```

```
Requirement already satisfied: cycler>=0.10 in
/usr/local/lib/python3.10/dist-packages (from matplotlib->wordcloud)
(0.12.1)
```

```
Requirement already satisfied: fonttools>=4.22.0 in
/usr/local/lib/python3.10/dist-packages (from matplotlib->wordcloud)
(4.53.0)
```

```
Requirement already satisfied: kiwisolver>=1.0.1 in
/usr/local/lib/python3.10/dist-packages (from matplotlib->wordcloud)
(1.4.5)
```

Requirement already satisfied: packaging>=20.0 in
/usr/local/lib/python3.10/dist-packages (from matplotlib->wordcloud)
(24.1)

Requirement already satisfied: pyparsing>=2.3.1 in
/usr/local/lib/python3.10/dist-packages (from matplotlib->wordcloud)
(3.1.2)

Requirement already satisfied: python-dateutil>=2.7 in
/usr/local/lib/python3.10/dist-packages (from matplotlib->wordcloud)
(2.9.0.post0)

Requirement already satisfied: six>=1.5 in
/usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.7->matplotlib->wordcloud) (1.16.0)

Requirement already satisfied: tensorflow in
/usr/local/lib/python3.10/dist-packages (2.15.0)

Requirement already satisfied: absl-py>=1.0.0 in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (1.4.0)

Requirement already satisfied: astunparse>=1.6.0 in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (1.6.3)

Requirement already satisfied: flatbuffers>=23.5.26 in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (24.3.25)

Requirement already satisfied: gast!=0.5.0,!0.5.1,!0.5.2,>=0.2.1
in /usr/local/lib/python3.10/dist-packages (from tensorflow) (0.6.0)

Requirement already satisfied: google-pasta>=0.1.1 in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (0.2.0)

Requirement already satisfied: h5py>=2.9.0 in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (3.11.0)

Requirement already satisfied: libclang>=13.0.0 in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (18.1.1)

Requirement already satisfied: ml-dtypes~=0.2.0 in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (0.2.0)

Requirement already satisfied: numpy<2.0.0,>=1.23.5 in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (1.25.2)

Requirement already satisfied: opt-einsum>=2.3.2 in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (3.3.0)

Requirement already satisfied: packaging in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (24.1)

Requirement already satisfied: protobuf!=4.21.0,!4.21.1,!4.21.2,!
=4.21.3,!4.21.4,!4.21.5,<5.0.0dev,>=3.20.3 in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (3.20.3)

Requirement already satisfied: setuptools in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (67.7.2)

Requirement already satisfied: six>=1.12.0 in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (1.16.0)

Requirement already satisfied: termcolor>=1.1.0 in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (2.4.0)

Requirement already satisfied: typing-extensions>=3.6.6 in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (4.12.2)

Requirement already satisfied: wrapt<1.15,>=1.11.0 in
/usr/local/lib/python3.10/dist-packages (from tensorflow) (1.14.1)

Requirement already satisfied: tensorflow-io-gcs-filesystem<=0.23.1 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (0.37.0)

Requirement already satisfied: grpcio<2.0,>=1.24.3 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (1.64.1)

Requirement already satisfied: tensorboard<2.16,>=2.15 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (2.15.2)

Requirement already satisfied: tensorflow-estimator<2.16,>=2.15.0 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (2.15.0)

Requirement already satisfied: keras<2.16,>=2.15.0 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (2.15.0)

Requirement already satisfied: wheel<1.0,>=0.23.0 in /usr/local/lib/python3.10/dist-packages (from astunparse>=1.6.0->tensorflow) (0.43.0)

Requirement already satisfied: google-auth<3,>=1.6.3 in /usr/local/lib/python3.10/dist-packages (from tensorboard<2.16,>=2.15->tensorflow) (2.27.0)

Requirement already satisfied: google-auth-oauthlib<2,>=0.5 in /usr/local/lib/python3.10/dist-packages (from tensorboard<2.16,>=2.15->tensorflow) (1.2.0)

Requirement already satisfied: markdown<=2.6.8 in /usr/local/lib/python3.10/dist-packages (from tensorboard<2.16,>=2.15->tensorflow) (3.6)

Requirement already satisfied: requests<3,>=2.21.0 in /usr/local/lib/python3.10/dist-packages (from tensorboard<2.16,>=2.15->tensorflow) (2.31.0)

Requirement already satisfied: tensorboard-data-server<0.8.0,>=0.7.0 in /usr/local/lib/python3.10/dist-packages (from tensorboard<2.16,>=2.15->tensorflow) (0.7.2)

Requirement already satisfied: werkzeug<=1.0.1 in /usr/local/lib/python3.10/dist-packages (from tensorboard<2.16,>=2.15->tensorflow) (3.0.3)

Requirement already satisfied: cachetools<6.0,>=2.0.0 in /usr/local/lib/python3.10/dist-packages (from google-auth<3,>=1.6.3->tensorboard<2.16,>=2.15->tensorflow) (5.3.3)

Requirement already satisfied: pyasn1-modules<=0.2.1 in /usr/local/lib/python3.10/dist-packages (from google-auth<3,>=1.6.3->tensorboard<2.16,>=2.15->tensorflow) (0.4.0)

Requirement already satisfied: rsa<5,>=3.1.4 in /usr/local/lib/python3.10/dist-packages (from google-auth<3,>=1.6.3->tensorboard<2.16,>=2.15->tensorflow) (4.9)

Requirement already satisfied: requests-oauthlib<=0.7.0 in /usr/local/lib/python3.10/dist-packages (from google-auth-oauthlib<2,>=0.5->tensorboard<2.16,>=2.15->tensorflow) (2.0.0)

Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.21.0->tensorboard<2.16,>=2.15->tensorflow) (3.3.2)

Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.21.0->tensorboard<2.16,>=2.15->tensorflow) (3.7)

```
Requirement already satisfied: urllib3<3,>=1.21.1 in
/usr/local/lib/python3.10/dist-packages (from requests<3,>=2.21.0-
>tensorboard<2.16,>=2.15->tensorflow) (2.0.7)
Requirement already satisfied: certifi>=2017.4.17 in
/usr/local/lib/python3.10/dist-packages (from requests<3,>=2.21.0-
>tensorboard<2.16,>=2.15->tensorflow) (2024.6.2)
Requirement already satisfied: MarkupSafe>=2.1.1 in
/usr/local/lib/python3.10/dist-packages (from werkzeug>=1.0.1-
>tensorboard<2.16,>=2.15->tensorflow) (2.1.5)
Requirement already satisfied: pyasn1<0.7.0,>=0.4.6 in
/usr/local/lib/python3.10/dist-packages (from pyasn1-modules>=0.2.1-
>google-auth<3,>=1.6.3->tensorboard<2.16,>=2.15->tensorflow) (0.6.0)
Requirement already satisfied: oauthlib>=3.0.0 in
/usr/local/lib/python3.10/dist-packages (from requests-
oauthlib>=0.7.0->google-auth-oauthlib<2,>=0.5-
>tensorboard<2.16,>=2.15->tensorflow) (3.2.2)
Requirement already satisfied: keras in
/usr/local/lib/python3.10/dist-packages (2.15.0)
Requirement already satisfied: numpy in
/usr/local/lib/python3.10/dist-packages (1.25.2)
Requirement already satisfied: pandas in
/usr/local/lib/python3.10/dist-packages (2.0.3)
Requirement already satisfied: scikit-learn in
/usr/local/lib/python3.10/dist-packages (1.2.2)
Requirement already satisfied: python-dateutil>=2.8.2 in
/usr/local/lib/python3.10/dist-packages (from pandas) (2.9.0.post0)
Requirement already satisfied: pytz>=2020.1 in
/usr/local/lib/python3.10/dist-packages (from pandas) (2024.1)
Requirement already satisfied: tzdata>=2022.1 in
/usr/local/lib/python3.10/dist-packages (from pandas) (2024.1)
Requirement already satisfied: scipy>=1.3.2 in
/usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.11.4)
Requirement already satisfied: joblib>=1.1.1 in
/usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.4.2)
Requirement already satisfied: threadpoolctl>=2.0.0 in
/usr/local/lib/python3.10/dist-packages (from scikit-learn) (3.5.0)
Requirement already satisfied: six>=1.5 in
/usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.8.2-
>pandas) (1.16.0)
```

LOADING DATA

```
# Loading the data set
```

```
raw = pd.read_csv('AMAZON REVIEWS.csv')
raw

{"type": "dataframe", "variable_name": "raw"}
```


DATA INSPECTION AND UNDERSTANDING

```
# Checking the data types and null values
```

```
raw.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 34660 entries, 0 to 34659
```

```
Data columns (total 21 columns):
```

#	Column	Non-Null Count	Dtype
0	id	34660 non-null	object
1	name	27900 non-null	object
2	asins	34658 non-null	object
3	brand	34660 non-null	object
4	categories	34660 non-null	object
5	keys	34660 non-null	object
6	manufacturer	34660 non-null	object
7	reviews.date	34621 non-null	object
8	reviews.dateAdded	24039 non-null	object
9	reviews.dateSeen	34660 non-null	object
10	reviews.didPurchase	1 non-null	object
11	reviews.doRecommend	34066 non-null	object
12	reviews.id	1 non-null	float64
13	reviews.numHelpful	34131 non-null	float64
14	reviews.rating	34627 non-null	float64
15	reviews.sourceURLs	34660 non-null	object
16	reviews.text	34659 non-null	object
17	reviews.title	34654 non-null	object
18	reviews.userCity	0 non-null	float64
19	reviews.userProvince	0 non-null	float64
20	reviews.username	34653 non-null	object

```
dtypes: float64(5), object(16)
```

```
memory usage: 5.6+ MB
```

Columns with 0 Non-Null Count

- This column has 0 non-null entries, meaning all 34,660 entries are missing or null.
- This column does not contain any useful data.

Columns with 1 Non-Null Count

- This column has only 1 non-null entry, meaning out of 34,660 rows, only one entry has a value and the rest are null.
- This column contains almost no useful data.

```
# Checking the data shape
```

```
raw.shape
```

```
(34660, 21)
```

```
#Summary statistics
```

```
raw.describe()
```

```
{"summary": "{\n  \"name\": \"raw\",\n  \"rows\": 8,\n  \"fields\": [\n    {\n      \"column\": \"reviews.id\",\n      \"properties\": {\n        \"dtype\": \"number\",\n        \"std\": 42094956.36805944,\n        \"min\": 1.0,\n        \"max\": 111372787.0,\n        \"num_unique_values\": 2,\n        \"samples\": [\n          111372787.0,\n          1.0\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\"\n      },\n      \"column\": \"reviews.numHelpful\",\n      \"properties\": {\n        \"dtype\": \"number\",\n        \"std\": 12028.672992019248,\n        \"min\": 0.0,\n        \"max\": 34131.0,\n        \"num_unique_values\": 5,\n        \"samples\": [\n          0.6302481614954147,\n          814.0\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\"\n      },\n      \"column\": \"reviews.rating\",\n      \"properties\": {\n        \"dtype\": \"number\",\n        \"std\": 12241.214519762767,\n        \"min\": 0.735652907647782,\n        \"max\": 34627.0,\n        \"num_unique_values\": 6,\n        \"samples\": [\n          34627.0,\n          4.584572732260953\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\"\n      },\n      \"column\": \"reviews.userCity\",\n      \"properties\": {\n        \"dtype\": \"number\",\n        \"std\": null,\n        \"min\": 0.0,\n        \"max\": 0.0,\n        \"num_unique_values\": 1,\n        \"samples\": [\n          0.0\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\"\n      },\n      \"column\": \"reviews.userProvince\",\n      \"properties\": {\n        \"dtype\": \"number\",\n        \"std\": null,\n        \"min\": 0.0,\n        \"max\": 0.0,\n        \"num_unique_values\": 1,\n        \"samples\": [\n          0.0\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\"\n      }\n    ]\n  },\n  \"type\": \"dataframe\"}
```

```
# Previewing the columns
```

```
raw.columns
```

```
Index(['id', 'name', 'asins', 'brand', 'categories', 'keys',  
      'manufacturer',  
      'reviews.date', 'reviews.dateAdded', 'reviews.dateSeen',  
      'reviews.didPurchase', 'reviews.doRecommend', 'reviews.id',  
      'reviews.numHelpful', 'reviews.rating', 'reviews.sourceURLs',  
      'reviews.text', 'reviews.title', 'reviews.userCity',  
      'reviews.userProvince', 'reviews.username'],  
      dtype='object')
```

```
# Renaming the columns to standard naming convention
```

```
column_names = {\n    'id': 'id',  
    'name': 'product_name',
```



```

0    id                34660 non-null object
1    product_name      27900 non-null object
2    asins              34658 non-null object
3    brand              34660 non-null object
4    product_categories 34660 non-null object
5    product_keys       34660 non-null object
6    manufacturer_name  34660 non-null object
7    review_date        34621 non-null datetime64[ns, UTC]
8    review_date_added  24039 non-null object
9    review_date_seen   34660 non-null object
10   review_did_purchase 1 non-null object
11   review_do_recommend 34066 non-null object
12   review_id           1 non-null float64
13   review_num_helpful  34131 non-null float64
14   review_rating       34627 non-null float64
15   review_source_urls  34660 non-null object
16   review_text         34659 non-null object
17   review_title        34654 non-null object
18   review_user_city    0 non-null float64
19   review_user_province 0 non-null float64
20   review_username     34653 non-null object
dtypes: datetime64[ns, UTC](1), float64(5), object(15)
memory usage: 5.6+ MB

```

Checking for proportion of missing values

```
raw.isnull().mean()
```

```

id                0.000000
product_name      0.195038
asins              0.000058
brand              0.000000
product_categories 0.000000
product_keys       0.000000
manufacturer_name  0.000000
review_date        0.001125
review_date_added  0.306434
review_date_seen   0.000000
review_did_purchase 0.999971
review_do_recommend 0.017138
review_id          0.999971
review_num_helpful 0.015263
review_rating      0.000952
review_source_urls 0.000000
review_text        0.000029
review_title       0.000173
review_user_city   1.000000
review_user_province 1.000000
review_username    0.000202
dtype: float64

```

```
# Checking the missing values
```

```
raw.isnull().sum()
```

```
id                0
product_name      6760
asins             2
brand             0
product_categories 0
product_keys      0
manufacturer_name 0
review_date       39
review_date_added 10621
review_date_seen  0
review_did_purchase 34659
review_do_recommend 594
review_id         34659
review_num_helpful 529
review_rating     33
review_source_urls 0
review_text       1
review_title      6
review_user_city  34660
review_user_province 34660
review_username   7
dtype: int64
```

```
#check percentage of missing values
```

```
# create a function to check the percentage of missing values
```

```
def missing_values(row):
    miss = raw.isnull().sum().sort_values(ascending = False)
    percentage_miss = (raw.isnull().sum() /
len(raw)).sort_values(ascending = False)
    missing = pd.DataFrame({"Missing Values": miss, "Percentage":
percentage_miss}).reset_index()
    missing.drop(missing[missing["Percentage"] == 0].index, inplace =
True)
    return missing
```

```
missing_data = missing_values(raw)
```

```
missing_data
```

```
{"summary":{"name": "missing_data", "rows": 14,
"fields": [{"column": "index",
"properties": {"dtype": "string",
"num_unique_values": 14, "samples": [
"review_rating", "review_title",
"review_user_city", ], "semantic_type": "",
"description": ""}
}, {"column":
"Missing Values", "properties": {"dtype":
```

```

\"number\", \n          \"std\": 15685, \n          \"min\": 1, \n
\"max\": 34660, \n          \"num_unique_values\": 12, \n
\"samples\": [\n          2, \n          6, \n          34660\
n          ], \n          \"semantic_type\": \"\", \n
\"description\": \"\" \n          } \n          }, \n          { \n          \"column\":
\"Percentage\", \n          \"properties\": { \n          \"dtype\":
\"number\", \n          \"std\": 0.45255361191177174, \n          \"min\":
2.8851702250432774e-05, \n          \"max\": 1.0, \n
\"num_unique_values\": 12, \n          \"samples\": [\n
5.770340450086555e-05, \n          0.00017311021350259665, \n
1.0 \n          ], \n          \"semantic_type\": \"\", \n
\"description\": \"\" \n          } \n          } \n          ] \n
n}], \"type\": \"dataframe\", \"variable_name\": \"missing_data\"}

```

Checking for unique values in all columns

Loop through each column and print unique values

```

for column_name in raw.columns:
    unique_values = raw[column_name].unique()
    num_unique_values = len(unique_values)
    print(f\"Unique Values in '{column_name}' (Total:
{num_unique_values}):\")
    print(unique_values)
    print(\"\\n\" + \"=\"*50 + \"\\n\")

```

Unique Values in 'id' (Total: 42):

```

['AVqkIhwDv8e3D10-lebb' 'AVqVGZ03nnclJgDc3jGK' 'AVpe9CMS1cnluZ0-aoC5'
'AVpfBEWcilAPnD_xTgb7' 'AVqkIiKWnnclJgDc3khH' 'AVqkIj9snnc1JgDc3khU'
'AVsRjfwAU2_QcyX9PHqe' 'AVqVGZNvQMLgs0JE6eUY' 'AVpfs_CLJeJML43DH5w'
'AVphgVaX1cnluZ0-DR74' 'AVqVGZN9QMLgs0JE6eUZ' 'AVpftoij1cnluZ0-p5n2'
'AVqkIhxunnc1JgDc3kg_' 'AVpioXbb1cnluZ0-PImd' 'AVpff7_VilAPnD_xc1E_'
'AVpjEN4jLJeJML43rpUe' 'AVpg3q4RLJeJML43TxA_' 'AVqVGWLKnnclJgDc3jF1'
'AV1YnRtnglJLPUi8IJmV' 'AVphPmHuilAPnD_x3E5h' 'AVzvXXxbvKc47QAVfRhy'
'AVpe7AsMilAPnD_xQ78G' 'AVph0EeEilAPnD_x9myq' 'AVqkIdntQMLgs0JE6fuB'
'AVzRlorb-jtxr-f3ygvQ' 'AVqVGWQDv8e3D10-ldFr' 'AVzvXXwEvKc47QAVfRhX'
'AVpgdkC8ilAPnD_xsuyi' 'AV1YnR7wglJLPUi8IJmi' 'AVpfl8cLLJeJML43AE3S'
'AVqkEM34QMLgs0JE6e8q' 'AVzoGHhAglJLPUi8GfzY' 'AVpfIfGA1cnluZ0-emyp'
'AVphLY7v1cnluZ0-_Ty0' 'AVpf_4sUilAPnD_xlwYV' 'AVpidLjVilAPnD_xEVpI'
'AVpfpK8KLJeJML43BCuD' 'AVpe8PEvilAPnD_xRYIi' 'AV1YE_muvKc47QAVgpeE'
'AVpf_znpilAPnD_xlvAF' 'AVpggqsrLJeJML4305zp' 'AVpfiBlyLJeJML43-4Tp']

```

=====

Unique Values in 'product_name' (Total: 49):

```

['All-New Fire HD 8 Tablet, 8 HD Display, Wi-Fi, 16 GB - Includes
Special Offers, Magenta'
'Kindle Oasis E-reader with Leather Charging Cover - Merlot, 6 High-
Resolution Display (300 ppi), Wi-Fi - Includes Special Offers,, '
'Amazon Kindle Lighted Leather Cover,,,\\r\\nAmazon Kindle Lighted

```

Leather Cover,,, '
'Amazon Kindle Lighted Leather Cover,,, \r\nKindle Keyboard,,, '
'Kindle Keyboard,,, \r\nKindle Keyboard,,, '
'All-New Fire HD 8 Tablet, 8 HD Display, Wi-Fi, 32 GB - Includes Special Offers, Magenta'
'Fire HD 8 Tablet with Alexa, 8 HD Display, 32 GB, Tangerine - with Special Offers, '
'Amazon 5W USB Official OEM Charger and Power Adapter for Fire Tablets and Kindle eReaders,,, \r\nAmazon 5W USB Official OEM Charger and Power Adapter for Fire Tablets and Kindle eReaders,,, '
'All-New Kindle E-reader - Black, 6 Glare-Free Touchscreen Display, Wi-Fi - Includes Special Offers,,, '
'Amazon Kindle Fire Hd (3rd Generation) 8gb,,, \r\nAmazon Kindle Fire Hd (3rd Generation) 8gb,,, '
'Fire Tablet, 7 Display, Wi-Fi, 8 GB - Includes Special Offers, Magenta'
'Kindle Oasis E-reader with Leather Charging Cover - Black, 6 High-Resolution Display (300 ppi), Wi-Fi - Includes Special Offers,,, '
'Amazon - Kindle Voyage - 4GB - Wi-Fi + 3G - Black,,, \r\nAmazon - Kindle Voyage - 4GB - Wi-Fi + 3G - Black,,, '
'Amazon - Kindle Voyage - 4GB - Wi-Fi + 3G - Black,,, \r\nFire HD 8 Tablet with Alexa, 8 HD Display, 16 GB, Tangerine - with Special Offers", '
'Fire HD 8 Tablet with Alexa, 8 HD Display, 16 GB, Tangerine - with Special Offers, '
'Amazon Standing Protective Case for Fire HD 6 (4th Generation) - Black,,, \r\nAmazon Standing Protective Case for Fire HD 6 (4th Generation) - Black,,, '
'Certified Refurbished Amazon Fire TV (Previous Generation - 1st),,,, \r\nCertified Refurbished Amazon Fire TV (Previous Generation - 1st),,,, '
'Brand New Amazon Kindle Fire 16gb 7 Ips Display Tablet Wifi 16 Gb Blue,,, '
'Amazon Kindle Touch Leather Case (4th Generation - 2011 Release), Olive Green,,, \r\nAmazon Kindle Touch Leather Case (4th Generation - 2011 Release), Olive Green,,, '
'Fire Kids Edition Tablet, 7 Display, Wi-Fi, 16 GB, Green Kid-Proof Case'
'Amazon Kindle Paperwhite - eBook reader - 4 GB - 6 monochrome Paperwhite - touchscreen - Wi-Fi - black,,, '
'Kindle Voyage E-reader, 6 High-Resolution Display (300 ppi) with Adaptive Built-in Light, PagePress Sensors, Wi-Fi - Includes Special Offers, '
'Certified Refurbished Amazon Fire TV Stick (Previous Generation - 1st),,,, \r\nCertified Refurbished Amazon Fire TV Stick (Previous Generation - 1st),,,, '
'Certified Refurbished Amazon Fire TV Stick (Previous Generation - 1st),,,, \r\nKindle Paperwhite,,, '
'Kindle Paperwhite,,, \r\nKindle Paperwhite,,, '

'Amazon Fire Kids Edition Tablet, 7 Display, Wi-Fi, 16 GB, Blue Kid-Proof Case - Blue'

'Kindle Paperwhite E-reader - White, 6 High-Resolution Display (300 ppi) with Built-in Light, Wi-Fi - Includes Special Offers,,,'

'Amazon Echo and Fire TV Power Adapter,,,\\r\\nAmazon Echo and Fire TV Power Adapter,,,'

'Amazon Fire Hd 8 8in Tablet 16gb Black B018s3t3bk 6th Gen (2016) Android,,,\\r\\nAmazon Fire Hd 8 8in Tablet 16gb Black B018s3t3bk 6th Gen (2016) Android,,,'

'Certified Refurbished Amazon Fire TV with Alexa Voice Remote,,,\\r\\nCertified Refurbished Amazon Fire TV with Alexa Voice Remote,,,'

'Amazon - Fire 16GB (5th Gen, 2015 Release) - Black,,,\\r\\nAmazon - Fire 16GB (5th Gen, 2015 Release) - Black,,,'

'Fire Tablet, 7 Display, Wi-Fi, 8 GB - Includes Special Offers, Black'

'Echo (White),,,,\\r\\nEcho (White),,,,'

'Echo (White),,,,\\r\\nFire Tablet, 7 Display, Wi-Fi, 8 GB - Includes Special Offers, Tangerine"'

'Echo (Black),,,,\\r\\nEcho (Black),,,,'

'Echo (Black),,,,\\r\\nAmazon 9W PowerFast Official OEM USB Charger and Power Adapter for Fire Tablets and Kindle eReaders,,,'

'Amazon 9W PowerFast Official OEM USB Charger and Power Adapter for Fire Tablets and Kindle eReaders,,,\\r\\nAmazon 9W PowerFast Official OEM USB Charger and Power Adapter for Fire Tablets and Kindle eReaders,,,'

'Amazon Fire Hd 6 Standing Protective Case(4th Generation - 2014 Release), Cayenne Red,,,\\r\\nAmazon Fire Hd 6 Standing Protective Case(4th Generation - 2014 Release), Cayenne Red,,,'

'Amazon Fire Hd 6 Standing Protective Case(4th Generation - 2014 Release), Cayenne Red,,,\\r\\nAmazon 5W USB Official OEM Charger and Power Adapter for Fire Tablets and Kindle eReaders,,,'

'Amazon Fire Hd 10 Tablet, Wi-Fi, 16 Gb, Special Offers - Silver Aluminum,,,\\r\\nAmazon Fire Hd 10 Tablet, Wi-Fi, 16 Gb, Special Offers - Silver Aluminum,,,'

'Amazon - Amazon Tap Portable Bluetooth and Wi-Fi Speaker - Black,,,\\r\\nAmazon - Amazon Tap Portable Bluetooth and Wi-Fi Speaker - Black,,,'

'Coconut Water Red Tea 16.5 Oz (pack of 12),,,,\\r\\nAmazon Fire Tv,,,'

'Amazon Fire Tv,,,\\r\\nAmazon Fire Tv,,,'

'Amazon Fire Tv,,,\\r\\nKindle Dx Leather Cover, Black (fits 9.7 Display, Latest and 2nd Generation Kindle Dxs)",,,,'

'Kindle Dx Leather Cover, Black (fits 9.7 Display, Latest and 2nd Generation Kindle Dxs),,,,'

'Amazon Kindle Fire Hd 9w Powerfast Adapter Charger + Micro Usb Angle Cable,,,\\r\\nAmazon Kindle Fire Hd 9w Powerfast Adapter Charger + Micro Usb Angle Cable,,,'

'New Amazon Kindle Fire Hd 9w Powerfast Adapter Charger + Micro Usb Angle Cable,,,'

'New Amazon Kindle Fire Hd 9w Powerfast Adapter Charger + Micro Usb
Angle Cable,,,\r\n'
nan]

=====

Unique Values in 'asins' (Total: 42):

['B01AHB9CN2' 'B00VINDBJK' 'B005PB2T0S' 'B002Y27P3M' 'B01AHB9CYG'
'B01AHB9C1E' 'B01J2G4VBG' 'B00ZV9PXP2' 'B0083Q04TA' 'B018Y2290U'
'B00REQKWA' 'B00IOYAM4I' 'B018T075DC' nan 'B00DU15MU4' 'B018Y225IA'
'B005PB2T2Q' 'B018Y23MNM' 'B000QVZDJM' 'B00IOY8XWQ' 'B00L029KXQ'
'B00QJDU3KY' 'B018Y22C2Y' 'B01BFIBRIE' 'B01J40RNHU' 'B018SZT3BK'
'B00UH4D8G2' 'B018Y22BI4' 'B00TSUGXKE' 'B00L9EPT80,B01E6A069U'
'B018Y23P7K' 'B00X4WHP5E' 'B00QFQRELG' 'B00LW9X0JM' 'B00QL1ZN3G'
'B0189XY0Q' 'B01BH8300M' 'B00BFJAHF8' 'B00U3FPN4U' 'B002Y27P6Y'
'B006GW05NE' 'B006GW05WK']

=====

Unique Values in 'brand' (Total: 6):

['Amazon' 'Amazon Fire' 'Amazon Echo' 'Amazon Coco T' 'Amazon Fire Tv'
'Amazon Digital Services Inc.']

=====

Unique Values in 'product_categories' (Total: 41):

['Electronics,iPad & Tablets,All Tablets,Fire
Tablets,Tablets,Computers & Tablets'
'eBook Readers,Kindle E-readers,Computers & Tablets,E-Readers &
Accessories,E-Readers'
'Electronics,eBook Readers & Accessories,Covers,Kindle Store,Amazon
Device Accessories,Kindle E-Reader Accessories,Kindle (5th Generation)
Accessories,Kindle (5th Generation) Covers'
'Kindle Store,Amazon Devices,Electronics'
'Tablets,Fire Tablets,Electronics,Computers,Computer Components,Hard
Drives & Storage,Computers & Tablets,All Tablets'
'Tablets,Fire Tablets,Computers & Tablets,All Tablets'
'Amazon Devices & Accessories,Amazon Device Accessories,Power
Adapters & Cables,Kindle Store,Kindle E-Reader Accessories,Kindle
Paperwhite Accessories'
'Electronics,iPad & Tablets,All Tablets,Computers/Tablets &
Networking,Tablets & eBook Readers,Computers & Tablets,E-Readers &
Accessories,E-Readers,Used:Computers
Accessories,Used:Tablets,Computers,iPads Tablets,Kindle E-
readers,Electronics Features'
'Computers/Tablets & Networking,Tablets & eBook
Readers,Electronics,eBook Readers & Accessories,eBook Readers'
'Fire Tablets,Tablets,Computers & Tablets,All Tablets,Electronics,
Tech Toys, Movies, Music,Electronics,iPad & Tablets,Android
Tablets,Frys'

'Kindle E-readers,Electronics Features,Computers & Tablets,E-Readers & Accessories,E-Readers,eBook Readers'

'Computers & Tablets,E-Readers & Accessories,eBook Readers,Kindle E-readers'

'Fire Tablets,Tablets,Computers & Tablets,All Tablets'

'Frys,Software & Books,eReaders & Accessories,Tablet Cases Covers,Tablet Accessories,Computer Accessories'

'Electronics,Categories,Streaming Media Players,Amazon Devices'

'Computers/Tablets & Networking,Tablets & eBook Readers,Computers & Tablets,Tablets,All Tablets'

'Amazon Device Accessories,Kindle Store,Kindle Touch (4th Generation) Accessories,Kindle E-Reader Accessories,Covers,Kindle Touch (4th Generation) Covers'

'Walmart for Business,Office Electronics,Tablets,Office,Electronics,iPad & Tablets,Windows Tablets,All Windows Tablets,Computers & Tablets,E-Readers & Accessories,E-Readers,eBook Readers,Kindle E-readers,Computers/Tablets & Networking,Tablets & eBook Readers,Electronics Features,Books & Magazines,Book Accessories,eReaders,TVs & Electronics,Computers & Laptops,Tablets & eReaders'

'Walmart for Business,Office Electronics,Tablets,Electronics,iPad & Tablets,All Tablets,Computers & Tablets,E-Readers & Accessories,Kindle E-readers,Electronics Features,eBook Readers,See more Amazon Kindle Voyage (Wi-Fi),See more Amazon Kindle Voyage 4GB, Wi-Fi 3G (Unlocked...'

'Electronics,Categories,Fire TV,Kindle Store'

'amazon.co.uk,Amazon Devices'

"Electronics,Computers,Computer Accessories,Cases & Bags,Fire Tablets,Electronics Features,Tablets,Computers & Tablets,Kids' Tablets,Electronics, Tech Toys, Movies, Music,iPad & Tablets,Top Rated"

'Electronics,iPad & Tablets,All Tablets,Computers & Tablets,Tablets,eBook Readers'

'Kindle Store,Categories,eBook Readers & Accessories,Fire TV Accessories,Electronics,Power Adapters & Cables,Amazon Device Accessories,Power Adapters'

'Fire Tablets,Tablets,Computers & Tablets,All Tablets,Computers/Tablets & Networking,Tablets & eBook Readers'

'Categories,Streaming Media Players,Electronics'

'Computers & Tablets,Tablets,All Tablets,Computers/Tablets & Networking,Tablets & eBook Readers,Fire Tablets,Frys'

'Electronics Features,Fire Tablets,Computers & Tablets,Tablets,All Tablets,Computers/Tablets & Networking,Tablets & eBook Readers'

'Stereos,Remote Controls,Amazon Echo,Audio Docks & Mini Speakers,Amazon Echo Accessories,Kitchen & Dining Features,Speaker Systems,Electronics,TVs Entertainment,Clearance,Smart Hubs & Wireless Routers,Featured Brands,Wireless Speakers,Smart Home & Connected Living,Home Security,Kindle Store,Home Automation,Home, Garage & Office,Home,Voice-Enabled Smart Assistants,Virtual Assistant

Speakers, Portable Audio & Headphones, Electronics Features, Amazon Device Accessories, iPod, Audio Player Accessories, Home & Furniture Clearance, Consumer Electronics, Smart Home, Surveillance, Home Improvement, Smart Home & Home Automation Devices, Smart Hubs, Home Safety & Security, Voice Assistants, Alarms & Sensors, Amazon Devices, Audio, Holiday Shop'

'Fire Tablets, Tablets, Computers & Tablets, All Tablets, Frys'

'TVs Entertainment, Wireless Speakers, Virtual Assistant Speakers, Featured Brands, Electronics, Amazon Devices, Home, Home Improvement, Home Safety & Security, Home Security, Alarms & Sensors, Smart Home & Home Automation Devices, Smart Hubs & Wireless Routers, Smart Hubs, Consumer Electronics, Voice-Enabled Smart Assistants, Smart Home & Connected Living, Home, Garage & Office, Smart Home, Voice Assistants, Surveillance, Home Automation, Speakers, Electronics Features, Holiday Shop, TV, Video & Home Audio, Internet & Media Streamers, Amazon Echo, Hubs & Controllers'

'Chargers & Adapters, Computers & Accessories, Tablet & E-Reader Accessories, Amazon Devices & Accessories, Fire Tablet Accessories, Electronics, Power Adapters & Cables, Cell Phones, Amazon Device Accessories, Cell Phone Accessories, Cell Phone Batteries & Power, Tablet Accessories, Featured Brands, Kindle Fire (2nd Generation) Accessories, Kindle Store, Home Improvement, Fire (5th Generation) Accessories, Electrical, Amazon Devices, Home, Tablets & E-Readers, Cables & Chargers'

'Cases, Kindle Store, Amazon Device Accessories, Accessories, Tablet Accessories'

'Electronics, eBook Readers & Accessories, Power Adapters, Computers/Tablets & Networking, Tablet & eBook Reader Accs, Chargers & Sync Cables, Power Adapters & Cables, Kindle Store, Amazon Device Accessories, Kindle Fire (2nd Generation) Accessories, Fire Tablet Accessories'

'Electronics, Tablets & E-Readers, Tablets, Back To College, College Electronics, College iPads & Tablets, Featured Brands, Amazon Devices, Electronics Deals, Computers & Tablets, All Tablets, Electronics Features, eBook Readers'

'Featured Brands, Electronics, Amazon Devices, Home, Home Improvement, Home Safety & Security, Home Security, Alarms & Sensors, Smart Home & Home Automation Devices, Mobile, Mobile Speakers, Mobile Bluetooth Speakers, Smart Hubs & Wireless Routers, Smart Hubs, Home, Garage & Office, Smart Home, Voice Assistants, Smart Home & Connected Living, Amazon Tap, Portable Audio, MP3 Accessories, Speakers, Amazon Echo, Electronics Features, TVs & Electronics, Portable Audio & Electronics, MP3 Player Accessories, Home Theater & Audio, Kindle Store, Frys, Electronic Components, Home Automation, Electronics, Tech Toys, Movies, Music, Audio, Bluetooth Speakers'

'Rice Dishes, Ready Meals, Beauty, Moisturizers, Lotions'

'Back To College, College Electronics, College TVs & Home Theater, Electronics, TVs & Home Theater, Streaming Devices, Featured

Brands,Amazon Devices,Holiday Shop,Ways To Shop,TV & Home Theater,Streaming Media Players,All Streaming Media Players,TVs Entertainment,Video Games,Kindle Store,Electronics Features,Kids & Family,Fire TV'

'Electronics,Amazon Device Accessories,Kindle Store,Covers,Kindle E-Reader Accessories,Kindle DX (2nd Generation, Global Wireless) Accessories'

'Power Adapters & Cables,Electronics,USB Cables'

'Computers/Tablets & Networking,Tablet & eBook Reader Accs,Chargers & Sync Cables,Power Adapters & Cables,Kindle Store,Amazon Device Accessories,Fire Tablet Accessories,Kindle Fire (2nd Generation) Accessories']

=====

Unique Values in 'product_keys' (Total: 42):

['841667104676,amazon/53004484,amazon/b01ahb9cn2,0841667104676,allnewfirehd8tablet8hddisplaywifi16gbincludesspecialoffersmagenta/5620406,allnewfirehd8tablet8hddisplaywifi16gbincludesspecialoffersmagenta/b01ahb9cn2'

'kindleoasisereaderwithleatherchargingcovermerlot6highresolutiondisplay300ppiwifiincludesspecialoffers/5234468,amazon/b00vindbjk,kindleoasisereaderwithleatherchargingcovermerlot6highresolutiondisplay300ppiwifiincludesspecialoffers/b00vindbjk,848719069587,0848719069587'

'amazonkindlelightedleathercover/b005pb2t0s'

'kindlekeyboard/b002y27p3m,amazon/d01101'

'841667104690,allnewfirehd8tablet8hddisplaywifi32gbincludesspecialoffersmagenta/5620408,0841667104690,allnewfirehd8tablet8hddisplaywifi32gbincludesspecialoffersmagenta/b01ahb9cyg,amazon/53004761'

'amazon/b01ahb9cle,0841667104577,firehd8tabletwithalexa8hddisplay32gbtangerinewithspecialoffers/b01ahb9cle,firehd8tabletwithalexa8hddisplay32gbtangerinewithspecialoffers/5620411,841667104577'

'0841667120171,841667120171,amazon5wusbofficialoemchargerpoweradapterforfiretabletskindleereaders/b01j2g4vbg'

'allnewkindleereaderblack6glarefreetouchscreendisplaywifiincludesspecialoffers/391843532825,allnewkindleereaderblack6glarefreetouchscreendisplaywifiincludesspecialoffers/b00zv9pxp2,0848719083774,allnewkindleereaderblack6glarefreetouchscreendisplaywifiincludesspecialoffers/252974470193,amazon/b00zv9pxp2,848719083774,allnewkindleereaderblack6glarefreetouchscreendisplaywifiincludesspecialoffers/

322538285013,allnewkindleereaderblack6glarefreetouchscreendisplaywifii
ncludesspecialoffers/
5442403,allnewkindleereaderblack6glarefreetouchscreendisplaywifiiinclu
desspecialoffers/
kier2016bk,allnewkindleereaderblack6glarefreetouchscreendisplaywifiiinc
ludesspecialoffers/
162691587356,allnewkindleereaderblack6glarefreetouchscreendisplaywifii
ncludesspecialoffers/1631053'

'amazon/53000386,amazonkindlefirehd3rdgeneration8gb/122605594245,amazo
nkindlefirehd3rdgeneration8gb/
152615237936,amazonkindlefirehd3rdgeneration8gb/
391871762463,amazonkindlefirehd3rdgeneration8gb/b0083q04ta'

'firetablet7displaywifi8gbincludesspecialoffersmagenta/5025800,8416671
03105,0841667103105,amazon/
b018y229ou,firetablet7displaywifi8gbincludesspecialoffersmagenta/
b018y229ou'

'0848719057331,kindleoasisereaderwithleatherchargingcoverblack6highres
olutiondisplay300ppiwifiincludesspecialoffers/b00reqkwga,amazon/
b00reqkwga,kindleoasisereaderwithleatherchargingcoverblack6highresolut
iondisplay300ppiwifiincludesspecialoffers/5195001,848719057331'

'amazonkindlevoyage4gbwifi3gblack/9301112,amazon/b00ioyam4i,0848719040
098,848719040098,amazonkindlevoyage4gbwifi3gblack/b00ioyam4i'

'amazon/b018t075dc,firehd8tabletwithalexa8hddisplay16gbtangerinewithsp
ecialoffers/
5620410,firehd8tabletwithalexa8hddisplay16gbtangerinewithspecialoffers
/b018t075dc,841667103068,0841667103068'

'848719047530,amazonstandingprotectivecaseforfirehd64thgenerationblack
/3610684,amazonstandingprotectivecaseforfirehd64thgenerationblack/
018w006857385001p,amazon/
b00kqe2qaw,amazonstandingprotectivecaseforfirehd64thgenerationblack/
018w006857385001'

'848719035551,0848719035551,certifiedrefurbishedamazonfiretvpreviousge
neration1st/b00dul5mu4'

'841667103143,0841667103143,brandnewamazonkindlefire16gb7ipsdisplaytab
letwifil6gbbblue/
5025500,brandnewamazonkindlefire16gb7ipsdisplaytabletwithwifil6gbbblue/
b018y225ia,brandnewamazonkindlefire16gb7ipsdisplaytabletwithwifil6gbbblue/
201625338826,brandnewamazonkindlefire16gb7ipsdisplaytabletwithwifil6gbbblue
/362123960192,amazon/b018y225ia'

'amazonkindletouchleathercase4thgeneration2011releaseolivegreen/b005pb
2t2q'

'firekidseditiontablet7displaywifil6gbgreenkidproofcase/b018y23mm,841667103402,0841667103402,firekidseditiontablet7displaywifil6gbgreenkidproofcase/5026300,amazon/b018y23mm'

'amazon/b00oqvzdm,848719056099,amazonkindlepaperrwhiteebookreader4gb6monochrome paperrwhite touchscreenwifiblack/
263087494445,amazonkindlepaperrwhiteebookreader4gb6monochrome paperrwhite touchscreenwifiblack/
9439005,amazonkindlepaperrwhiteebookreader4gb6monochrome paperrwhite touchscreenwifiblack/
b00oqvzdm,0848719056099,amazonkindlepaperrwhiteebookreader4gb6monochrome paperrwhite touchscreenwifiblack/00355266000p'

'848719040104,kindlevoyageereader6highresolutiondisplay300ppiwithadaptivebuiltinlightpagepressensorswifiincludesspecialoffers/
b00ioy8xwq,0848719040104,kindlevoyageereader6highresolutiondisplay300ppiwithadaptivebuiltinlightpagepressensorswifiincludesspecialoffers/
321689278417,kindlevoyageereader6highresolutiondisplay300ppiwithadaptivebuiltinlightpagepressensorswifiincludesspecialoffers/
9302088,amazon/53002680'

'certifiedrefurbishedamazonfiretvstickpreviousgeneration1st/b00lo29kxq,0848719052121,848719052121'
'kindlepaperrwhite/b00qjdu3ky'

'amazon/b018y22c2y,841667103389,0841667103389,firekidseditiontablet7displaywifil6gbbluekidproofcase/
b018y22c2y,amazonfirekidsedition16gb5thgen2015releaseblue/
5026000,amazonfirekidsedition7tablet16gbblue/
5026000,amazonkidsedition7inch16gbfiretabletblue/kifk716cblu'

'841667107868,amazon/53004915,amazonkindlepaperrwhitewhite/5435104,0841667107868,kindlepaperrwhiteereaderwhite6highresolutiondisplay300ppiwithbuiltinlightwifiincludesspecialoffers/b01bfibrie'
'amazonechofiretvpoweradapter/b01j4ornhu,0841667120829,841667120829'

'amazonfirehd88intablet16gbblackb018szt3bk6thgen2016android/5538501,amazonfirehd88intablet16gbblackb018szt3bk6thgen2016android/
b018szt3bk,amazonfirehd88intablet16gbblackb018szt3bk6thgen2016android/
182378029308,amazonfirehd88intablet16gbblackb018szt3bk6thgen2016android/
322430145717,841667103037,0841667103037,amazonfirehd88intablet16gbblackb018szt3bk6thgen2016android/152627691815,amazon/
b018szt3bk,amazonfirehd88intablet16gbblackb018szt3bk6thgen2016android/
332403091354,amazonfirehd88intablet16gbblackb018szt3bk6thgen2016android/322598029639'

'certifiedrefurbishedamazonfiretvwithalexa voiceremote/b00uh4d8g2,0848719063264,848719063264'

'amazonfire16gb5thgen2015releaseblack/272201222631,amazonfire16gb5thgen2015releaseblack/
b018y22bi4,841667103129,0841667103129,amazonfire16gb5thgen2015releaseblack/5023200,amazonfire16gb5thgen2015releaseblack/
332273296844,amazonfire16gb5thgen2015releaseblack/
232443003172,amazon/b018y22bi4'

'amazon/b00tsugxke,0848719062854,firetablet7displaywifi8gbincludesspecialoffersblack/
b00tsugxke,848719062854,firetablet7displaywifi8gbincludesspecialoffersblack/4390200,firetablet7displaywifi8gbincludesspecialoffersblack/
322581680105'

'echowwhite/263039693056,echowwhite/152558276095,echowwhite/292178880467,echowwhite/222588935706,echowwhite/253120140398,echowwhite/
322577436254,echowwhite/122597356284,echowwhite/132263972952,echowwhite/322586415668,echowwhite/152626395386,echowwhite/272724680159,echowwhite/222587602421,echowwhite/122474318097,echowwhite/5588528,echowwhite/
112567699636,echowwhite/272768463386,echowwhite/332175902683,echowwhite/311908601694,echowwhite/292041139369,echowwhite/192239032596,echowwhite/272768869474,0841667112862,echowwhite/222507973621,echowwhite/
112391858963,echowwhite/291992370210,echowwhite/b00l9ept8o,echowwhite/112480241614,echowwhite/b01e6ao69u,echowwhite/322589755316,echowwhite/
322574315372,echowwhite/253051886606,echowwhite/382165760287,echowwhite/222582493180,echowwhite/282581384521,echowwhite/112479310908,echowwhite/302201691992,echowwhite/201761456849,echowwhite/amechow2k,echowwhite/
132262816901,echowwhite/282571823011,echowwhite/322511136772,841667112862,echowwhite/232407174148,echowwhite/
322441917397,echowwhite/amechow,echowwhite/332296207643,echowwhite/152610914446,echowwhite/222578584785,echowwhite/162591117080,echowwhite/
162593787621,echowwhite/232407374203,echowwhite/162595518416,echowwhite/152623638099,amazon/b01e6ao69u'

'firetablet7displaywifi8gbincludesspecialofferstangerine/b018y23p7k,amazon/
b018y23p7k,841667103112,0841667103112,firetablet7displaywifi8gbincludesspecialofferstangerine/5025600'

'echoblack/122590756021,echoblack/302383519724,echoblack/152610063234,echoblack/263085373650,echoblack/292186264538,echoblack/
142444819720,echoblack/4747312,echoblack/332049242998,echoblack/122605571555,echoblack/262723627567,echoblack/322574315387,echoblack/
272745799372,echoblack/222580642553,echoblack/152621788949,echoblack/amecho,echoblack/172781603329,echoblack/282572756222,echoblack/
122610899605,echoblack/162559901119,echoblack/142452148997,echoblack/332206230624,echoblack/b00x4whp5e,echoblack/112485743326,echoblack/
322375155959,echoblack/292170574038,echoblack/201960373389,echoblack/142444811424,echoblack/172628665232,echoblack/162593779001,echoblack/
362035886227,echoblack/232211327517,echoblack/172561946998,echoblack/

253042076236,echoblack/263078789473,echoblack/192224346233,echoblack/
272514028697,echoblack/142447784347,echoblack/222572724175,echoblack/
182666767685,echoblack/152624829595,echoblack/302379034663,echoblack/
332285103532,echoblack/321843010218,echoblack/162580754896,echoblack/
122610162473,echoblack/112480244746,echoblack/352113335165,echoblack/
263076769150,echoblack/253016522126,echoblack/122345243004,echoblack/
291843555681,echoblack/152632082666,echoblack/322443963965,echoblack/
282547976259,echoblack/152558240605,echoblack/302385706126,echoblack/
332303083809,echoblack/232402457168,echoblack/263088372235,echoblack/
122607176005,echoblack/282582537447,echoblack/112479564543,echoblack/
132180303866,echoblack/232415728345,echoblack/272770487395,echoblack/
332310066903,echoblack/112483045441,echoblack/232377282103,echoblack/
232234984279,amazon/b00x4whp5e,echoblack/162598705262,echoblack/
332043727777,echoblack/112474305384,echoblack/122610140759,echoblack/
152527501046,echoblack/152626398155,echoblack/263085892574,echoblack/
142443837098,echoblack/332387407823,echoblack/332308051158,echoblack/
152621786578,echoblack/152628376969,echoblack/
282578839976,848719071733,echoblack/152622010270,echoblack/
amechob2k,echoblack/282570718216,echoblack/282572809095,echoblack/
142443963349,echoblack/201986162871,echoblack/201990846674,echoblack/
322589784359,echoblack/152605359984,echoblack/192121594431,echoblack/
272758321435,echoblack/05740865000p,echoblack/222515042938,echoblack/
152710137365,0848719071733,echoblack/332119861214,echoblack/
302390440897,echoblack/253045357845,echoblack/142364243234'

'amazon/51752067,amazon9wpowerfastofficialoemusbchargerpoweradapterfor
firetabletskindleereaders/
b00qfqrelg,848719056556,0848719056556,amazon9wpowerfastofficialoemusb
chargerpoweradapterforfiretabletskindleereaders/4467600'

'amazonfirehd6standingprotectivecase4thgeneration2014releasecayennered
/b00lw9xojm'

'amazon5wusbofficialoemchargerpoweradapterforfiretabletskindleereaders
/
b00ql1zn3g,848719056532,0848719056532,amazon5wusbofficialoemchargerpow
eradapterforfiretabletskindleereaders/272684045946,amazon/55000660'

'841667101743,amazonfire/51441641,amazonfirehd1016gb5thgen2015releases
ilveraluminum/
5386601,0841667101743,firehd10tablet101hddisplaywifi16gbincludesspecia
lofferssilveraluminum/b0189xyy0q,amazon/
51441641,amazonfirehd10101tablet16gbsilveraluminum/5386601'

'amazonecho/b01bh83oom,amazon/b01bh83oom,amazonamazontapportableblueto
othwifispeakerblack/
5097300,amazonamazontapportablebluetoothwifispeakerblack/
1001803403,841667107929,0841667107929,amazonamazontapportablebluetooth
wifispeakerblack/
05743627000p,amazonamazontapportablebluetoothwifispeakerblack/


```
amtap,amazonamazontapportablebluetoothwifispeakerblack/b01bh83oom'

'amazoncocot/spk73188,coconutwaterredteal65ozpackof12/b00bfjahf8,amazon/spk73188'

'848719057492,amazonfiretv/51454342,amazonfiretv2015modelblack/4370400,amazon/b00u3fnp4u,amazonfiretvstreamingmediaplayer/b00u3fnp4u,0848719057492,firetvstreamingmediaplayer2015model/amftv2,firetv/05740864000p'
'kindledxleathercoverblackfits97displaylatest2ndgenerationkindledxs/b002y27p6y'
'amazondigitalservices/53000407,0814916017379,814916017379,amazon/53000407,amazonkindlefire5ftusbtoemicrousb cableworkswithmostmicrousbtablets/b006gwo5ne'
'newamazonkindlefirehd9wpowerfastadapterchargeremicrousbanglecable/272582562733,amazon/53000136,newamazonkindlefirehd9wpowerfastadapterchargeremicrousbanglecable/b006gwo5wk,amazondigitalservices/53000136,639767206372,0639767206372']
```

```
=====

Unique Values in 'manufacturer_name' (Total: 2):
['Amazon' 'Amazon Digital Services, Inc']

=====
```

```
Unique Values in 'review_date' (Total: 1055):
<DatetimeArray>
['2017-01-13 00:00:00+00:00', '2017-01-12 00:00:00+00:00',
 '2017-01-23 00:00:00+00:00', '2017-01-24 00:00:00+00:00',
 '2017-01-27 00:00:00+00:00', '2017-02-03 00:00:00+00:00',
 '2017-02-06 00:00:00+00:00', '2017-02-05 00:00:00+00:00',
 '2017-03-20 00:00:00+00:00', '2017-03-19 00:00:00+00:00',
 ...
 '2013-02-11 00:00:00+00:00', '2017-12-03 00:00:00+00:00',
 '2012-11-13 00:00:00+00:00', '2012-11-02 00:00:00+00:00',
 '2012-10-16 00:00:00+00:00', '2012-09-18 00:00:00+00:00',
 '2012-11-21 00:00:00+00:00', '2012-10-19 00:00:00+00:00',
 '2012-10-31 00:00:00+00:00', '2012-12-23 00:00:00+00:00']
Length: 1055, dtype: datetime64[ns, UTC]

=====
```

```
Unique Values in 'review_date_added' (Total: 1942):
['2017-07-03T23:33:15Z' '2017-07-03T23:28:24Z' '2017-07-03T23:27:54Z' ...
 '2017-08-29T16:58:30Z' '2017-08-13T08:15:30Z' '2017-07-18T23:57:10Z']

=====
```

```
Unique Values in 'review_date_seen' (Total: 3911):
['2017-06-07T09:04:00.000Z,2017-04-30T00:45:00.000Z'
 '2017-06-07T09:04:00.000Z,2017-04-30T00:44:00.000Z'
 '2017-06-07T09:04:00.000Z,2017-04-30T00:42:00.000Z' ...
 '2015-09-02T00:00:00Z' '2015-09-04T00:00:00Z' '2015-09-01T00:00:00Z']
```

```
=====
Unique Values in 'review_did_purchase' (Total: 2):
[nan True]
```

```
=====
Unique Values in 'review_do_recommend' (Total: 3):
[True False nan]
```

```
=====
Unique Values in 'review_id' (Total: 2):
[          nan 1.11372787e+08]
```

```
=====
Unique Values in 'review_num_helpful' (Total: 98):
[ 0.  1.  2.  3. 55.  4. 24. 11. 42. 62.  7.  8.  6. 10.
 36. 16. 15. 13.  5. 271. 730. 221. 53. nan  9. 105. 19. 25.
 21. 14. 20. 22. 12. 96. 102. 34. 17. 73. 109. 27. 39. 57.
 18. 40. 33. 112. 355. 60. 263. 37. 28. 103. 26. 32. 43. 64.
 23. 650. 780. 740. 139. 126. 69. 75. 48. 292. 144. 93. 49. 95.
 31. 63. 204. 270. 82. 174. 98. 84. 629. 163. 422. 261. 185. 205.
 132. 170. 814. 434. 302. 54. 30. 46. 660. 195. 744. 384. 238.
217.]
```

```
=====
Unique Values in 'review_rating' (Total: 6):
[ 5.  4.  2.  1.  3. nan]
```

```
=====
Unique Values in 'review_source_urls' (Total: 11929):
['http://reviews.bestbuy.com/3545/5620406/reviews.htm?
format=embedded&page=200,http://reviews.bestbuy.com/3545/5620406/
reviews.htm?format=embedded&page=166'
 'http://reviews.bestbuy.com/3545/5620406/reviews.htm?
format=embedded&page=200,http://reviews.bestbuy.com/3545/5620406/
reviews.htm?format=embedded&page=167'
 'http://reviews.bestbuy.com/3545/5620406/reviews.htm?
format=embedded&page=154,http://reviews.bestbuy.com/3545/5620406/
```

reviews.htm?format=embedded&page=120'

...

'http://www.amazon.com/Amazon-Kindle-Micro-USB-Cable-Tablets/dp/B006GW05NE'

'https://www.ebay.com/itm/NEW-Amazon-Kindle-Fire-HD-9W-Powerfast-Adapter-Charger-Micro-USB-Angle-Cable/272582562733'

'http://www.amazon.com/Amazon-PowerFast-Adapter-Accelerated-Charging/dp/B006GW05WK']

=====

Unique Values in 'review_text' (Total: 34660):

['This product so far has not disappointed. My children love to use it and I like the ability to monitor control what content they see with ease.'

'great for beginner or experienced person. Bought as a gift and she loves it'

'Inexpensive tablet for him to use and learn on, step up from the NABI. He was thrilled with it, learn how to Skype on it already...'

...

"Love my Kindle Fire but I am really disappointed in the Kindle Power Fast Charging Unit. I've had it two months and I've used it many times - The first two times it worked okay but failed on the third and many subsequent tries. I've disposed of it and use my wife's iPad Nano charger which always works just fine."

"I was surprised to find it did not come with any type of charging cords so I had to purchase one and then found my Sprint HTC 3D charger is faster. I would not purchase again- 1st item I've ever not liked I've purchased from Amazon"

"to spite the fact that i have nothing but good things to say about amazon and anthing i've ever gotten from them. and that i love my fire. i find it greedy that the wall charger doesn't come with the kindle. not everyone, ok most people, but still not everyone has a usb port to plug into. i'm taking my charger back. i think amazon should make things right and let anyone who purchased a kindle without a charger have one for free, or credit those who had to buy one."]

=====

Unique Values in 'review_title' (Total: 19767):

['Kindle' 'very fast' 'Beginner tablet for our 9 year old son.' ...

'Should be included' 'Disappointing Charger' 'as with everyone else']

=====

Unique Values in 'review_user_city' (Total: 1):

[nan]

=====

```
Unique Values in 'review_user_province' (Total: 1):  
[nan]
```

```
=====
```

```
Unique Values in 'review_username' (Total: 26789):  
['Adapter' 'truman' 'DaveZ' ... 'Jonathan Stewart' 'J Lawson'  
 'Just the Buyer']
```

```
=====
```

DATA CLEANING

Handling Missing values

```
#drop all columns with high percentage of missing values and columns  
not needed
```

```
raw.drop(columns = ['review_date_added', 'review_date_seen',  
'review_did_purchase' , 'review_user_city',  
'review_user_province','review_id' , 'product_name' ,  
'review_source_urls'], inplace = True)
```

```
# drop rows with missing values  
raw.dropna(inplace = True)
```

```
# Verify that there are no more missing values  
print(raw.isnull().sum().sum()) # Should print 0
```

```
# Get the shape of the cleaned data  
print(raw.shape)
```

```
# Display the first few rows of the cleaned data  
raw.head(2)
```

```
0  
(34054, 13)
```

```
{"summary":{"\n  \"name\": \"raw\",\n  \"rows\": 34054,\n  \"fields\":  
[\n    {\n      \"column\": \"id\",\n      \"properties\": {\n        \"dtype\": \"category\",\n        \"num_unique_values\": 24,\n        \"samples\": [\n          \"AVqkIhxunnc1JgDc3kg_\",\n          \"AVpgdkC8ilAPnD_xsvyi\",\n          \"AVqkIhwDv8e3D10-lebb\"\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\"\n      }\n    },\n    {\n      \"column\": \"asins\",\n      \"properties\": {\n        \"dtype\": \"category\",\n        \"num_unique_values\":  
24,\n        \"samples\": [\n          \"B018T075DC\",\n          \"B018Y22BI4\",\n          \"B01AHB9CN2\"\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\"\n      }\n    }\n  ]\n}}
```

```

n    },\n    {\n        \"column\": \"brand\", \n        \"properties\": {\n            \"dtype\": \"category\", \n            \"num_unique_values\": 4, \n            \"samples\": [\n                \"Amazon Fire\", \n                \"Amazon Fire Tv\", \n                \"Amazon\", \n            ], \n            \"semantic_type\": \"\", \n            \"description\": \"\", \n        }, \n    {\n        \"column\": \"product_categories\", \n        \"properties\": {\n            \"dtype\": \"category\", \n            \"num_unique_values\": 23, \n            \"samples\": [\n                \"Computers & Tablets, Tablets, All Tablets, Computers/Tablets & Networking, Tablets & eBook Readers, Fire Tablets, Frys\", \n                \"Computers/Tablets & Networking, Tablets & eBook Readers, Computers & Tablets, Tablets, All Tablets\", \n                \"Electronics, iPad & Tablets, All Tablets, Fire Tablets, Tablets, Computers & Tablets\", \n            ], \n            \"semantic_type\": \"\", \n            \"description\": \"\", \n        }, \n    {\n        \"column\": \"product_keys\", \n        \"properties\": {\n            \"dtype\": \"category\", \n            \"num_unique_values\": 24, \n            \"samples\": [\n                \"amazon/b018t075dc, firehd8tabletwithalexa8hddisplay16gbtangerinewithspecialoffers/5620410, firehd8tabletwithalexa8hddisplay16gbtangerinewithspecialoffers/b018t075dc, 841667103068, 0841667103068\", \n                \"amazonfire16gb5thgen2015releaseblack/272201222631, amazonfire16gb5thgen2015releaseblack/b018y22bi4, 841667103129, 0841667103129, amazonfire16gb5thgen2015releaseblack/5023200, amazonfire16gb5thgen2015releaseblack/332273296844, amazonfire16gb5thgen2015releaseblack/232443003172, amazon/b018y22bi4\", \n                \"841667104676, amazon/53004484, amazon/b01ahb9cn2, 0841667104676, allnewfirehd8tablet8hddisplaywifil6gbincludesspecialoffersmagenta/5620406, allnewfirehd8tablet8hddisplaywifil6gbincludesspecialoffersmagenta/b01ahb9cn2\", \n            ], \n            \"semantic_type\": \"\", \n            \"description\": \"\", \n        }, \n    {\n        \"column\": \"\", \n        \"properties\": {\n            \"dtype\": \"\", \n            \"manufacturer_name\": \"\", \n            \"category\": \"\", \n            \"num_unique_values\": 1, \n            \"samples\": [\n                \"Amazon\", \n            ], \n            \"semantic_type\": \"\", \n            \"description\": \"\", \n        }, \n    {\n        \"column\": \"review_date\", \n        \"properties\": {\n            \"dtype\": \"date\", \n            \"min\": \"2014-10-24 00:00:00+00:00\", \n            \"max\": \"2018-04-18 00:00:00+00:00\", \n            \"num_unique_values\": 941, \n            \"samples\": [\n                \"2015-11-29 00:00:00+00:00\", \n            ], \n            \"semantic_type\": \"\", \n            \"description\": \"\", \n        }, \n    {\n        \"column\": \"\", \n        \"properties\": {\n            \"dtype\": \"\", \n            \"review_do_recommend\": \"\", \n            \"category\": \"\", \n            \"num_unique_values\": 2, \n            \"samples\": [\n                false, \n            ], \n            \"semantic_type\": \"\", \n            \"description\": \"\", \n        }, \n    {\n        \"column\": \"\", \n        \"properties\": {\n            \"dtype\": \"\", \n            \"review_num_helpful\": \"\", \n            \"number\": \"\", \n            \"std\": 2.194084771528854, \n            \"min\": 0.0, \n            \"max\": 109.0, \n            \"num_unique_values\": 57, \n        }, \n    }

```

```

\"samples\": [\n          0.0\n        ],\n        \"semantic_type\": \n        \"\", \n        \"description\": \"\"\n      },\n      {\n        \"column\": \"review_rating\", \n        \"properties\": {\n          \"dtype\": \"number\", \n          \"std\": 0.7217255917862178, \n          \"min\": 1.0, \n          \"max\": 5.0, \n          \"num_unique_values\": 5, \n          \"samples\": [\n            4.0\n          ], \n          \"semantic_type\": \"\", \n          \"description\": \"\"\n        }, \n        {\n          \"column\": \"review_text\", \n          \"properties\": {\n            \"dtype\": \"string\", \n            \"num_unique_values\": 34054, \n            \"samples\": [\n              \"My kids love this product, as do I. Parental restrictions can be set and they know when they have to shut them off. Good battery life too.\"\n            ], \n            \"semantic_type\": \"\", \n            \"description\": \"\"\n          }, \n          {\n            \"column\": \"review_title\", \n            \"properties\": {\n              \"dtype\": \"string\", \n              \"num_unique_values\": 19448, \n              \"samples\": [\n                \"Money's Worth\"\n              ], \n              \"semantic_type\": \"\", \n              \"description\": \"\"\n            }, \n            {\n              \"column\": \"review_username\", \n              \"properties\": {\n                \"dtype\": \"string\", \n                \"num_unique_values\": 26309, \n                \"samples\": [\n                  \"RED3\"\n                ], \n                \"semantic_type\": \"\", \n                \"description\": \"\"\n              }, \n              ]\n            }, \n            \"type\": \"dataframe\", \"variable_name\": \"raw\"}

```

Checking for duplicates

```

# Checking duplicated rows
num_duplicated = raw.duplicated().sum()
print(f\"Number of duplicated rows: {num_duplicated}\")

Number of duplicated rows: 0

# Checking for duplicates using the 'CustomerId' column
raw[raw.duplicated(subset=[\"asins\"])]

{\"summary\": \"{\\n  \"name\": \"raw[raw\", \\n  \"rows\": 34030, \\n
  \"fields\": [\\n    {\\n      \"column\": \"id\", \\n      \"properties\":
    {\\n        \"dtype\": \"category\", \\n        \"num_unique_values\":
    23, \\n        \"samples\": [\\n          \"AVqVGWQDv8e3D10-ldFr\", \\n
          \"AVpjEN4jLJeJML43rpUe\", \\n          \"AVqkIhwDv8e3D10-lebb\"
        ], \\n        \"semantic_type\": \"\", \\n        \"description\": \"\"
      }, \\n      {\\n        \"column\": \"asins\", \\n        \"properties\":
    {\\n        \"dtype\": \"category\", \\n        \"num_unique_values\":
    23, \\n        \"samples\": [\\n          \"B018SZT3BK\", \\n
          \"B018Y225IA\", \\n          \"B01AHB9CN2\"
        ], \\n        \"semantic_type\": \"\", \\n        \"description\": \"\"
      }, \\n      {\\n        \"column\": \"brand\", \\n        \"properties\": {
      {\\n        \"dtype\": \"category\", \\n        \"num_unique_values\": 4, \\n
      \"samples\": [\\n        \"Amazon Fire\", \\n        \"Amazon Fire
      Tv\", \\n        \"Amazon\"
    ], \\n        \"semantic_type\":

```

```
\",\n      \"description\": \"\",\n    },\n    {\n      \"column\": \"product_categories\",,\n      \"properties\": {\n        \"dtype\": \"category\",,\n        \"num_unique_values\": 22,\n        \"samples\": [\n          \"Electronics,iPad & Tablets,All\nTablets,Fire Tablets,Tablets,Computers & Tablets\",,\n          \"Electronics,iPad & Tablets,All Tablets,Computers &\nTablets,Tablets,eBook Readers\",,\n          \"Fire\nTablets,Tablets,Computers & Tablets,All Tablets\",,\n        ],,\n        \"semantic_type\": \"\",,\n        \"description\": \"\",,\n      },,\n      {\n        \"column\": \"product_keys\",,\n        \"properties\": {\n          \"dtype\": \"category\",,\n          \"num_unique_values\": 23,\n          \"samples\": [\n            \"amazonfirehd88intablet16gbblackb018szt3bk6thgen2016android/5538501,a\nmazonfirehd88intablet16gbblackb018szt3bk6thgen2016android/\nb018szt3bk,amazonfirehd88intablet16gbblackb018szt3bk6thgen2016android/\n182378029308,amazonfirehd88intablet16gbblackb018szt3bk6thgen2016androi\nd/\n322430145717,841667103037,0841667103037,amazonfirehd88intablet16gbblac\nkb018szt3bk6thgen2016android/152627691815,amazon/\nb018szt3bk,amazonfirehd88intablet16gbblackb018szt3bk6thgen2016android/\n332403091354,amazonfirehd88intablet16gbblackb018szt3bk6thgen2016androi\nd/322598029639\",,\n            \"841667103143,0841667103143,brandnewamazonkindlefire16gb7ipsdisplayta\nbletwifi16gbbblue/\n5025500,brandnewamazonkindlefire16gb7ipsdisplaytabletewifi16gbbblue/\nb018y225ia,brandnewamazonkindlefire16gb7ipsdisplaytabletewifi16gbbblue/\n201625338826,brandnewamazonkindlefire16gb7ipsdisplaytabletewifi16gbbblue\n/362123960192,amazon/b018y225ia\",,\n            \"841667104676,amazon/53004484,amazon/b01ahb9cn2,0841667104676,allnewf\nirehd8tablet8hddisplaywifil6gbincludesspecialoffersmagenta/\n5620406,allnewfirehd8tablet8hddisplaywifil6gbincludesspecialoffersmage\nnta/b01ahb9cn2\",,\n            ],,\n            \"semantic_type\": \"\",,\n            \"description\": \"\",,\n          },,\n          {\n            \"column\":\n\"manufacturer_name\",,\n            \"properties\": {\n              \"dtype\":\n\"category\",,\n              \"num_unique_values\": 1,\n              \"samples\":\n[\n                \"Amazon\",,\n              ],,\n              \"semantic_type\": \"\",,\n              \"description\": \"\",,\n            },,\n            {\n              \"column\": \"review_date\",,\n              \"properties\": {\n                \"dtype\": \"date\",,\n                \"min\": \"2014-10-24 00:00:00+00:00\",,\n                \"max\": \"2018-04-18 00:00:00+00:00\",,\n                \"num_unique_values\": 939,\n                \"samples\": [\n                  \"2016-04-15 00:00:00+00:00\",,\n                ],,\n                \"semantic_type\": \"\",,\n                \"description\": \"\",,\n              },,\n              {\n                \"column\":\n\"review_do_recommend\",,\n                \"properties\": {\n                  \"dtype\":\n\"category\",,\n                  \"num_unique_values\": 2,\n                  \"samples\":\n[\n                    false,\n                  ],,\n                  \"semantic_type\": \"\",,\n                  \"description\": \"\",,\n                },,\n                {\n                  \"column\":\n\"review_num_helpful\",,\n                  \"properties\": {\n                    \"dtype\":\n\"number\",,\n                    \"std\": 2.194837919560385,\n                    \"min\":
```



```

if not found_placeholder:
    print("No potential placeholders found in the DataFrame.")

Column 'review_title': Found 1 occurrences of potential placeholder
'Na'
Column 'review_username': Found 2 occurrences of potential placeholder
'none'
Column 'review_username': Found 3 occurrences of potential placeholder
'Unknown'

# Checking our column names
raw.columns

Index(['id', 'asins', 'brand', 'product_categories', 'product_keys',
      'manufacturer_name', 'review_date', 'review_do_recommend',
      'review_num_helpful', 'review_rating', 'review_text',
      'review_title',
      'review_username'],
      dtype='object')

#Checking the null values and data types after changes made
raw.info()

<class 'pandas.core.frame.DataFrame'>
Index: 34054 entries, 0 to 34624
Data columns (total 13 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   id                                    34054 non-null  object
1   asins                                34054 non-null  object
2   brand                                34054 non-null  object
3   product_categories                    34054 non-null  object
4   product_keys                          34054 non-null  object
5   manufacturer_name                     34054 non-null  object
6   review_date                           34054 non-null  datetime64[ns, UTC]
7   review_do_recommend                   34054 non-null  object
8   review_num_helpful                    34054 non-null  float64
9   review_rating                         34054 non-null  float64
10  review_text                           34054 non-null  object
11  review_title                          34054 non-null  object
12  review_username                       34054 non-null  object
dtypes: datetime64[ns, UTC](1), float64(2), object(10)
memory usage: 3.6+ MB

```

After cleaning the data set, we now have 34,054 rows and no missing values.

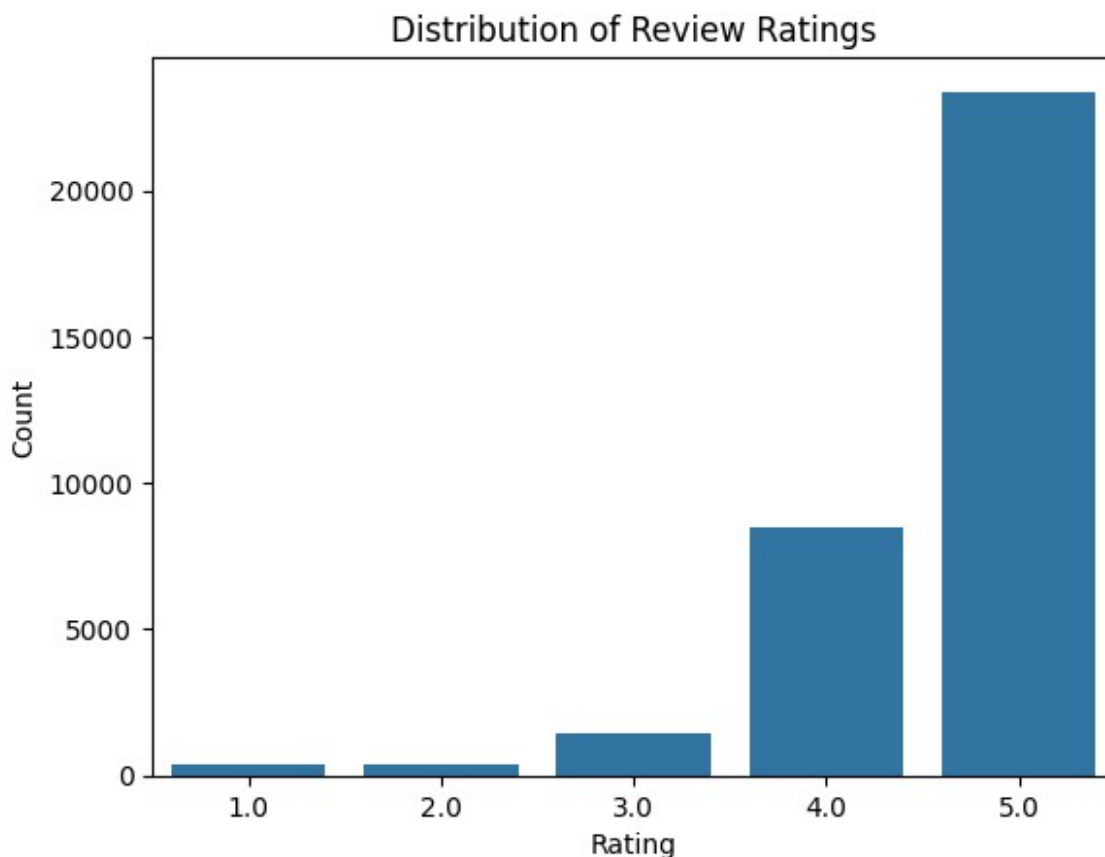
The data set is ready for EDA.

EXPLORATORY DATA ANALYSIS

UNIVARIATE ANALYSIS

1. Distribution of ratings Word frequency, Word cloud and Sentiment Distribution

```
# Distribution of ratings  
import matplotlib.pyplot as plt  
# Sentiment distribution (simple visualization based on ratings)  
sns.countplot(x='review_rating', data=raw)  
plt.title('Distribution of Review Ratings')  
plt.xlabel('Rating')  
plt.ylabel('Count')  
plt.show()
```



- The distribution of review ratings shows that most reviews tend to be positive, with higher counts towards ratings 4 and 5.

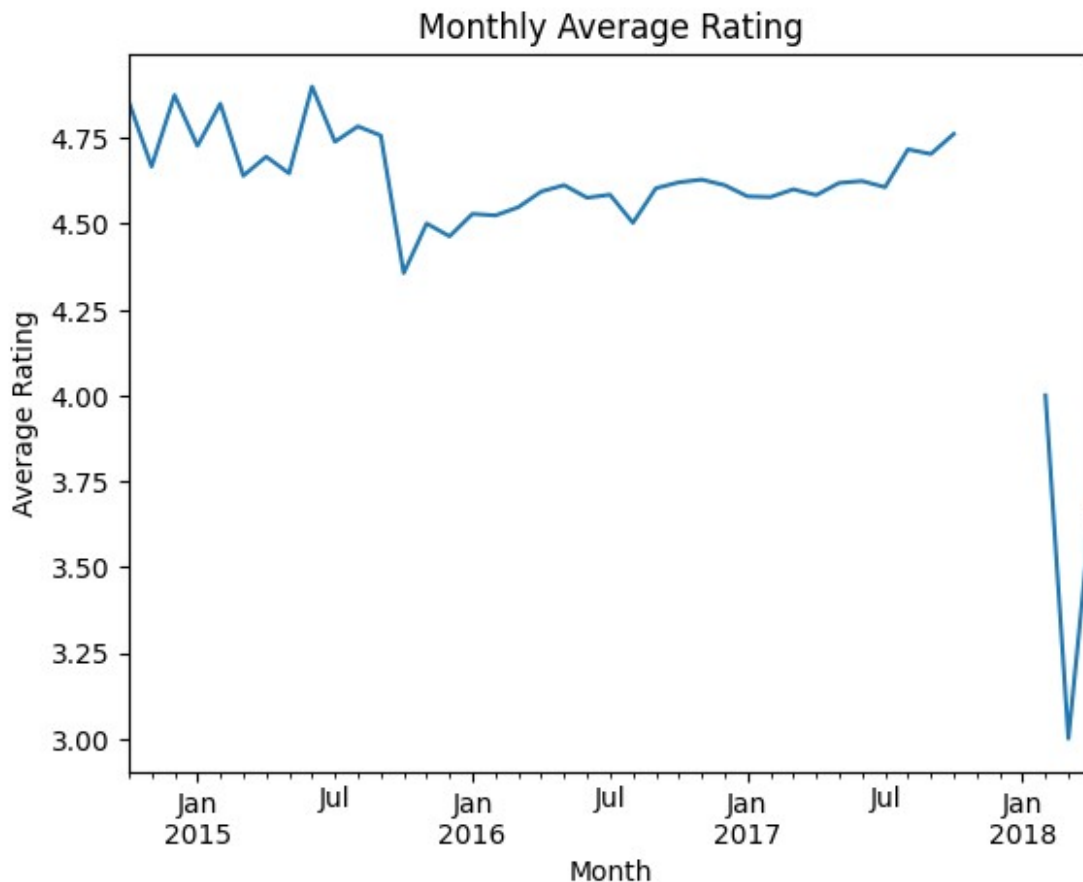
2.Temporal Analysis

```
# Temporal Analysis of rating over time
```

```

raw['review_date'] = pd.to_datetime(raw['review_date'])
raw.set_index('review_date', inplace=True)
raw['review_rating'].resample('M').mean().plot()
plt.title('Monthly Average Rating')
plt.xlabel('Month')
plt.ylabel('Average Rating')
plt.show()

```



- There is a slight fluctuation in average ratings over time, but no clear trend is evident from the monthly average ratings plot.

3. Reviews by product category

```

# Count occurrences of each category
category_counts = raw['product_categories'].value_counts().head(20)

# Extract top 20 categories and their counts
top_categories = category_counts.index

print("Top 20 Product Categories:")
print(category_counts)

```

```

# Assuming categories are separated by commas and need to be split
# Convert the 'product_categories' column to string type
raw['product_categories'] = raw['product_categories'].astype(str)

# Split the categories by commas
raw['product_categories'] = raw['product_categories'].str.split(',')

# Explode the list of categories
exploded_raw = raw.explode('product_categories')

# Group by 'product_categories' and calculate the mean review rating
mean_ratings = exploded_raw.groupby('product_categories')
['review_rating'].mean().sort_values(ascending=False)

mean_ratings

```

Top 20 Product Categories:

product_categories

Fire Tablets,Tablets,Computers & Tablets,All Tablets,Electronics, Tech
Toys, Movies, Music,Electronics,iPad & Tablets,Android Tablets,Frys
10965

Stereos,Remote Controls,Amazon Echo,Audio Docks & Mini Speakers,Amazon
Echo Accessories,Kitchen & Dining Features,Speaker

Systems,Electronics,TVs Entertainment,Clearance,Smart Hubs & Wireless
Routers,Featured Brands,Wireless Speakers,Smart Home & Connected
Living,Home Security,Kindle Store,Home Automation,Home, Garage &
Office,Home,Voice-Enabled Smart Assistants,Virtual Assistant

Speakers,Portable Audio & Headphones,Electronics Features,Amazon
Device Accessories,iPod, Audio Player Accessories,Home & Furniture

Clearance,Consumer Electronics,Smart Home,Surveillance,Home
Improvement,Smart Home & Home Automation Devices,Smart Hubs,Home
Safety & Security,Voice Assistants,Alarms & Sensors,Amazon

Devices,Audio,Holiday Shop 6606

Back To College,College Electronics,College TVs & Home

Theater,Electronics,TVs & Home Theater,Streaming Devices,Featured
Brands,Amazon Devices,Holiday Shop,Ways To Shop,TV & Home

Theater,Streaming Media Players,All Streaming Media Players,TVs
Entertainment,Video Games,Kindle Store,Electronics Features,Kids &
Family,Fire TV

5051

Walmart for Business,Office

Electronics,Tablets,Office,Electronics,iPad & Tablets,Windows

Tablets,All Windows Tablets,Computers & Tablets,E-Readers &
Accessories,E-Readers,eBook Readers,Kindle E-readers,Computers/Tablets

& Networking,Tablets & eBook Readers,Electronics Features,Books &
Magazines,Book Accessories,eReaders,TVs & Electronics,Computers &

Laptops,Tablets & eReaders

3175

Electronics,iPad & Tablets,All Tablets,Fire Tablets,Tablets,Computers

& Tablets

2812

Tablets,Fire Tablets,Computers & Tablets,All Tablets

1698

Computers/Tablets & Networking,Tablets & eBook Readers,Computers & Tablets,Tablets,All Tablets

1038

Featured Brands,Electronics,Amazon Devices,Home,Home Improvement,Home Safety & Security,Home Security,Alarms & Sensors,Smart Home & Home Automation Devices,Mobile,Mobile Speakers,Mobile Bluetooth Speakers,Smart Hubs & Wireless Routers,Smart Hubs,Home, Garage & Office,Smart Home,Voice Assistants,Smart Home & Connected Living,Amazon Tap,Portable Audio,MP3 Accessories,Speakers,Amazon Echo,Electronics Features,TVs & Electronics,Portable Audio & Electronics,MP3 Player Accessories,Home Theater & Audio,Kindle Store,Frys,Electronic Components,Home Automation,Electronics, Tech Toys, Movies, Music,Audio,Bluetooth Speakers

633

Walmart for Business,Office Electronics,Tablets,Electronics,iPad & Tablets,All Tablets,Computers & Tablets,E-Readers & Accessories,Kindle E-readers,Electronics Features,eBook Readers,See more Amazon Kindle Voyage (Wi-Fi),See more Amazon Kindle Voyage 4GB, Wi-Fi 3G (Unlocked...

580

Electronics Features,Fire Tablets,Computers & Tablets,Tablets,All Tablets,Computers/Tablets & Networking,Tablets & eBook Readers

371

Fire Tablets,Tablets,Computers & Tablets,All Tablets,Computers/Tablets & Networking,Tablets & eBook Readers

269

Electronics,Tablets & E-Readers,Tablets,Back To College,College Electronics,College Ipads & Tablets,Featured Brands,Amazon Devices,Electronics Deals,Computers & Tablets,All Tablets,Electronics Features,eBook Readers

254

Electronics,iPad & Tablets,All Tablets,Computers/Tablets & Networking,Tablets & eBook Readers,Computers & Tablets,E-Readers & Accessories,E-Readers,Used:Computers Accessories,Used:Tablets,Computers,iPads Tablets,Kindle E-readers,Electronics Features

212

Tablets,Fire Tablets,Electronics,Computers,Computer Components,Hard Drives & Storage,Computers & Tablets,All Tablets

158

eBook Readers,Kindle E-readers,Computers & Tablets,E-Readers & Accessories,E-Readers

67

Chargers & Adapters,Computers & Accessories,Tablet & E-Reader Accessories,Amazon Devices & Accessories,Fire Tablet

Accessories,Electronics,Power Adapters & Cables,Cell Phones,Amazon Device Accessories,Cell Phone Accessories,Cell Phone Batteries & Power,Tablet Accessories,Featured Brands,Kindle Fire (2nd Generation) Accessories,Kindle Store,Home Improvement,Fire (5th Generation) Accessories,Electrical,Amazon Devices,Home,Tablets & E-Readers,Cables & Chargers

54

Computers & Tablets,E-Readers & Accessories,eBook Readers,Kindle E-readers

51

Electronics,iPad & Tablets,All Tablets,Computers & Tablets,Tablets,eBook Readers

30

Computers & Tablets,Tablets,All Tablets,Computers/Tablets & Networking,Tablets & eBook Readers,Fire Tablets,Frys

10

Fire Tablets,Tablets,Computers & Tablets,All Tablets

7

Name: count, dtype: int64

product_categories

Top Rated 4.833333

Computer Accessories 4.833333

Cases & Bags 4.833333

Kids' Tablets 4.833333

Book Accessories 4.772283

...

Frys 4.458524

Movies 4.458463

Tech Toys 4.458463

Music 4.458463

Android Tablets 4.454172

Name: review_rating, Length: 120, dtype: float64

```
plt.figure(figsize=(20, 18))
```

```
# Create a bar plot with a color gradient
```

```
bars = sns.barplot(y=top_categories, x=category_counts.values,  
palette="viridis")
```

```
# Add value labels to the bars
```

```
for bar, count in zip(bars.patches, category_counts.values):
```

```
    plt.text(count + 10, # x-coordinate position  
             bar.get_y() + bar.get_height() / 2, # y-coordinate
```

```
position
```

```
             f'{count}', # formatted label text  
             ha='center', va='center', # horizontal and vertical
```

```
alignment
```

```
             fontsize=10, color='black') # text properties
```

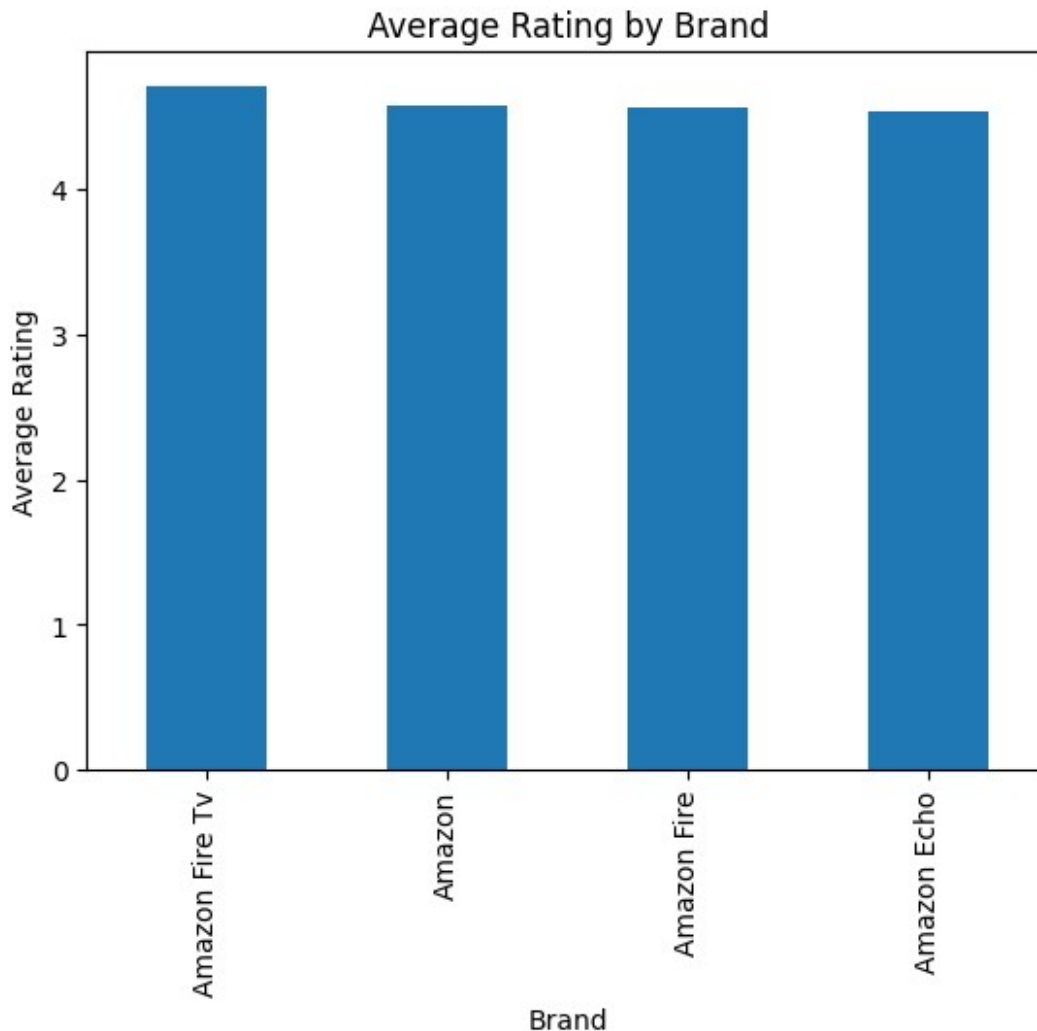
```
plt.title('Top 20 Product Categories by Count of Reviews',
fontsize=16)
plt.xlabel('Count', fontsize=14)
plt.ylabel('Product Category', fontsize=14)

plt.xticks(fontsize=12)
plt.yticks(fontsize=12)

plt.tight_layout()
plt.show()
```



```
# Plot review rating by brand
raw.groupby('brand')
['review_rating'].mean().sort_values(ascending=False).plot(kind='bar')
plt.title('Average Rating by Brand')
plt.xlabel('Brand')
plt.ylabel('Average Rating')
plt.show()
```



```
# Assuming categories are separated by commas and need to be split
# Convert the 'product_categories' column to string type
raw['product_categories'] = raw['product_categories'].astype(str)

# Split the categories by commas
raw['product_categories'] = raw['product_categories'].str.split(',')

# Explode the list of categories
exploded_raw = raw.explode('product_categories')

# Group by 'product_categories' and calculate the mean review rating
mean_ratings = exploded_raw.groupby('product_categories')
['review_rating'].mean().sort_values(ascending=False)

mean_ratings

product_categories
'Kindle E-readers']      4.862745
```



```

['Computers & Tablets'      4.836066
 'Top Rated']              4.833333
 'Cases & Bags'             4.833333
 'Kids' Tablets'           4.833333
...
 ' Movies'                 4.458463
['Fire Tablets'            4.456947
 'Frys']                   4.454446
 'Android Tablets'         4.454172
['Electronics Features'    4.425876
Name: review_rating, Length: 139, dtype: float64

```

Conclusions

- Fire Tablets, Tablets, Computers & Tablets: Dominates with 10,965 reviews, indicating a strong presence in consumer feedback.
- Stereos, Remote Controls, Amazon Echo: Follows with 6,606 reviews, highlighting significant interest in home electronics and smart devices.
- Back To College, College Electronics: Shows strong engagement in electronics geared towards college students, with 5,051 reviews.

4. Most helpful Votes

```

# Most helpful reviews
raw.sort_values(by='review_num_helpful', ascending=False).head(10)

{"summary":{"\n  \"name\": \"raw\", \n  \"rows\": 10, \n  \"fields\": [\n    {\n      \"column\": \"review_date\", \n      \"properties\": {\n        \"dtype\": \"date\", \n        \"min\": \"2014-11-16 00:00:00+00:00\", \n        \"max\": \"2016-10-05 00:00:00+00:00\", \n        \"num_unique_values\": 10, \n        \"samples\": [\n          \"2014-11-16 00:00:00+00:00\", \n          \"2016-10-05 00:00:00+00:00\", \n          \"2016-05-22 00:00:00+00:00\" \n        ], \n        \"semantic_type\": \"\", \n        \"description\": \"\", \n        \"column\": \"id\", \n        \"properties\": {\n          \"dtype\": \"string\", \n          \"num_unique_values\": 6, \n          \"samples\": [\n            \"AVphgVaX1cnluZ0-DR74\", \n            \"AVqkIiKWnnc1JgDc3khH\", \n            \"AVphPmHuilaPnD_x3E5h\" \n          ], \n          \"semantic_type\": \"\", \n          \"description\": \"\", \n          \"column\": \"asins\", \n          \"properties\": {\n            \"dtype\": \"string\", \n            \"num_unique_values\": 6, \n            \"samples\": [\n              \"B018Y2290U\", \n              \"B01AHB9CYG\", \n              \"B00IOY8XWQ\" \n            ], \n            \"semantic_type\": \"\", \n            \"description\": \"\", \n            \"column\": \"brand\", \n            \"properties\": {\n              \"dtype\": \"category\", \n              \"num_unique_values\": 2, \n              \"samples\": [\n                \"Amazon Fire Tv\", \n                \"Amazon\" \n              ], \n              \"semantic_type\": \"\", \n              \"description\": \"\" \n            } \n          } \n        ] \n      } \n    ] \n  } \n}

```

```

n    },\n    {\n        \"column\": \"product_categories\", \n        \"properties\": {\n            \"dtype\": \"object\", \n            \"semantic_type\": \"\", \n            \"description\": \"\" \n        } \n    },\n    {\n        \"column\": \"product_keys\", \n        \"properties\": {\n            \"dtype\": \"string\", \n            \"num_unique_values\": 6, \n            \"samples\": [\n                \"firetablet7displaywifi8gbincludesspecialoffersmagenta/5025800,841667103105,0841667103105,amazon/b018y229ou,firetablet7displaywifi8gbincludesspecialoffersmagenta/b018y229ou\", \n                \"841667104690,allnewfirehd8tablet8hddisplaywifi32gbincludesspecialoffersmagenta/5620408,0841667104690,allnewfirehd8tablet8hddisplaywifi32gbincludesspecialoffersmagenta/b01ahb9cyg,amazon/53004761\" \n            ], \n            \"semantic_type\": \"\", \n            \"description\": \"\" \n        } \n    },\n    {\n        \"column\": \"manufacturer_name\", \n        \"properties\": {\n            \"dtype\": \"category\", \n            \"num_unique_values\": 1, \n            \"samples\": [\n                \"Amazon\" \n            ], \n            \"semantic_type\": \"\", \n            \"description\": \"\" \n        } \n    },\n    {\n        \"column\": \"review_do_recommend\", \n        \"properties\": {\n            \"dtype\": \"category\", \n            \"num_unique_values\": 2, \n            \"samples\": [\n                false \n            ], \n            \"semantic_type\": \"\", \n            \"description\": \"\" \n        } \n    },\n    {\n        \"column\": \"review_num_helpful\", \n        \"properties\": {\n            \"dtype\": \"number\", \n            \"std\": 18.41165090069027, \n            \"min\": 63.0, \n            \"max\": 109.0, \n            \"num_unique_values\": 10, \n            \"samples\": [\n                64.0 \n            ], \n            \"semantic_type\": \"\", \n            \"description\": \"\" \n        } \n    },\n    {\n        \"column\": \"review_rating\", \n        \"properties\": {\n            \"dtype\": \"number\", \n            \"std\": 0.9944289260117531, \n            \"min\": 2.0, \n            \"max\": 5.0, \n            \"num_unique_values\": 4, \n            \"samples\": [\n                5.0 \n            ], \n            \"semantic_type\": \"\", \n            \"description\": \"\" \n        } \n    },\n    {\n        \"column\": \"review_text\", \n        \"properties\": {\n            \"dtype\": \"string\", \n            \"num_unique_values\": 10, \n            \"samples\": [\n                \"I am a big fan of e-readers. I prefer the e-ink screens over tablet screen when reading books. I decided to pick up the new Kindle Voyage. Here are my thoughts.First, the Kindle is much sleeker and lighter than the Paperwhite model. It's very easy to hold for a long length of time without getting tired.Secondly, the addition of the page turn buttons is a welcomed addition. The buttons make it so easy to hold the Kindle with one hand and turn pages.Lastly the screen. I marked off a star because I had to return my first Voyage directly to Amazon because the top half of the screen had a yellow tint to it that was very distracting while reading. The replacement Kindle Voyage that I received had a perfect screen and the 300 ppi looks amazing. That being said, Amazon needs to really focus on quality control because

```

```
I've read that the yellow tint on the screen is a common issue. If
you're spending $200 for an e-reader, then the screen should be
perfect. Especially since this is the 3rd generation of the lighted
screens for the e-readers.If you're an avid reader, then the Kindle
Voyage would be a worthy upgrade to make. If you're a casual reader
who has a Paperwhite, you're probably okay with what you have.\\n
],\\n      \\\"semantic_type\\\": \\\"\\\",\\n      \\\"description\\\": \\\"\\\"\\n
}\\n    },\\n    {\\n      \\\"column\\\": \\\"review_title\\\",\\n
\\\"properties\\\": {\\n      \\\"dtype\\\": \\\"string\\\",\\n
\\\"num_unique_values\\\": 9,\\n      \\\"samples\\\": [\\n
\\\"Almost...\\\"\\n      ],\\n      \\\"semantic_type\\\": \\\"\\\",\\n
\\\"description\\\": \\\"\\\"\\n      }\\n    },\\n    {\\n      \\\"column\\\":
\\\"review_username\\\",\\n      \\\"properties\\\": {\\n      \\\"dtype\\\":
\\\"string\\\",\\n      \\\"num_unique_values\\\": 10,\\n      \\\"samples\\\":
[\\n      \\\"Cooper25\\\"\\n      ],\\n      \\\"semantic_type\\\":
\\\"\\\",\\n      \\\"description\\\": \\\"\\\"\\n      }\\n    }\\n  ]\\
n}\\\", \"type\": \"dataframe\"}
```

BIVARIATE ANALYSIS

1. Helpful votes vs rating

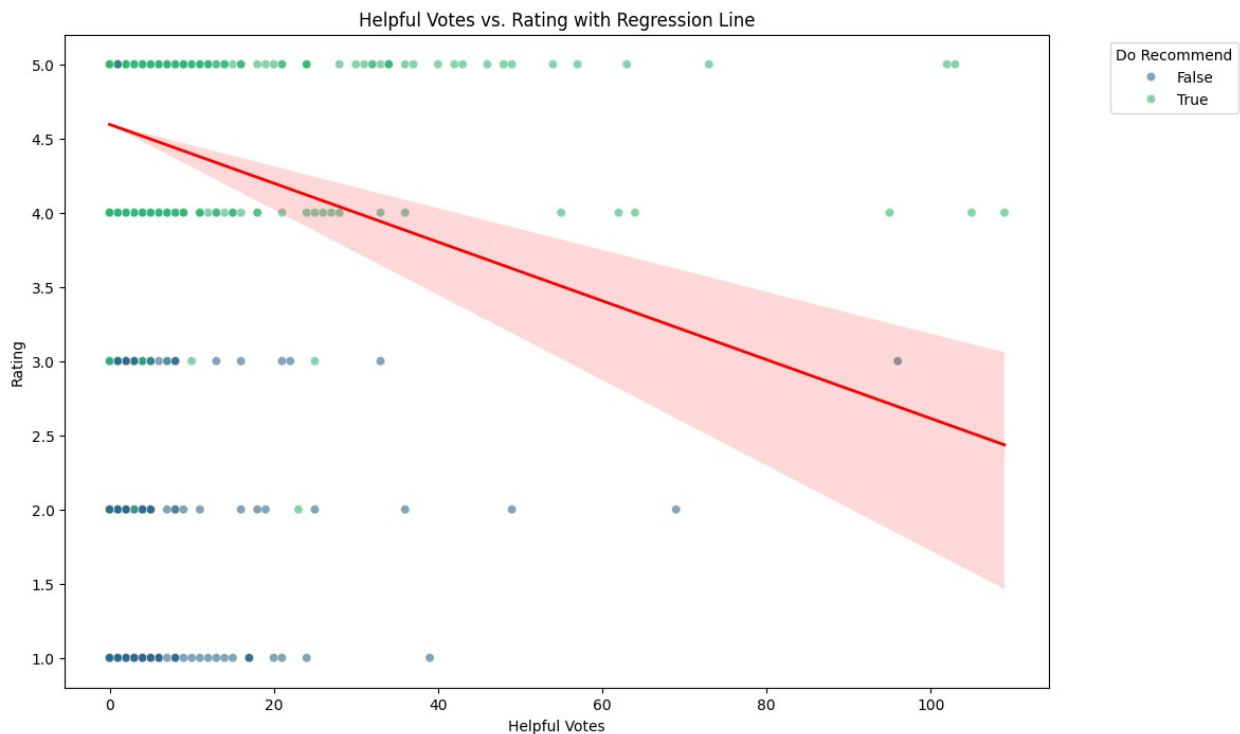
```
plt.figure(figsize=(12, 8))

# Scatter plot with color coding, size encoding, and transparency
scatter = sns.scatterplot(
    x='review_num_helpful',
    y='review_rating',
    hue='review_do_recommend',
    sizes=(20, 200), # Minimum and maximum size of points
    alpha=0.6,
    palette='viridis', # Using a different color palette
    data=raw
)

# Add a regression line
sns.regplot(
    x='review_num_helpful',
    y='review_rating',
    scatter=False,
    color='red',
    line_kws={"linewidth": 2},
    data=raw
)

plt.title('Helpful Votes vs. Rating with Regression Line')
plt.xlabel('Helpful Votes')
plt.ylabel('Rating')
plt.legend(title='Do Recommend', loc='upper right',
```

```
bbox_to_anchor=(1.2, 1))
plt.show()
```

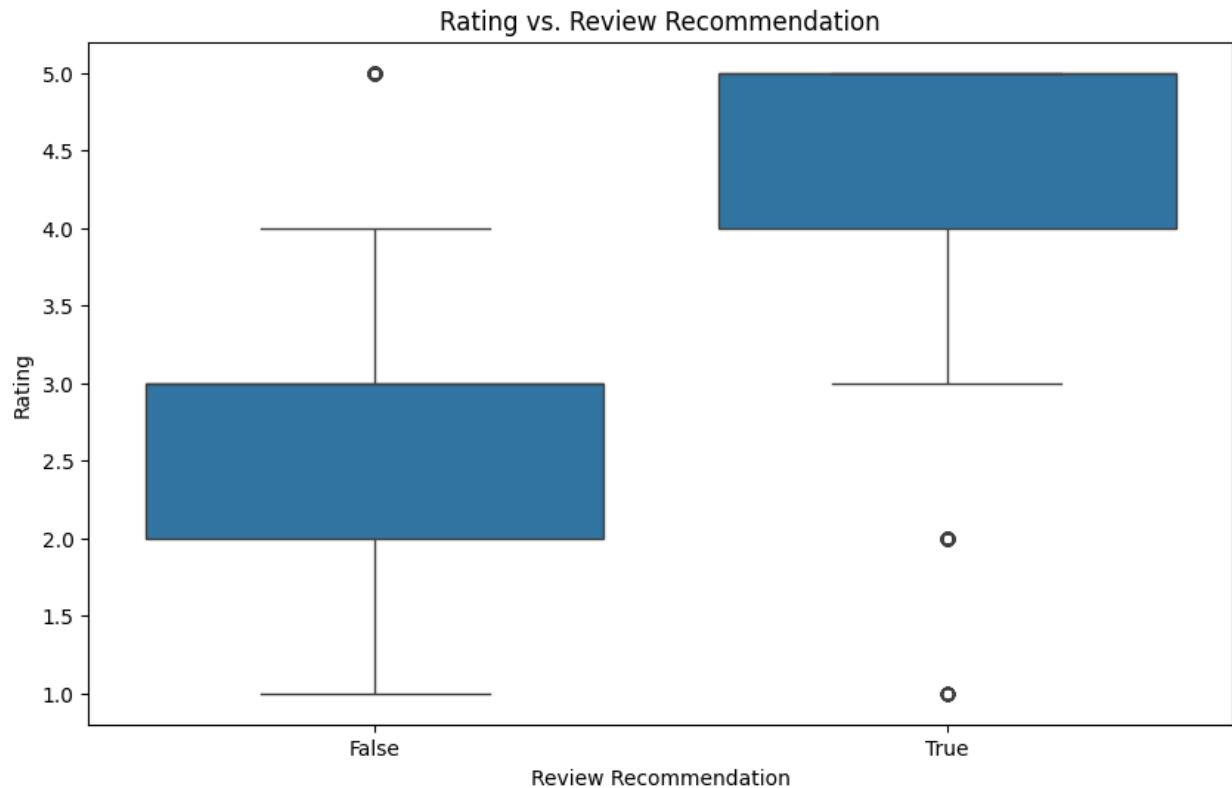


- The scatter plot and regression analysis of helpful votes versus rating illustrate a positive correlation, indicating that more helpful reviews tend to have higher ratings.
- This suggests that customers find high-rated reviews more useful

2. Rating vs. Review recommendation

```
# Convert review_do_recommend to a categorical type
raw['review_do_recommend'] =
raw['review_do_recommend'].astype('category')

# Box plot of rating vs. review recommendation
plt.figure(figsize=(10, 6))
sns.boxplot(x='review_do_recommend', y='review_rating', data=raw)
plt.title('Rating vs. Review Recommendation')
plt.xlabel('Review Recommendation')
plt.ylabel('Rating')
plt.show()
```

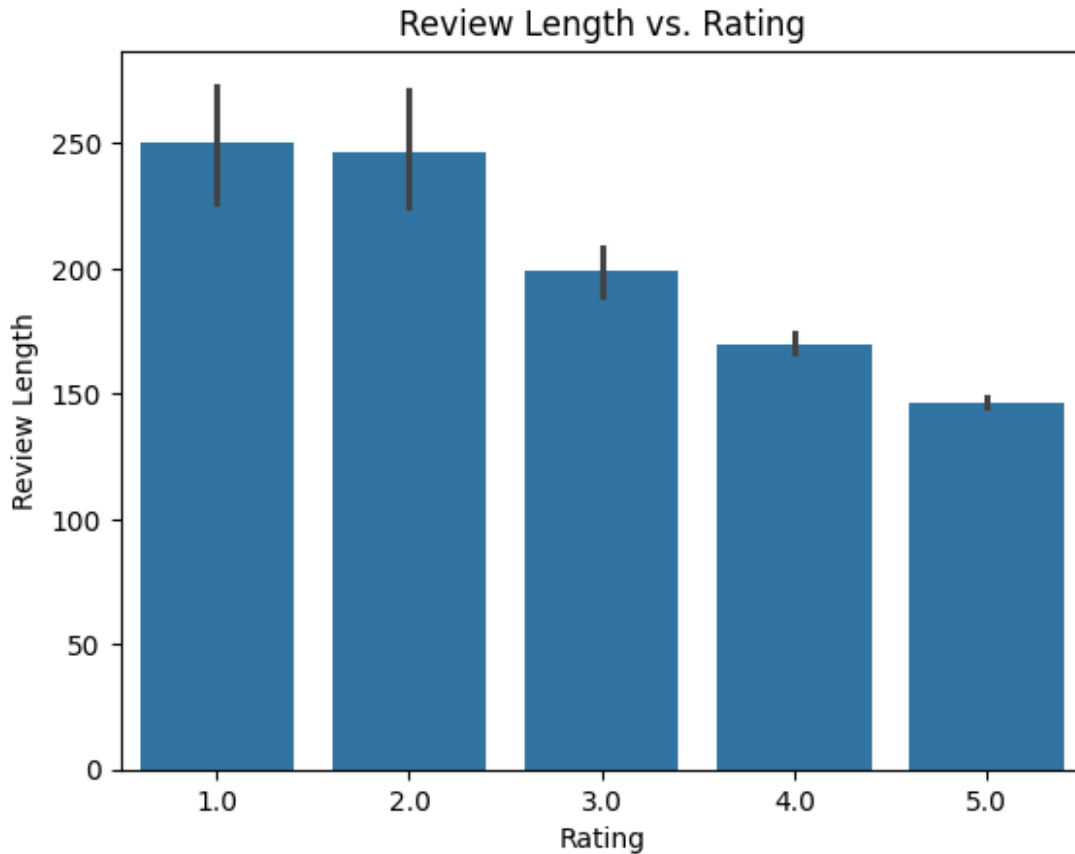


- The analysis shows that reviews with a positive recommendation (review_do_recommend = True) generally have higher ratings compared to those without a recommendation.
- This highlights the influence of product satisfaction on recommendation.

3. Rating vs Length

```
raw['review_length'] = raw['review_text'].apply(len)

sns.barplot(x='review_rating', y='review_length', data=raw)
plt.title('Review Length vs. Rating')
plt.xlabel('Rating')
plt.ylabel('Review Length')
plt.show()
```

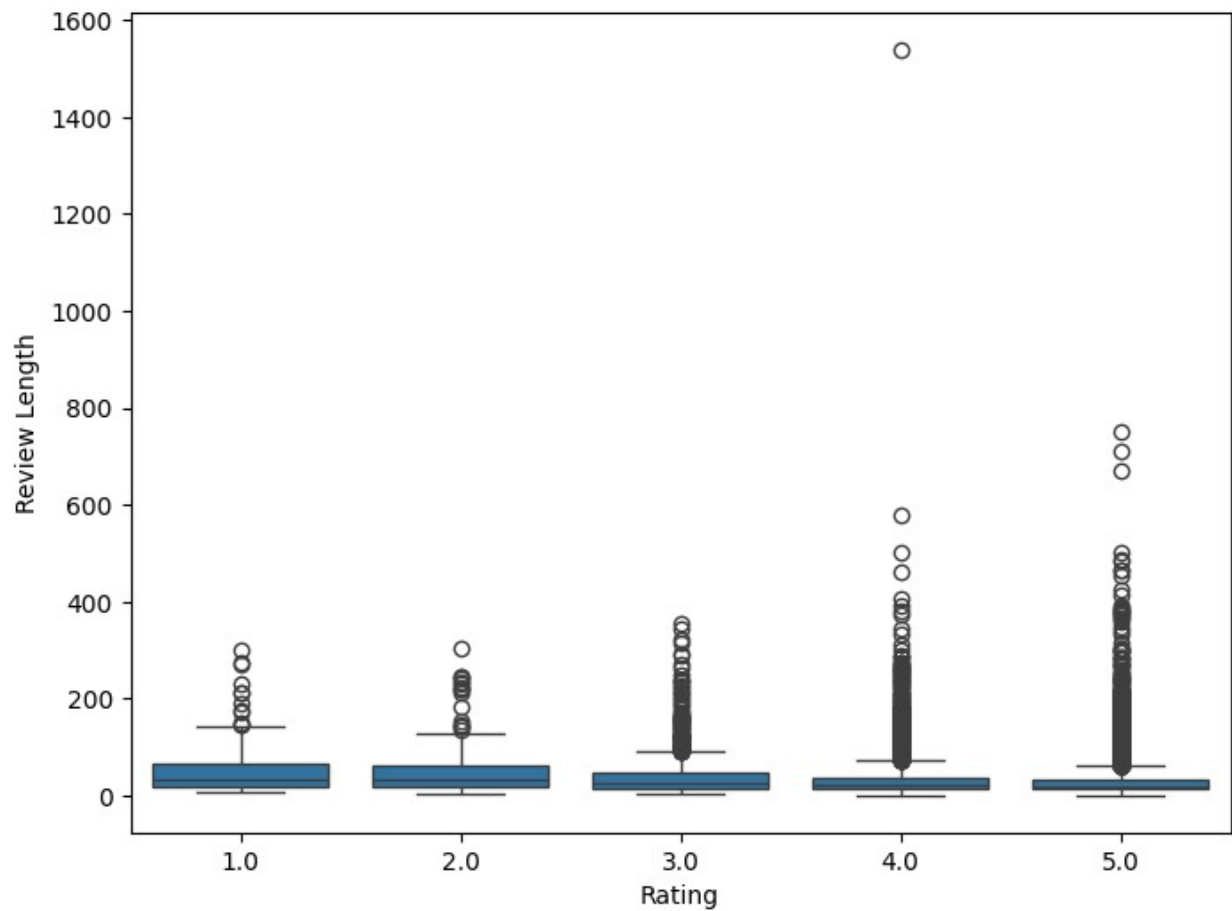


- This visualization illustrates the relationship between review length and review rating. It is evident that shorter reviews tend to receive higher ratings.

```
word_count=[]
for s1 in raw.review_text:
    word_count.append(len(str(s1).split()))
plt.figure(figsize = (8,6))

import seaborn as sns
import matplotlib.pyplot as plt
sns.boxplot(x="review_rating",y=word_count,data=raw)
plt.xlabel('Rating')
plt.ylabel('Review Length')

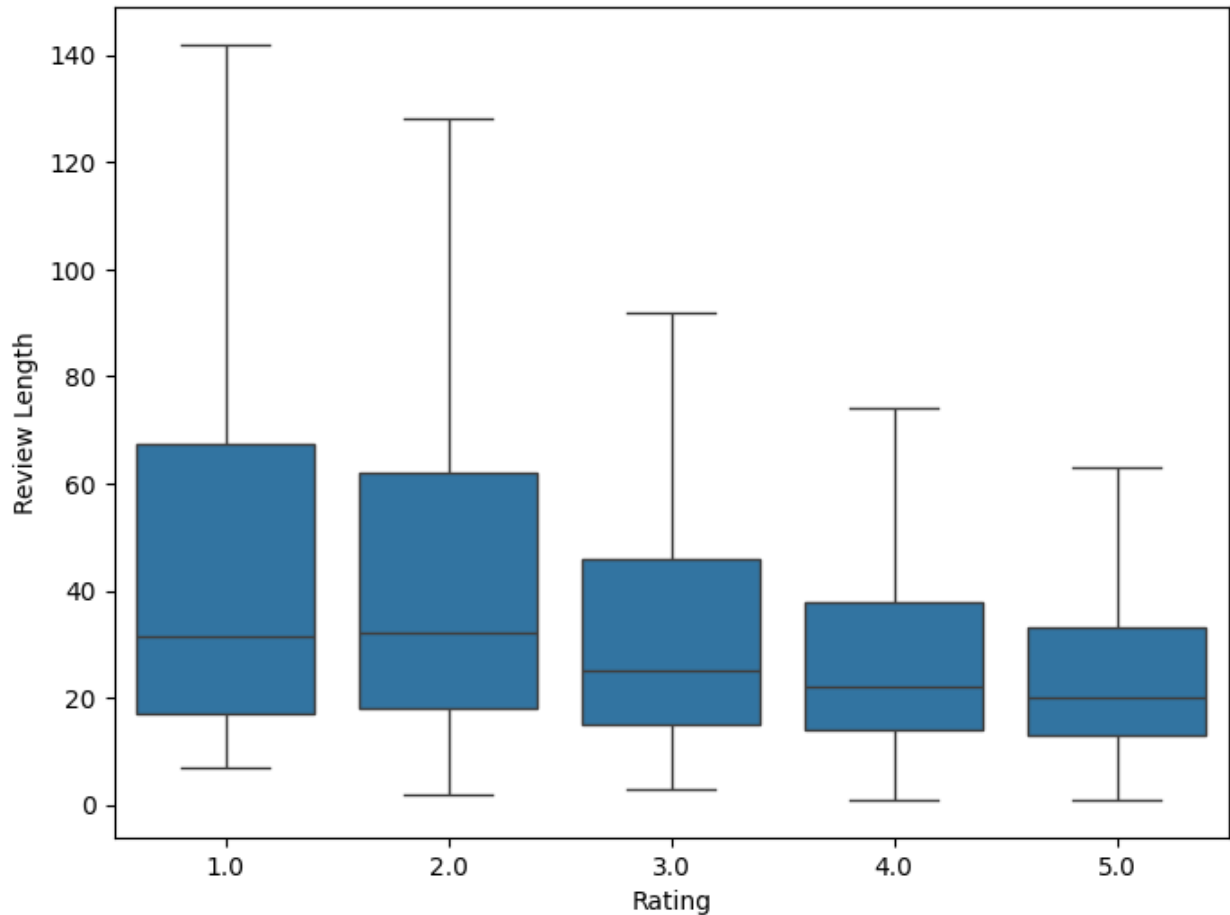
plt.show()
```



- Due to the presence of outliers shown in the box plot, our visualization is currently obscured. To improve clarity, we will proceed by removing these outliers from the dataset.

```
# Generate box plots excluding outliers
```

```
plt.figure(figsize = (8,6))  
sns.boxplot(x="review_rating",y=word_count,data=raw,showfliers=False)  
plt.xlabel('Rating')  
plt.ylabel('Review Length')  
plt.show()
```



- We can now see that shorter reviews tend to receive higher ratings much better.

Conclusions

The bar plot and box plot analyses show the relationship between review ratings and the length of reviews:

Bar Plot Analysis: Indicates that longer reviews are generally associated with lower ratings. This suggests that while longer reviews can provide richer insights, their association with lower ratings indicates that customers who invest more time in detailing their experiences often do so when they feel particularly disappointed or dissatisfied.

Box Plot Analysis: Initially showed outliers affecting clarity in visualization. After excluding outliers, the relationship between review length and rating became clearer

Lower ratings tend to have a wider range of review lengths, suggesting variability in experiences or dissatisfaction reasons.

Higher ratings are associated with a more concentrated range of review lengths, possibly indicating clearer satisfaction or positive experiences with the product.

These insights provide a deeper understanding of how review characteristics such as recommendation status and review length correlate with customer ratings, contributing valuable insights for product evaluation and improvement strategies.

3. Multivariate Analysis

1. Scatter plot of reviews

```
# Ensure the column names are correct
review_rating_col = 'review_rating'
review_num_helpful_col = 'review_num_helpful'
total_votes_col = 'total_votes'
review_did_purchase_col = 'review_did_purchase'

# Check if 'review_did_purchase' exists, if not create it with a default value
if review_did_purchase_col not in raw.columns:
    raw[review_did_purchase_col] = False

# Ensure 'total_votes' column exists, if not create it with a default value
if total_votes_col not in raw.columns:
    raw[total_votes_col] = 0

# Plotting

fig = plt.figure(figsize=(10, 8))
ax = fig.add_subplot(111, projection='3d')

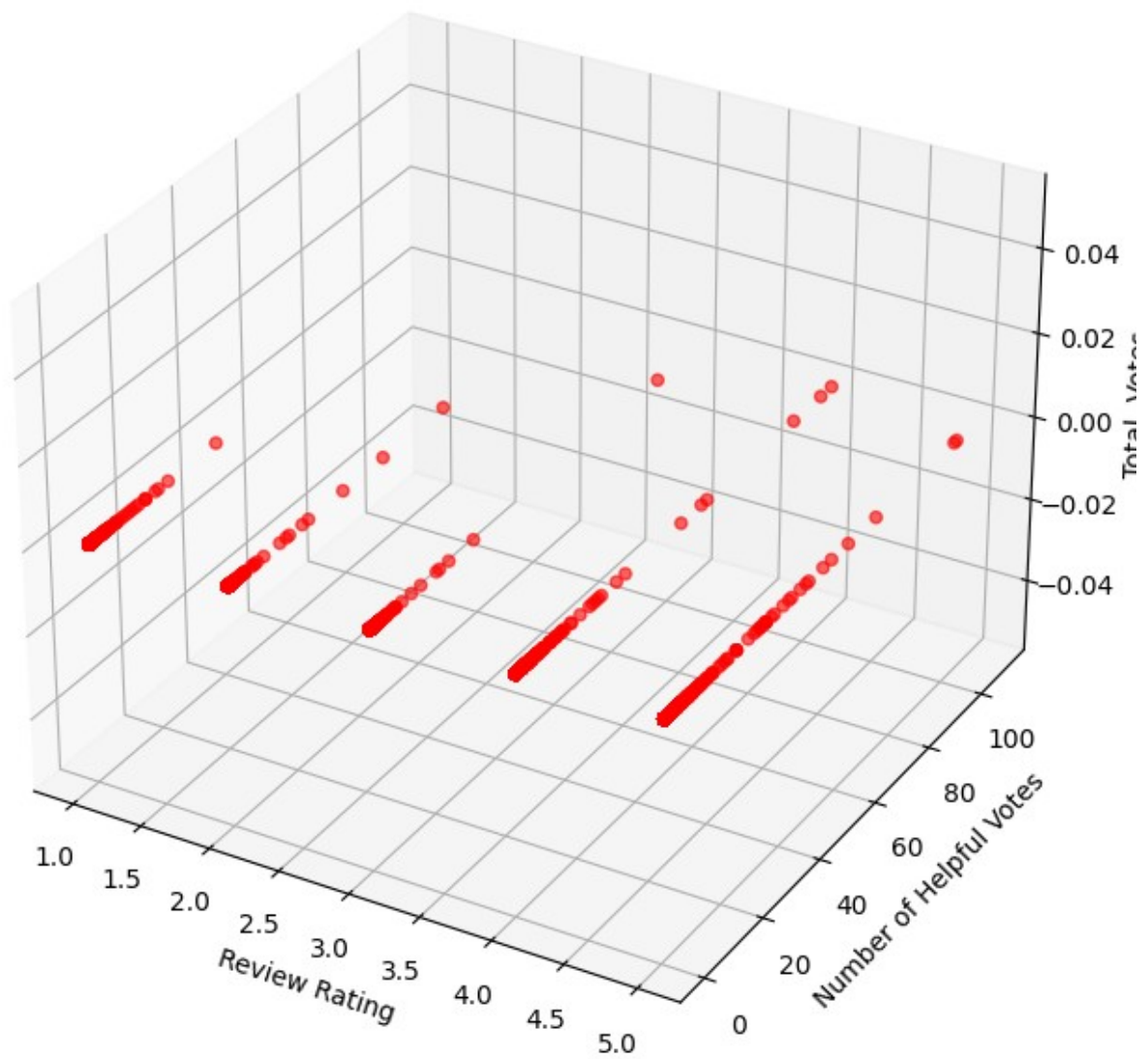
# Map verified purchase to colors
colors = raw[review_did_purchase_col].map({True: 'blue', False: 'red'})

sc = ax.scatter(raw[review_rating_col], raw[review_num_helpful_col],
               raw[total_votes_col], c=colors, alpha=0.6)

# Adding labels and title
ax.set_xlabel('Review Rating')
ax.set_ylabel('Number of Helpful Votes')
ax.set_zlabel('Total Votes')
plt.title('3D Scatter Plot of Reviews')

plt.show()
```

3D Scatter Plot of Reviews

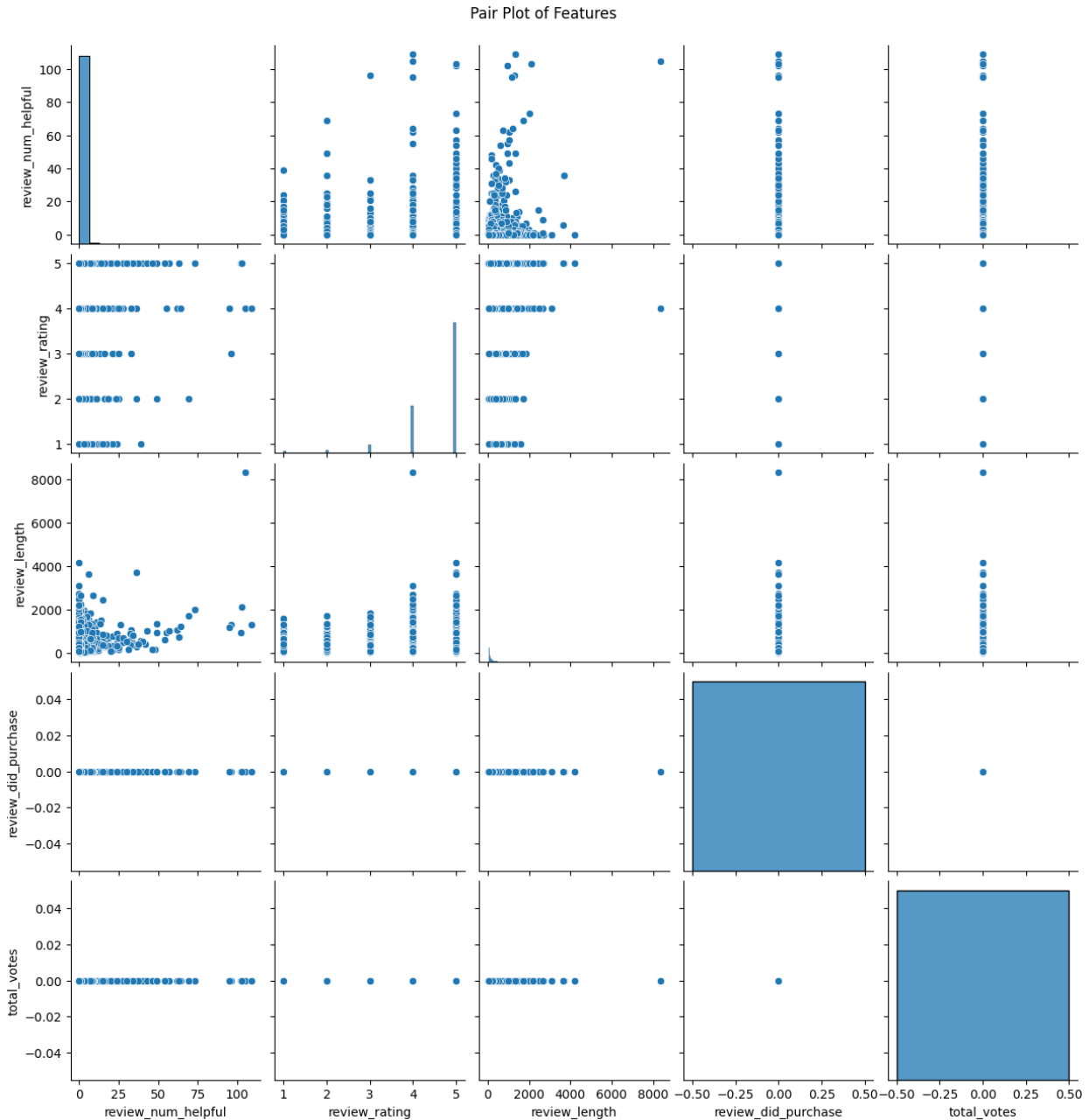


Conclusions

- Visualizing reviews based on rating, helpful votes, and total votes shows various patterns, but it doesn't clearly reveal distinct groups based on whether the purchase was verified.

2. Pair Plot of Features

```
sns.pairplot(raw)
plt.suptitle('Pair Plot of Features', y=1.02)
plt.show()
```



Conclusions

Pair Plot: The pair plot visually explored relationships between different numerical features in the dataset. It provides a quick overview of potential correlations and distributions among variables, aiding in identifying patterns or trends that might warrant further investigation.

Data pre-processing

```
# Check the column names
print(raw.columns)
```

```
Index(['id', 'asins', 'brand', 'product_categories', 'product_keys',
      'manufacturer_name', 'review_do_recommend',
      'review_num_helpful',
      'review_rating', 'review_text', 'review_title',
      'review_username',
      'review_length', 'review_did_purchase', 'total_votes'],
      dtype='object')
```

- Let's preview the first sentence in our text

```
# Previewing the first sentence in our text

first_document = raw.iloc[2]['review_text']
first_document

{"type": "string"}

# Changing the name of our dataframe

data = pd.DataFrame(raw)
```

- For NLP preprocessing, we'll eliminate stopwords, punctuation, and numbers, and convert text to lowercase.
- Subsequently, tokenizing our data is essential because it breaks down text into individual words or tokens, enabling deeper analysis and understanding of the textual content.

```
# Download NLTK stopwords and punctuation
nltk.download('stopwords')
nltk.download('punkt')

# Load stopwords and punctuation
stop_words = set(stopwords.words('english'))

# Function to clean and preprocess text
def clean_text(text):
    # Ensure text is a string and lowercase
    text = str(text).lower()

    # Remove numbers
    text = re.sub(r'\d+', '', text)

    # Remove punctuation
    text = text.translate(str.maketrans('', '', string.punctuation))

    # Tokenization using regex pattern
    pattern = "([a-zA-Z]+(?:'[a-z]+)?)"
    tokens = nltk.regexp_tokenize(text, pattern)

    # Remove stopwords
    clean_tokens = [token for token in tokens if token not in
```

```

stop_words]

    return ' '.join(clean_tokens)

data['clean_text'] = raw['review_text'].apply(clean_text)
data['clean_title'] = raw['review_title'].apply(clean_text)

# Display the cleaned text along with original columns
data[['review_text', 'review_title', 'clean_text', 'clean_title']]

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!

{"summary":{"\n  \"name\": \"data[['review_text', 'review_title',
'clean_text', 'clean_title']]\", \n  \"rows\": 34054, \n  \"fields\": [\n
    {\n      \"column\": \"review_date\", \n      \"properties\": {\n
        \"dtype\": \"date\", \n        \"min\": \"2014-10-24 00:00:00+00:00\", \n
        \"max\": \"2018-04-18 00:00:00+00:00\", \n
        \"num_unique_values\": 941, \n        \"samples\": [\n          \"2015-
11-29 00:00:00+00:00\", \n          \"2016-03-14 00:00:00+00:00\", \n
          \"2016-08-21 00:00:00+00:00\" \n        ], \n        \"semantic_type\":
        \"\", \n        \"description\": \"\" \n      }, \n      { \n
        \"column\": \"review_text\", \n        \"properties\": {\n
          \"dtype\": \"string\", \n          \"num_unique_values\": 34054, \n
          \"samples\": [\n            \"My kids love this product, as do I.
Parental restrictions can be set and they know when they have to shut
them off. Good battery life too.\", \n            \"This is an excellent
replacement for my Apple TV. I love it. Quick and easy to use. My
whole family enjoys it.\", \n            \"Excellent for what I or family
want to use it for and the price IS nice!\" \n          ], \n          \"semantic_type\": \"\", \n          \"description\": \"\" \n        }, \n        { \n
          \"column\": \"review_title\", \n          \"properties\": {\n
            \"dtype\": \"string\", \n            \"num_unique_values\": 19448, \n
            \"samples\": [\n              \"Money's Worth\", \n              \"bought for sling\", \n              \"Easy
to used.\" \n            ], \n            \"semantic_type\": \"\", \n
            \"description\": \"\" \n          }, \n          { \n
            \"column\":
            \"clean_text\", \n            \"properties\": {\n
              \"dtype\":
              \"string\", \n              \"num_unique_values\": 33957, \n
              \"samples\": [\n                \"im currently love item setting reminders
also playing music spotify bluetooth\", \n                \"nice size nice
looklove feel turning page great reading experience\", \n
                \"good quality android tablet also makes nice gift\" \n              ], \n
              \"semantic_type\": \"\", \n              \"description\": \"\" \n            }, \n            { \n
              \"column\": \"clean_title\", \n              \"properties\": {\n
                \"dtype\": \"category\", \n                \"num_unique_values\": 13899, \n
                \"samples\": [\n

```

```

\"hate\", \n          \"great purchase first time kindle buyer\", \n
\"great sounding musicand much\" \n          ], \n
\"semantic_type\": \"\", \n          \"description\": \"\" \n          } \n
n      } \n      ], \"type\": \"dataframe\"}

# Dropping the original columns as we now have the clean ones
data.drop(columns = ['review_text', 'review_title'] , inplace = True)
data.head(2)

{\"summary\": { \n      \"name\": \"data\", \n      \"rows\": 34054, \n
\"fields\": [ \n          { \n              \"column\": \"review_date\", \n
\"properties\": { \n                  \"dtype\": \"date\", \n                  \"min\":
\"2014-10-24 00:00:00+00:00\", \n                  \"max\": \"2018-04-18
00:00:00+00:00\", \n                  \"num_unique_values\": 941, \n
\"samples\": [ \n                      \"2015-11-29 00:00:00+00:00\", \n
\"2016-03-14 00:00:00+00:00\", \n                      \"2016-08-21
00:00:00+00:00\" \n                  ], \n                  \"semantic_type\": \"\", \n
\"description\": \"\" \n              }, \n          { \n              \"column\":
\"id\", \n              \"properties\": { \n                  \"dtype\": \"category\", \n
\"num_unique_values\": 24, \n                  \"samples\": [ \n
\"AVqkIhxunnc1JgDc3kg_\", \n                  \"AVpgdkC8ilAPnD_xsvyi\", \n
\"AVqkIhwDv8e3D10-lebb\" \n                  ], \n                  \"semantic_type\":
\"\", \n                  \"description\": \"\" \n              }, \n          { \n
\"column\": \"asins\", \n              \"properties\": { \n                  \"dtype\":
\"category\", \n                  \"num_unique_values\": 24, \n                  \"samples\": [ \n
\"B018T075DC\", \n                  \"B018Y22BI4\", \n                  \"B01AHB9CN2\" \n                  ], \n                  \"semantic_type\": \"\", \n
\"description\": \"\" \n              }, \n          { \n              \"column\":
\"brand\", \n              \"properties\": { \n                  \"dtype\": \"category\", \n
\"num_unique_values\": 4, \n                  \"samples\": [ \n
\"Amazon Fire\", \n                  \"Amazon Fire Tv\", \n                  \"Amazon\" \n                  ], \n                  \"semantic_type\": \"\", \n
\"description\": \"\" \n              }, \n          { \n              \"column\":
\"product_categories\", \n              \"properties\": { \n                  \"dtype\":
\"object\", \n                  \"semantic_type\": \"\", \n                  \"description\": \"\" \n              }, \n          { \n              \"column\":
\"product_keys\", \n              \"properties\": { \n                  \"dtype\":
\"category\", \n                  \"num_unique_values\": 24, \n                  \"samples\": [ \n
\"amazon/b018t075dc,firehd8tabletwithalexa8hddisplay16gbtangerinewiths
pecialoffers/
5620410,firehd8tabletwithalexa8hddisplay16gbtangerinewithspecialoffers
/b018t075dc,841667103068,0841667103068\", \n
\"amazonfire16gb5thgen2015releaseblack/272201222631,amazonfire16gb5thg
en2015releaseblack/
b018y22bi4,841667103129,0841667103129,amazonfire16gb5thgen2015releaseb
lack/5023200,amazonfire16gb5thgen2015releaseblack/
332273296844,amazonfire16gb5thgen2015releaseblack/
232443003172,amazon/b018y22bi4\", \n
\"841667104676,amazon/53004484,amazon/b01ahb9cn2,0841667104676,allnewf

```

```
irehd8tablet8hddisplaywifil6gbincludesspecialoffersmagenta/
5620406,allnewfirehd8tablet8hddisplaywifil6gbincludesspecialoffersmage
nta/b01ahb9cn2"\n      ],\n      \"semantic_type\": \"\", \n
\"description\": \"\"\n    }\n  },\n  {\n    \"column\":
\"manufacturer_name\", \n    \"properties\": {\n      \"dtype\":
\"category\", \n      \"num_unique_values\": 1, \n      \"samples\":
[\n        \"Amazon\", \n      ], \n      \"semantic_type\": \"\", \n
    \"description\": \"\"\n    }\n  },\n  {\n    \"column\": \"review_do_recommend\", \n    \"properties\": {\n
\"dtype\": \"category\", \n      \"num_unique_values\": 2, \n
\"samples\": [\n        false\n      ], \n      \"semantic_type\": \"\", \n
    \"description\": \"\"\n    }\n  },\n  {\n    \"column\": \"review_num_helpful\", \n    \"properties\": {\n
\"dtype\": \"number\", \n      \"std\": 2.194084771528854, \n
\"min\": 0.0, \n      \"max\": 109.0, \n
\"num_unique_values\": 57, \n      \"samples\": [\n        0.0\n      ], \n      \"semantic_type\": \"\", \n
    \"description\": \"\"\n    }\n  },\n  {\n    \"column\": \"review_rating\", \n    \"properties\": {\n
\"dtype\": \"number\", \n      \"std\": 0.7217255917862178, \n
\"min\": 1.0, \n      \"max\": 5.0, \n
\"num_unique_values\": 5, \n      \"samples\": [\n        4.0\n      ], \n      \"semantic_type\": \"\", \n
    \"description\": \"\"\n    }\n  },\n  {\n    \"column\": \"review_username\", \n    \"properties\": {\n
\"dtype\": \"string\", \n      \"num_unique_values\": 26309, \n      \"samples\": [\n
\"RED3\", \n      ], \n      \"semantic_type\": \"\", \n
    \"description\": \"\"\n    }\n  },\n  {\n    \"column\":
\"review_length\", \n    \"properties\": {\n      \"dtype\":
\"number\", \n      \"std\": 167, \n      \"min\": 6, \n
\"max\": 8351, \n      \"num_unique_values\": 984, \n
\"samples\": [\n        1479\n      ], \n      \"semantic_type\":
\"\", \n      \"description\": \"\"\n    }\n  },\n  {\n    \"column\":
\"review_did_purchase\", \n    \"properties\": {\n      \"dtype\":
\"boolean\", \n      \"num_unique_values\": 1, \n
\"samples\": [\n        false\n      ], \n      \"semantic_type\": \"\", \n
    \"description\": \"\"\n    }\n  },\n  {\n    \"column\": \"total_votes\", \n    \"properties\": {\n
\"dtype\": \"number\", \n      \"std\": 0, \n
\"min\": 0, \n      \"max\": 0, \n
\"num_unique_values\": 1, \n      \"samples\": [\n        0\n      ], \n      \"semantic_type\": \"\", \n
    \"description\": \"\"\n    }\n  },\n  {\n    \"column\": \"clean_text\", \n    \"properties\": {\n
\"dtype\": \"string\", \n      \"num_unique_values\": 33957, \n      \"samples\": [\n
\"im currently love item setting reminders also playing music spotify bluetooth\", \n      ], \n      \"semantic_type\": \"\", \n
    \"description\": \"\"\n    }\n  },\n  {\n    \"column\":
\"clean_title\", \n    \"properties\": {\n      \"dtype\":
\"category\", \n      \"num_unique_values\": 13899, \n
```

```

\"samples\": [\n          \"hate\"\n        ],\n\"semantic_type\": \"\", \n          \"description\": \"\" \n        }\n      ]\n    }, \"type\": \"dataframe\", \"variable_name\": \"data\"}

# Rename the columns with the original column names
data.rename(columns={'clean_text': 'review_text', 'clean_title':
'review_title'}, inplace=True)

# Display the new DataFrame
data.head(1)

{"summary": "{\n  \"name\": \"data\",\n  \"rows\": 34054,\n  \"fields\": [\n    {\n      \"column\": \"review_date\",\n      \"properties\": {\n        \"dtype\": \"date\",\n        \"min\": \"2014-10-24 00:00:00+00:00\",\n        \"max\": \"2018-04-18 00:00:00+00:00\",\n        \"num_unique_values\": 941,\n        \"samples\": [\n          \"2015-11-29 00:00:00+00:00\",\n          \"2016-03-14 00:00:00+00:00\",\n          \"2016-08-21 00:00:00+00:00\"\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\" \n      },\n      \"column\": \"id\",\n      \"properties\": {\n        \"dtype\": \"category\",\n        \"num_unique_values\": 24,\n        \"samples\": [\n          \"AVqkIhxunnc1JgDc3kg_\",\n          \"AVpgdkC8ilAPnD_xsvyi\",\n          \"AVqkIhwDv8e3D10-lebb\"\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\" \n      },\n      \"column\": \"asins\",\n      \"properties\": {\n        \"dtype\": \"category\",\n        \"num_unique_values\": 24,\n        \"samples\": [\n          \"B018T075DC\",\n          \"B018Y22BI4\",\n          \"B01AHB9CN2\"\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\" \n      },\n      \"column\": \"brand\",\n      \"properties\": {\n        \"dtype\": \"category\",\n        \"num_unique_values\": 4,\n        \"samples\": [\n          \"Amazon Fire\",\n          \"Amazon Fire Tv\",\n          \"Amazon\"\n        ],\n        \"semantic_type\": \"\",\n        \"description\": \"\" \n      },\n      \"column\": \"product_categories\",\n      \"properties\": {\n        \"dtype\": \"object\",\n        \"semantic_type\": \"\",\n        \"description\": \"\" \n      },\n      \"column\": \"product_keys\",\n      \"properties\": {\n        \"dtype\": \"category\",\n        \"num_unique_values\": 24,\n        \"samples\": [\n          \"amazon/b018t075dc,firehd8tabletwithalexa8hddisplay16gbtangerinewithspecialoffers/5620410,firehd8tabletwithalexa8hddisplay16gbtangerinewithspecialoffers/b018t075dc,841667103068,0841667103068\",\n          \"amazonfire16gb5thgen2015releaseblack/272201222631,amazonfire16gb5thgen2015releaseblack/b018y22bi4,841667103129,0841667103129,amazonfire16gb5thgen2015releaseblack/5023200,amazonfire16gb5thgen2015releaseblack/332273296844,amazonfire16gb5thgen2015releaseblack/\"

```



```

232443003172,amazon/b018y22bi4",\n
\"841667104676,amazon/53004484,amazon/b01ahb9cn2,0841667104676,allnewf
irehd8tablet8hddisplaywifil6gbincludesspecialoffersmagenta/
5620406,allnewfirehd8tablet8hddisplaywifil6gbincludesspecialoffersmage
nta/b01ahb9cn2\",,\n      \"semantic_type\": \"\",,\n
\"description\": \"\",,\n      },,\n      {\n      \"column\":
\"manufacturer_name\",,\n      \"properties\": {\n      \"dtype\":
\"category\",,\n      \"num_unique_values\": 1,\n      \"samples\":
[\n      \"Amazon\",,\n      ],,\n      \"semantic_type\": \"\",,\n
      \"description\": \"\",,\n      },,\n      {\n
\"column\": \"review_do_recommend\",,\n      \"properties\": {\n
\"dtype\": \"category\",,\n      \"num_unique_values\": 2,\n
\"samples\": [\n      false,\n      ],,\n
\"semantic_type\": \"\",,\n      \"description\": \"\",,\n
      },,\n      {\n      \"column\": \"review_num_helpful\",,\n
\"properties\": {\n      \"dtype\": \"number\",,\n      \"std\":
2.194084771528854,\n      \"min\": 0.0,\n      \"max\": 109.0,\n
\"num_unique_values\": 57,\n      \"samples\": [\n      0.0,\n
      ],,\n      \"semantic_type\": \"\",,\n      \"description\": \"\",,\n
      },,\n      {\n      \"column\": \"review_rating\",,\n
\"properties\": {\n      \"dtype\": \"number\",,\n      \"std\":
0.7217255917862178,\n      \"min\": 1.0,\n      \"max\": 5.0,\n
\"num_unique_values\": 5,\n      \"samples\": [\n      4.0,\n
      ],,\n      \"semantic_type\": \"\",,\n      \"description\": \"\",,\n
      },,\n      {\n      \"column\": \"review_username\",,\n
\"properties\": {\n      \"dtype\": \"string\",,\n
\"num_unique_values\": 26309,\n      \"samples\": [\n
\"RED3\",,\n      ],,\n      \"semantic_type\": \"\",,\n
\"description\": \"\",,\n      },,\n      {\n      \"column\":
\"review_length\",,\n      \"properties\": {\n      \"dtype\":
\"number\",,\n      \"std\": 167,\n      \"min\": 6,\n
\"max\": 8351,\n      \"num_unique_values\": 984,\n
\"samples\": [\n      1479,\n      ],,\n      \"semantic_type\":
\",,\n      \"description\": \"\",,\n      },,\n      {\n
\"column\": \"review_did_purchase\",,\n      \"properties\": {\n
\"dtype\": \"boolean\",,\n      \"num_unique_values\": 1,\n
\"samples\": [\n      false,\n      ],,\n
\"semantic_type\": \"\",,\n      \"description\": \"\",,\n
      },,\n      {\n      \"column\": \"total_votes\",,\n
\"properties\": {\n      \"dtype\": \"number\",,\n      \"std\":
0,\n      \"min\": 0,\n      \"max\": 0,\n
\"num_unique_values\": 1,\n      \"samples\": [\n      0,\n
      ],,\n      \"semantic_type\": \"\",,\n      \"description\": \"\",,\n
      },,\n      {\n      \"column\": \"review_text\",,\n
\"properties\": {\n      \"dtype\": \"string\",,\n
\"num_unique_values\": 33957,\n      \"samples\": [\n      \"im
currently love item setting reminders also playing music spotify
bluetooth\",,\n      ],,\n      \"semantic_type\": \"\",,\n
\"description\": \"\",,\n      },,\n      {\n      \"column\":

```

```

{"review_title","\n      \n"properties\": {\n      \n"dtype\":
{"category","\n      \n"num_unique_values\": 13899,\n
{"samples\": [\n      \n"hate"\n      ],\n
{"semantic_type\": "\"",\n      \n"description\": "\""\n      }\n      }\n      ]\n      }","type":"dataframe","variable_name":"data"}

# Download NLTK WordNet
nltk.download('wordnet')

[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Package wordnet is already up-to-date!

True

```

- We will now perform lemmatization, which reduces words to their base form while still preserving their meaning to ensure consistency and improve the accuracy of our analysis.

```

# Initialize the WordNet lemmatizer
lemmatizer = WordNetLemmatizer()

# Initialize the WordNet lemmatizer
lemmatizer = WordNetLemmatizer()

# Function to perform lemmatization on text
def lemmatize_text(text):
    words = text.split()

    # Lemmatization
    lemmatized_words = [lemmatizer.lemmatize(word) for word in words]

    return ' '.join(lemmatized_words)

# Apply lemmatization to review_text and review_title separately
data['lemmatized_text'] = data['review_text'].apply(lemmatize_text)
data['lemmatized_title'] = data['review_title'].apply(lemmatize_text)

# Display the lemmatized text along with original columns
data[['review_text', 'review_title', 'lemmatized_text',
'lemmatized_title']]

{"summary":{"\n  \n"name\": \"data[['review_text', 'review_title',
'lemmatized_text', 'lemmatized_title']]\",\n  \n"rows\": 34054,\n
\n"fields\": [\n    {\n      \n"column\": \"review_date\",\n
\n"properties\": {\n      \n"dtype\": \"date\",\n      \n"min\":
\n"2014-10-24 00:00:00+00:00\",\n      \n"max\": \"2018-04-18
00:00:00+00:00\",\n      \n"num_unique_values\": 941,\n
\n"samples\": [\n      \n"2015-11-29 00:00:00+00:00\",\n
\n"2016-03-14 00:00:00+00:00\",\n      \n"2016-08-21
00:00:00+00:00"\n      ],\n      \n"semantic_type\": "\"",\n
\n"description\": "\""\n      }\n      },\n    {\n      \n"column\":
\n"review_text\",\n      \n"properties\": {\n      \n"dtype\":

```

```

\"string\", \n          \"num_unique_values\": 33957, \n
\"samples\": [\n          \"im currently love item setting reminders
also playing music spotify bluetooth\", \n          \"nice size nice
looklove feel turning page great reading experience\", \n
\"good quality android tablet also makes nice gift\" \n          ], \n
\"semantic_type\": \"\", \n          \"description\": \"\" \n          } \n
n      }, \n      { \n          \"column\": \"review_title\", \n
\"properties\": { \n          \"dtype\": \"category\", \n
\"num_unique_values\": 13899, \n          \"samples\": [\n
\"hate\", \n          \"great purchase first time kindle buyer\", \n
\"great sounding musicand much\" \n          ], \n
\"semantic_type\": \"\", \n          \"description\": \"\" \n          } \n
n      }, \n      { \n          \"column\": \"lemmatized_text\", \n
\"properties\": { \n          \"dtype\": \"string\", \n
\"num_unique_values\": 33952, \n          \"samples\": [\n
\"soul purpose buying item read completely satisfied howevef camera
quality isnt great quality\", \n          \"want ereader reading
book\", \n          \"love kindle best one yet ive owned original
kindle fire fire hd bought paper white preferred reading actual book
ereaders ereader convenient deal book travel need find place store im
advocate paper white like reading book\" \n          ], \n
\"semantic_type\": \"\", \n          \"description\": \"\" \n          } \n
n      }, \n      { \n          \"column\": \"lemmatized_title\", \n
\"properties\": { \n          \"dtype\": \"category\", \n
\"num_unique_values\": 13649, \n          \"samples\": [\n
\"great ebook reader around house free wifi\", \n          \"great gift
granddaughter\", \n          \"bad child\" \n          ], \n
\"semantic_type\": \"\", \n          \"description\": \"\" \n          } \n
n      } \n      ] \n      }, \"type\": \"dataframe\"}

```

dropping the columns not lemmatized

```

data.drop(columns = ['review_text', 'review_title'] , inplace = True)
data.head(1)

```

```

{\"summary\": { \n      \"name\": \"data\", \n      \"rows\": 34054, \n
\"fields\": [\n          { \n          \"column\": \"review_date\", \n
\"properties\": { \n          \"dtype\": \"date\", \n          \"min\":
\"2014-10-24 00:00:00+00:00\", \n          \"max\": \"2018-04-18
00:00:00+00:00\", \n          \"num_unique_values\": 941, \n
\"samples\": [\n          \"2015-11-29 00:00:00+00:00\", \n
\"2016-03-14 00:00:00+00:00\", \n          \"2016-08-21
00:00:00+00:00\" \n          ], \n          \"semantic_type\": \"\", \n
\"description\": \"\" \n          } \n          }, \n          { \n          \"column\":
\"id\", \n          \"properties\": { \n          \"dtype\": \"category\", \n
\"num_unique_values\": 24, \n          \"samples\": [\n
\"AVqkIhxunnc1JgDc3kg_\", \n          \"AVpgdkC8ilAPnD_xsvyi\", \n
\"AVqkIhwDv8e3D10-lebb\" \n          ], \n          \"semantic_type\":
\"\", \n          \"description\": \"\" \n          } \n          }, \n          { \n
\"column\": \"asins\", \n          \"properties\": { \n          \"dtype\":

```

```

\"category\",\\n          \"num_unique_values\": 24,\\n
\"samples\": [\\n          \"B018T075DC\",\\n          \"B018Y22BI4\",\\n
\"B01AHB9CN2\",\\n          ],\\n          \"semantic_type\": \"\",\\n
\"description\": \"\"\\n          }\\n          },\\n          {\\n          \"column\":
\"brand\",\\n          \"properties\": {\\n          \"dtype\": \"category\",\\n
          \"num_unique_values\": 4,\\n          \"samples\": [\\n
\"Amazon Fire\",\\n          \"Amazon Fire Tv\",\\n          \"Amazon\"\\n
          ],\\n          \"semantic_type\": \"\",\\n
\"description\": \"\"\\n          }\\n          },\\n          {\\n          \"column\":
\"product_categories\",\\n          \"properties\": {\\n          \"dtype\":
\"object\",\\n          \"semantic_type\": \"\",\\n
\"description\": \"\"\\n          }\\n          },\\n          {\\n          \"column\":
\"product_keys\",\\n          \"properties\": {\\n          \"dtype\":
\"category\",\\n          \"num_unique_values\": 24,\\n
\"samples\": [\\n
\"amazon/b018t075dc,firehd8tabletwithalexa8hddisplay16gbtangerinewiths
pecialoffers/
5620410,firehd8tabletwithalexa8hddisplay16gbtangerinewithspecialoffers
/b018t075dc,841667103068,0841667103068\",\\n
\"amazonfire16gb5thgen2015releaseblack/272201222631,amazonfire16gb5thg
en2015releaseblack/
b018y22bi4,841667103129,0841667103129,amazonfire16gb5thgen2015releaseb
lack/5023200,amazonfire16gb5thgen2015releaseblack/
332273296844,amazonfire16gb5thgen2015releaseblack/
232443003172,amazon/b018y22bi4\",\\n
\"841667104676,amazon/53004484,amazon/b01ahb9cn2,0841667104676,allnewf
irehd8tablet8hddisplaywifil6gbincludesspecialoffersmagenta/
5620406,allnewfirehd8tablet8hddisplaywifil6gbincludesspecialoffersmage
nta/b01ahb9cn2\"\\n          ],\\n          \"semantic_type\": \"\",\\n
\"description\": \"\"\\n          }\\n          },\\n          {\\n          \"column\":
\"manufacturer_name\",\\n          \"properties\": {\\n          \"dtype\":
\"category\",\\n          \"num_unique_values\": 1,\\n          \"samples\":
[\\n          \"Amazon\"\\n          ],\\n          \"semantic_type\": \"\",\\n
          \"description\": \"\"\\n          }\\n          },\\n          {\\n
          \"column\": \"review_do_recommend\",\\n          \"properties\": {\\n
          \"dtype\": \"category\",\\n          \"num_unique_values\": 2,\\n
          \"samples\": [\\n          false\\n          ],\\n
          \"semantic_type\": \"\",\\n          \"description\": \"\"\\n          }\\n
          },\\n          {\\n          \"column\": \"review_num_helpful\",\\n
          \"properties\": {\\n          \"dtype\": \"number\",\\n          \"std\":
2.194084771528854,\\n          \"min\": 0.0,\\n          \"max\": 109.0,\\n
          \"num_unique_values\": 57,\\n          \"samples\": [\\n          0.0\\n
          ],\\n          \"semantic_type\": \"\",\\n          \"description\": \"\"\\n
          }\\n          },\\n          {\\n          \"column\": \"review_rating\",\\n
          \"properties\": {\\n          \"dtype\": \"number\",\\n          \"std\":
0.7217255917862178,\\n          \"min\": 1.0,\\n          \"max\": 5.0,\\n
          \"num_unique_values\": 5,\\n          \"samples\": [\\n          4.0\\n
          ],\\n          \"semantic_type\": \"\",\\n          \"description\": \"\"\\n
          }\\n          },\\n          {\\n          \"column\": \"review_username\",\\n

```



```

\"AVqkIhwDv8e3D10-lebb\"\\n      ],\\n      \"semantic_type\\\":
\\\"\",\\n      \"description\\\": \\\"\\\"\\n      }\\n      },\\n      {\\n
\"column\\\": \"asins\\\",\\n      \"properties\\\": {\\n      \"dtype\\\":
\"category\\\",\\n      \"num_unique_values\\\": 24,\\n
\"samples\\\": [\\n      \"B018T075DC\\\",\\n      \"B018Y22BI4\\\",\\n
\"B01AHB9CN2\\\"\\n      ],\\n      \"semantic_type\\\": \\\"\\\",\\n
\"description\\\": \\\"\\\"\\n      }\\n      },\\n      {\\n      \"column\\\":
\"brand\\\",\\n      \"properties\\\": {\\n      \"dtype\\\": \"category\\\",\\n
      \"num_unique_values\\\": 4,\\n      \"samples\\\": [\\n
\"Amazon Fire\\\",\\n      \"Amazon Fire Tv\\\",\\n      \"Amazon\\\"
\\n      ],\\n      \"semantic_type\\\": \\\"\\\",\\n
\"description\\\": \\\"\\\"\\n      }\\n      },\\n      {\\n      \"column\\\":
\"product_categories\\\",\\n      \"properties\\\": {\\n      \"dtype\\\":
\"object\\\",\\n      \"semantic_type\\\": \\\"\\\",\\n
\"description\\\": \\\"\\\"\\n      }\\n      },\\n      {\\n      \"column\\\":
\"product_keys\\\",\\n      \"properties\\\": {\\n      \"dtype\\\":
\"category\\\",\\n      \"num_unique_values\\\": 24,\\n
\"samples\\\": [\\n
\"amazon/b018t075dc,firehd8tabletwithalexa8hddisplay16gbtangerinewiths
pecialoffers/
5620410,firehd8tabletwithalexa8hddisplay16gbtangerinewithspecialoffers
/b018t075dc,841667103068,0841667103068\\\",\\n
\"amazonfire16gb5thgen2015releaseblack/272201222631,amazonfire16gb5thg
en2015releaseblack/
b018y22bi4,841667103129,0841667103129,amazonfire16gb5thgen2015releaseb
lack/5023200,amazonfire16gb5thgen2015releaseblack/
332273296844,amazonfire16gb5thgen2015releaseblack/
232443003172,amazon/b018y22bi4\\\",\\n
\"841667104676,amazon/53004484,amazon/b01ahb9cn2,0841667104676,allnewf
irehd8tablet8hddisplaywifil6gbincludesspecialoffersmagenta/
5620406,allnewfirehd8tablet8hddisplaywifil6gbincludesspecialoffersmage
nta/b01ahb9cn2\\\"\\n      ],\\n      \"semantic_type\\\": \\\"\\\",\\n
\"description\\\": \\\"\\\"\\n      }\\n      },\\n      {\\n      \"column\\\":
\"manufacturer_name\\\",\\n      \"properties\\\": {\\n      \"dtype\\\":
\"category\\\",\\n      \"num_unique_values\\\": 1,\\n      \"samples\\\":
[\\n      \"Amazon\\\"\\n      ],\\n      \"semantic_type\\\": \\\"\\\",\\n
      \"description\\\": \\\"\\\"\\n      }\\n      },\\n      {\\n
\"column\\\": \"review_do_recommend\\\",\\n      \"properties\\\": {\\n
\"dtype\\\": \"category\\\",\\n      \"num_unique_values\\\": 2,\\n
\"samples\\\": [\\n      false\\n      ],\\n
\"semantic_type\\\": \\\"\\\",\\n      \"description\\\": \\\"\\\"\\n      }\\n
      },\\n      {\\n      \"column\\\": \"review_num_helpful\\\",\\n
\"properties\\\": {\\n      \"dtype\\\": \"number\\\",\\n      \"std\\\":
2.194084771528854,\\n      \"min\\\": 0.0,\\n      \"max\\\": 109.0,\\n
\"num_unique_values\\\": 57,\\n      \"samples\\\": [\\n      0.0\\n
      ],\\n      \"semantic_type\\\": \\\"\\\",\\n      \"description\\\": \\\"\\\"\\n
      }\\n      },\\n      {\\n      \"column\\\": \"review_rating\\\",\\n
\"properties\\\": {\\n      \"dtype\\\": \"number\\\",\\n      \"std\\\":
0.7217255917862178,\\n      \"min\\\": 1.0,\\n      \"max\\\": 5.0,\\n

```

```

{"num_unique_values": 5, "samples": [4.0],
 "semantic_type": "",
 "description": ""}
{"column": "review_username",
 "properties": {"dtype": "string",
 "num_unique_values": 26309, "samples": ["RED3",
 "description": ""}
 "column": "review_length",
 "properties": {"dtype": "number",
 "std": 167, "min": 6, "max": 8351,
 "num_unique_values": 984, "samples": [1479,
 "semantic_type": "",
 "description": ""}
 "column": "review_did_purchase",
 "properties": {"dtype": "boolean",
 "num_unique_values": 1, "samples": [false,
 "semantic_type": "",
 "description": ""}
 "column": "total_votes",
 "properties": {"dtype": "number",
 "std": 0, "min": 0, "max": 0,
 "num_unique_values": 1, "samples": [0,
 "semantic_type": "",
 "description": ""}
 "column": "review_text",
 "properties": {"dtype": "string",
 "num_unique_values": 33952, "samples": ["soul purpose buying item read completely satisfied howevef camera
quality isnt great quality",
 "semantic_type": "",
 "description": ""}
 "column": "review_title",
 "properties": {"dtype": "category",
 "num_unique_values": 13649, "samples": ["great ebook reader around house free
wifi",
 "semantic_type": "",
 "description": ""}
 "type": "dataframe", "variable name": "data"}

```

```
# Removing white spaces
```

```
# Function to remove extra spaces from text
```

```
def remove_extra_spaces(text):  
    return ' '.join(text.strip().split())
```

```
# Apply function to the 'lemmatized review text' column
```

```
data['clean text'] = data['review text'].apply(remove_extra_spaces)
```

```
# Apply function to the 'lemmatized review title' column
```

```
data['clean title'] = data['review title'].apply(remove_extra_spaces)
```

```
# Display cleaned text along with original columns
```

```
data[['review text', 'review title', 'clean text', 'clean title']]
```



```
{
  "summary": {
    "name": "data[['review_text',
                  'review_title', 'clean_text', 'clean_title']]",
    "rows": 34054,
    "fields": [
      {
        "column": "review_date",
        "properties": {
          "dtype": "date",
          "min": "2014-10-24 00:00:00+00:00",
          "max": "2018-04-18 00:00:00+00:00",
          "num_unique_values": 941,
          "samples": [
            "2015-11-29 00:00:00+00:00",
            "2016-03-14 00:00:00+00:00",
            "2016-08-21 00:00:00+00:00"
          ],
          "semantic_type": "\"\"",
          "description": "\"\"\"",
          "column": "review_text",
          "properties": {
            "dtype": "string",
            "num_unique_values": 33952,
            "samples": [
              "soul purpose buying item read completely satisfied howevef camera quality isnt great quality",
              "want ereader reading book",
              "love kindle best one yet ive owned original kindle fire fire hd bought paper white preferred reading actual book ereaders ereader convenient deal book travel need find place store im advocate paper white like reading book"
            ],
            "semantic_type": "\"\"",
            "description": "\"\"\"",
            "column": "review_title",
            "properties": {
              "dtype": "category",
              "num_unique_values": 13649,
              "samples": [
                "great ebook reader around house free wifi",
                "great gift granddaughter",
                "bad child"
              ],
              "semantic_type": "\"\"",
              "description": "\"\"\"",
              "column": "clean_text",
              "properties": {
                "dtype": "string",
                "num_unique_values": 33952,
                "samples": [
                  "soul purpose buying item read completely satisfied howevef camera quality isnt great quality",
                  "want ereader reading book",
                  "love kindle best one yet ive owned original kindle fire fire hd bought paper white preferred reading actual book ereaders ereader convenient deal book travel need find place store im advocate paper white like reading book"
                ],
                "semantic_type": "\"\"",
                "description": "\"\"\"",
                "column": "clean_title",
                "properties": {
                  "dtype": "category",
                  "num_unique_values": 13649,
                  "samples": [
                    "great ebook reader around house free wifi",
                    "great gift granddaughter",
                    "bad child"
                  ],
                  "semantic_type": "\"\"",
                  "description": "\"\"\"",
                  "column": "clean_title"
                }
              }
            }
          ]
        }
      }
    ],
    "type": "dataframe"
  }
}
```

Feature Engineering

In the feature engineering section, we process and transform the textual data for further analysis and modeling:

The methods used are;

- **Sentiment Analysis** to determine the sentiment of each review.
- **Visualization with Word Clouds** to visualize the most frequent words in positive and negative reviews
- **Text Vectorization** to convert textual data into numerical form using TF-IDF and Count Vectorization.
- **Word Embedding** to capture the semantic relationships between words by representing them in a continuous vector space.
- **Extraction of Bigrams and Trigrams**

Sentiment Analysis

This was done using the `SentimentIntensityAnalyzer` from the `vaderSentiment` library to calculate a sentiment score for each review.

Each review was labeled with a sentiment score, and reviews were classified as either 'positive' or 'negative' based on this score.

```
from nltk.sentiment.vader import SentimentIntensityAnalyzer

# Download the VADER lexicon
nltk.download('vader_lexicon')

# Initialize the VADER sentiment analyzer
sid = SentimentIntensityAnalyzer()
# Define the sentiment function to calculate the compound score
def sentiment(x):
    score = sid.polarity_scores(x)
    return score['compound']

# Apply the sentiment function to the text column to get sentiment scores
data['sentiment'] = data['clean_text'].apply(lambda x: sentiment(x))

# Print the DataFrame with the sentiment scores
data[['clean_text', 'sentiment', 'review_rating']]

[nltk_data] Downloading package vader_lexicon to /root/nltk_data...
[nltk_data]   Package vader_lexicon is already up-to-date!

{"summary": "{\n  \"name\": \"data[['clean_text', 'sentiment', 'review_rating']]\", \n  \"rows\": 34054, \n  \"fields\": [\n    {\n      \"column\": \"review_date\", \n      \"properties\": {\n        \"dtype\": \"date\", \n        \"min\": \"2014-10-24 00:00:00+00:00\", \n        \"max\": \"2018-04-18 00:00:00+00:00\", \n        \"num_unique_values\": 941, \n        \"samples\": [\n          \"2015-11-29 00:00:00+00:00\", \n          \"2016-03-14 00:00:00+00:00\", \n          \"2016-08-21 00:00:00+00:00\" \n        ], \n        \"semantic_type\": \"date\" \n      } \n    } \n  ] \n}
```

```

{"",\n      "description": "\n      },\n      {\n      "column": "clean_text",\n      "properties": {\n      "dtype": "string",\n      "num_unique_values": 33952,\n      "samples": [\n      "soul purpose buying item read completely\n      satisfied however camera quality isnt great quality",\n      "want ereader reading book",\n      "love kindle best one yet\n      ive owned original kindle fire fire hd bought paper white preferred\n      reading actual book ereaders ereader convenient deal book travel need\n      find place store im advocate paper white like reading book"\n      ],\n      "semantic_type": "",\n      "description": "\n      },\n      {\n      "column":\n      "sentiment",\n      "properties": {\n      "dtype":\n      "number",\n      "std": 0.3331587742851247,\n      "min": -\n      0.9574,\n      "max": 0.9978,\n      "num_unique_values":\n      3076,\n      "samples": [\n      0.9895,\n      0.9686,\n      0.3561\n      ],\n      "semantic_type": "",\n      "description": "\n      },\n      {\n      "column":\n      "review_rating",\n      "properties": {\n      "dtype":\n      "number",\n      "std": 0.7217255917862178,\n      "min":\n      1.0,\n      "max": 5.0,\n      "num_unique_values": 5,\n      "samples": [\n      4.0,\n      3.0,\n      2.0\n      ],\n      "semantic_type": "",\n      "description": "\n      }\n      }\n      ],\n      "type": "dataframe"}

```

```

# Filter the original data DataFrame for negative and positive reviews
negative_reviews_text = data[data['sentiment'].apply(lambda x: 0 <= x
<= 0.6)][['clean_text']]
positive_reviews_text = data[data['sentiment'].apply(lambda x: x >
0.6)][['clean_text']]

```

```

# Create labels for negative and positive reviews
data.loc[data['sentiment'] <= 0.5, 'label'] = 'negative'
data.loc[data['sentiment'] > 0.5, 'label'] = 'positive'

```

```

# Print the updated DataFrame to verify
data[['clean_text', 'sentiment', 'label']]

```

```

{"summary":{\n  "name": "data[['clean_text', 'sentiment',\n  'label']]",\n  "rows": 34054,\n  "fields": [\n    {\n      "column": "review_date",\n      "properties": {\n      "dtype": "date",\n      "min": "2014-10-24 00:00:00+00:00",\n      "max": "2018-04-18 00:00:00+00:00",\n      "num_unique_values": 941,\n      "samples": [\n      "2015-11-29 00:00:00+00:00",\n      "2016-03-14 00:00:00+00:00",\n      "2016-08-21 00:00:00+00:00"\n      ],\n      "semantic_type":\n      "",\n      "description": "\n      },\n      {\n      "column": "clean_text",\n      "properties": {\n      "dtype": "string",\n      "num_unique_values": 33952,\n
```

```

\"samples\": [\n          \"soul purpose buying item read completely\nsatisfied howevef camera quality isnt great quality\", \n\n          \"want ereader reading book\", \n          \"love kindle best one yet\nive owned original kindle fire fire hd bought paper white preferred\nreading actual book ereaders ereader convenient deal book travel need\nfind place store im advocate paper white like reading book\"\n\n          ], \n          \"semantic_type\": \"\", \n          \"description\": \"\" \n          }, \n          {\n          \"column\":\n          \"sentiment\", \n          \"properties\": {\n          \"dtype\":\n          \"number\", \n          \"std\": 0.3331587742851247, \n          \"min\": -\n0.9574, \n          \"max\": 0.9978, \n          \"num_unique_values\":\n3076, \n          \"samples\": [\n          0.9895, \n          0.9686, \n          0.3561\n          ], \n          \"semantic_type\": \"\", \n          \"description\": \"\" \n          }, \n          {\n          \"column\":\n          \"label\", \n          \"properties\": {\n          \"dtype\": \"category\", \n          \"num_unique_values\": 2, \n          \"samples\": [\n          \"negative\", \n          \"positive\"\n          ], \n          \"semantic_type\": \"\", \n          \"description\": \"\" \n          } \n          } \n          ], \n          \"type\": \"dataframe\"}

```

Labelling the reviews using the sentiment scores

- Scores ranging from 0 - 0.5 will be labeled as **negative**
- Scores ranging from 0.6 - 1 will be labeled as **positive**

```

# Filter the original data DataFrame for negative and positive reviews
negative_reviews_text = data[data['sentiment'].apply(lambda x: 0 <= x
<= 0.5)][['clean_text']]
positive_reviews_text = data[data['sentiment'].apply(lambda x: x >
0.5)][['clean_text']]

# Create labels for negative and positive reviews
data.loc[data['sentiment'] <= 0.5, 'label'] = 'negative'
data.loc[data['sentiment'] > 0.5, 'label'] = 'positive'

# Print the updated DataFrame to verify
# Print the DataFrame with the sentiment scores
data[['clean_text', 'sentiment', 'label']]

{"summary": "{\n  \"name\": \"data[['clean_text', 'sentiment',\n'label']]\",\n  \"rows\": 34054,\n  \"fields\": [\n    {\n      \"column\": \"review_date\", \n      \"properties\": {\n        \"dtype\": \"date\", \n        \"min\": \"2014-10-24 00:00:00+00:00\", \n        \"max\": \"2018-04-18 00:00:00+00:00\", \n        \"num_unique_values\": 941, \n        \"samples\": [\n          \"2015-11-29 00:00:00+00:00\", \n          \"2016-03-14 00:00:00+00:00\", \n          \"2016-08-21 00:00:00+00:00\"\n        ], \n        \"semantic_type\":\n        \"\", \n        \"description\": \"\" \n        }, \n        {\n          \"column\": \"clean_text\", \n          \"properties\": {\n

```

```

\ "dtype\ ": \ "string\ ",\n          \ "num_unique_values\ ": 33952,\n
\ "samples\ ": [\n          \ "soul purpose buying item read completely
satisfied howevef camera quality isnt great quality\ ",\n
\ "want ereader reading book\ ",\n          \ "love kindle best one yet
ive owned original kindle fire fire hd bought paper white preferred
reading actual book ereaders ereader convenient deal book travel need
find place store im advocate paper white like reading book\ "\n
n          ],\n          \ "semantic_type\ ": \ "\",\n
\ "description\ ": \ "\",\n          },\n          {\n          \ "column\ ":
\ "sentiment\ ",\n          \ "properties\ ": {\n          \ "dtype\ ":
\ "number\ ",\n          \ "std\ ": 0.3331587742851247,\n          \ "min\ ": -
0.9574,\n          \ "max\ ": 0.9978,\n          \ "num_unique_values\ ":
3076,\n          \ "samples\ ": [\n          0.9895,\n          0.9686,\n
0.3561\n          ],\n          \ "semantic_type\ ": \ "\",\n
\ "description\ ": \ "\",\n          },\n          {\n          \ "column\ ":
\ "label\ ",\n          \ "properties\ ": {\n          \ "dtype\ ": \ "category\ ",\n
n          \ "num_unique_values\ ": 2,\n          \ "samples\ ": [\n
\ "negative\ ",\n          \ "positive\ "\n          ],\n
\ "semantic_type\ ": \ "\",\n          \ "description\ ": \ "\",\n          }\n
n          }\n          ]\n          }", "type": "dataframe"}

```

```

print("Number of negative reviews:", negative_reviews_text.shape[0])
print("Number of positive reviews:", positive_reviews_text.shape[0])

```

```

Number of negative reviews: 6054
Number of positive reviews: 26271

```

- We can observe from this that we have class imbalance.

```

from sklearn.feature_extraction.text import CountVectorizer

# # DataFrame setup
# data = pd.DataFrame({
#     'clean_text': ["I love this product", "This is the worst thing
ever", "Not bad", "Absolutely fantastic", "Terrible experience"],
#     'sentiment': [0.9, 0.2, 0.6, 0.8, 0.3],
# })

# Create labels for negative and positive reviews
data.loc[data['sentiment'] <= 0.5, 'label'] = 'negative'
data.loc[data['sentiment'] > 0.5, 'label'] = 'positive'

# Filter the original data for negative and positive reviews
negative_reviews_text = data[data['sentiment'].apply(lambda x: 0 <= x
<= 0.5)][['clean_text']]
positive_reviews_text = data[data['sentiment'].apply(lambda x: x >
0.5)][['clean_text']]

# Create a CountVectorizer to count word frequencies
vectorizer = CountVectorizer()

```

```

# Fit and transform the 'clean_text' data for negative and positive
reviews
X_negative = vectorizer.fit_transform(negative_reviews_text)
X_positive = vectorizer.fit_transform(positive_reviews_text)

# Sum up the counts of each vocabulary word
word_frequencies_negative = X_negative.sum(axis=0).A1
word_frequencies_positive = X_positive.sum(axis=0).A1

# Create a dictionary of word frequencies
vocab = vectorizer.get_feature_names_out()
word_frequencies_negative = dict(zip(vocab,
word_frequencies_negative))
word_frequencies_positive = dict(zip(vocab,
word_frequencies_positive))

# Create word clouds for negative and positive reviews
wordcloud_negative = WordCloud(width=800, height=400,
background_color='white').generate_from_frequencies(word_frequencies_n
egative)
wordcloud_positive = WordCloud(width=800, height=400,
background_color='white').generate_from_frequencies(word_frequencies_p
ositive)

# Display the word clouds in separate figures
plt.figure(figsize=(10, 5))
plt.imshow(wordcloud_negative, interpolation='bilinear')
plt.title('Negative Reviews')
plt.axis('off')
plt.show()

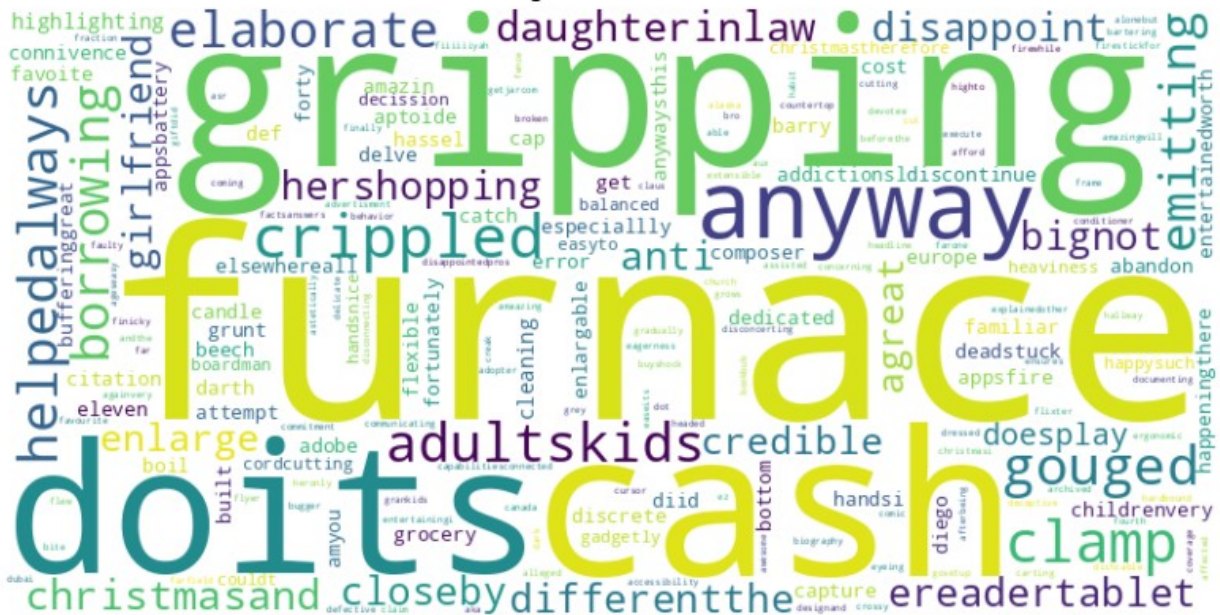
print() # Separating the word clouds display for clarity

plt.figure(figsize=(10, 5))
plt.imshow(wordcloud_positive, interpolation='bilinear')
plt.title('Positive Reviews')
plt.axis('off')
plt.show()

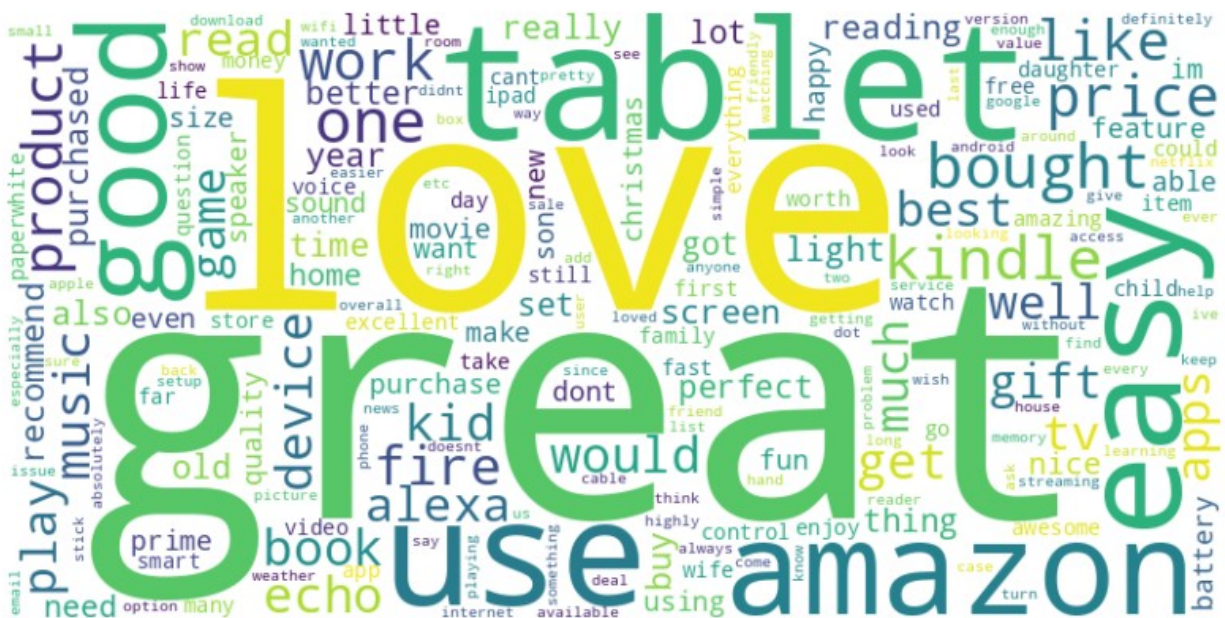
# Add a sentiment_label column for the countplot
data['sentiment_label'] = data['label']

```


Negative Reviews



Positive Reviews



```
# Previewing our column names
data.columns
```

```
Index(['id', 'asins', 'brand', 'product_categories', 'product_keys',
      'manufacturer_name', 'review_do_recommend',
      'review_num_helpful',
```

```

        'review_rating', 'review_username', 'review_length',
        'review_did_purchase', 'total_votes', 'review_text',
'review_title',
        'clean_text', 'clean_title', 'sentiment', 'label',
'sentiment_label'],
        dtype='object')

```

- Let's visualize the distribution of sentiment scores and review ratings.
- We will now convert our labels into numerical data for modeling

```
# Perform label encoding
```

```
label_encoder = LabelEncoder()
```

```
data['labeled'] = label_encoder.fit_transform(data['label'])
```

```
print(data[['clean_text', 'sentiment', 'labeled']])
```

```
clean_text \
review_date
```

```

2017-01-13 00:00:00+00:00  product far disappointed child love use
like a...
2017-01-13 00:00:00+00:00  great beginner experienced person bought
gift ...
2017-01-13 00:00:00+00:00  inexpensive tablet use learn step nabi
thrille...
2017-01-13 00:00:00+00:00  ive fire hd two week love tablet great
valuewe...
2017-01-12 00:00:00+00:00  bought grand daughter come visit set user
ente...
...
...
2016-05-07 00:00:00+00:00      able stream tv movie around world work
great
2016-05-07 00:00:00+00:00      best streaming device portable amazing
picture
2016-05-07 00:00:00+00:00  simply best watch tv series movie work even
be...
2016-07-05 00:00:00+00:00  looking way cut cost raising cable bill
friend...
2015-12-03 00:00:00+00:00      enjoy kindle tv beat paying cable every
month

```

	sentiment	labeled
review_date		
2017-01-13 00:00:00+00:00	0.8126	1
2017-01-13 00:00:00+00:00	0.9042	1
2017-01-13 00:00:00+00:00	0.4404	0
2017-01-13 00:00:00+00:00	0.9899	1
2017-01-12 00:00:00+00:00	0.9371	1

...
2016-05-07 00:00:00+00:00	0.6249	1
2016-05-07 00:00:00+00:00	0.8402	1
2016-05-07 00:00:00+00:00	0.9022	1
2016-07-05 00:00:00+00:00	0.6808	1
2015-12-03 00:00:00+00:00	0.4939	0

[34054 rows x 3 columns]

Feature Extraction

- In this step, we will extract bigrams from the text data and analyze their frequency.

#Extraction of Bigrams

Function to generate n-grams

```
from collections import defaultdict
from nltk import ngrams # Import the ngrams function
```

Function to generate n-grams

```
def generate_ngrams(clean_text, n):
    words = clean_text.split()
    return list(ngrams(words, n))
```

Initialize a defaultdict for frequency counts

```
freq_dict = defaultdict(int)
```

Calculate bigram frequency

```
for sent in data["clean_text"]:
    for word in generate_ngrams(sent, 2):
        freq_dict[word] += 1
```

Sort the frequency dictionary and create a DataFrame

```
fd_sorted = pd.DataFrame(sorted(freq_dict.items(), key=lambda x: x[1],
reverse=True))
fd_sorted.columns = ["word", "wordcount"]
print(fd_sorted.head(25))
```

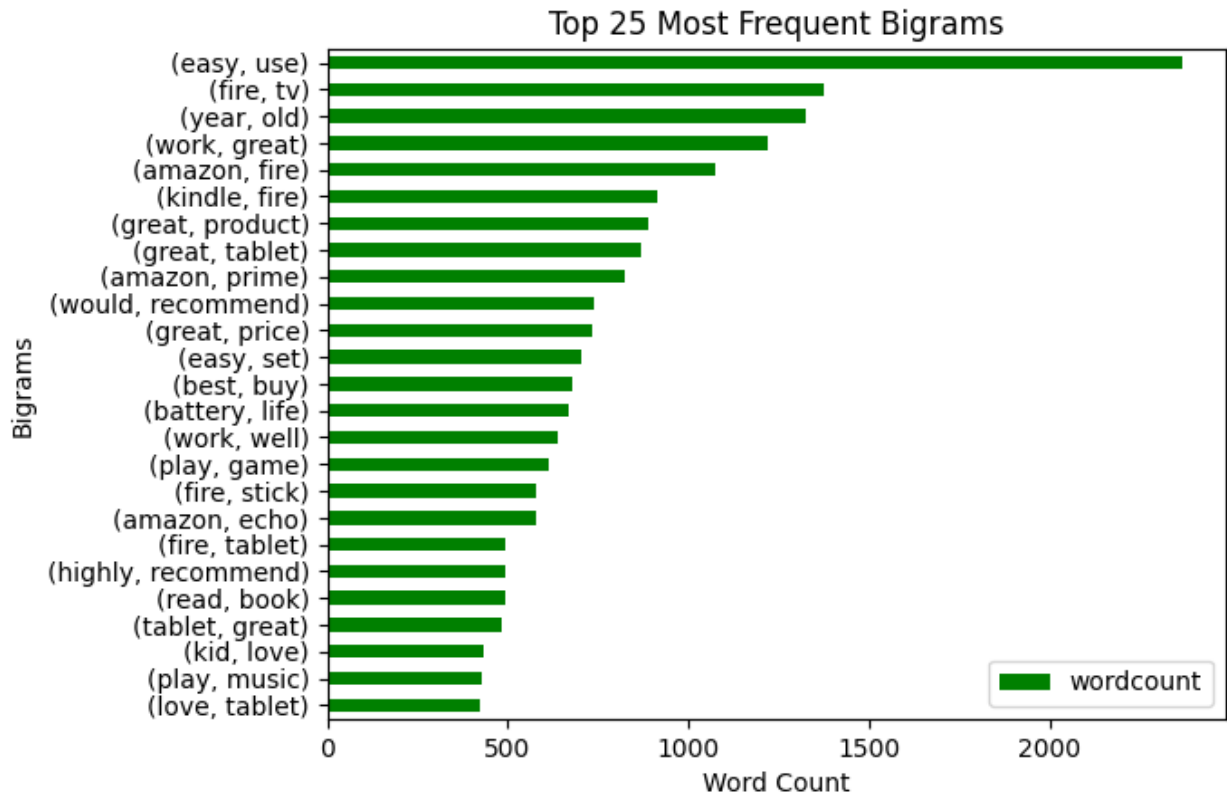
	word	wordcount
0	(easy, use)	2367
1	(fire, tv)	1373
2	(year, old)	1327
3	(work, great)	1219
4	(amazon, fire)	1076
5	(kindle, fire)	916
6	(great, product)	889
7	(great, tablet)	871
8	(amazon, prime)	823
9	(would, recommend)	738
10	(great, price)	732
11	(easy, set)	706
12	(best, buy)	679

13	(battery, life)	667
14	(work, well)	639
15	(play, game)	614
16	(fire, stick)	580
17	(amazon, echo)	579
18	(fire, tablet)	495
19	(highly, recommend)	495
20	(read, book)	492
21	(tablet, great)	485
22	(kid, love)	435
23	(play, music)	429
24	(love, tablet)	423

- Let's visualize the top 25 most frequent bigrams

```
# Function to plot a horizontal bar chart
def horizontal_bar_chart(data, color):
    data.plot(kind='barh', x='word', y='wordcount', color=color)
    plt.xlabel('Word Count')
    plt.ylabel('Bigrams')
    plt.title('Top 25 Most Frequent Bigrams')
    plt.gca().invert_yaxis() # Invert y-axis to have the highest
count on top
    plt.show()

# Plot the top 25 most frequent bigrams
horizontal_bar_chart(fd_sorted.head(25), 'green')
```



#Extraction of Trigrams

```
# Calculate trigram frequency
for sent in data["clean_text"]:
    for word in generate_ngrams(sent,3):
        freq_dict[word] += 1
# Sort the frequency dictionary and create a DataFrame
fd_sorted = pd.DataFrame(sorted(freq_dict.items(), key=lambda x: x[1],
reverse=True))
fd_sorted.columns = ["word", "wordcount"]
print(fd_sorted.head(25))
```

	word	wordcount
0	(easy, use)	2367
1	(fire, tv)	1373
2	(year, old)	1327
3	(work, great)	1219
4	(amazon, fire)	1076
5	(kindle, fire)	916
6	(great, product)	889
7	(great, tablet)	871
8	(amazon, prime)	823
9	(would, recommend)	738
10	(great, price)	732

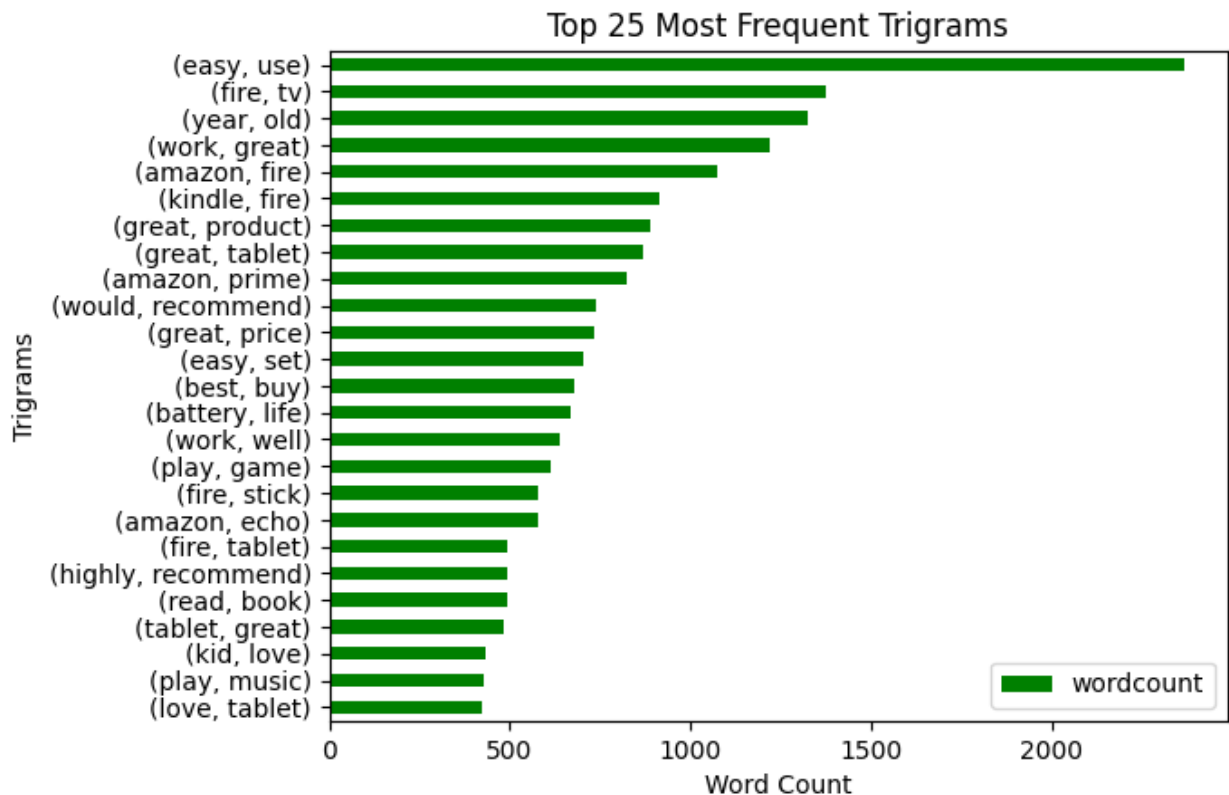
11	(easy, set)	706
12	(best, buy)	679
13	(battery, life)	667
14	(work, well)	639
15	(play, game)	614
16	(fire, stick)	580
17	(amazon, echo)	579
18	(fire, tablet)	495
19	(highly, recommend)	495
20	(read, book)	492
21	(tablet, great)	485
22	(kid, love)	435
23	(play, music)	429
24	(love, tablet)	423

Function to plot a horizontal bar chart

```
def horizontal_bar_chart(data, color):  
    data.plot(kind='barh', x='word', y='wordcount', color=color)  
    plt.xlabel('Word Count')  
    plt.ylabel('Trigrams')  
    plt.title('Top 25 Most Frequent Trigrams')  
    plt.gca().invert_yaxis() # Invert y-axis to have the highest  
count on top  
    plt.show()
```

Plot the top 25 most frequent trigrams

```
horizontal_bar_chart(fd_sorted.head(25), 'green')
```



Word Vectorization

Methods used are:

TF-IDF Vectorization

The TF-IDF (Term Frequency-Inverse Document Frequency) vectorizer transforms the text into a weighted matrix, where each term's importance is adjusted based on its frequency in the document and across all documents.

Count Vectorization

The Count Vectorizer converts the text into a matrix of token counts, representing the raw frequency of each term.

The result

Two matrices one with TF-IDF weights and another with raw token counts, each representing the reviews in a numerical format.

```
from sklearn.feature_extraction.text import CountVectorizer

clean_text = data['clean_text']

# Initialize CountVectorizer
vectorizer = CountVectorizer()
```

```

# Fit and transform the clean_text column
X_count = vectorizer.fit_transform(clean_text)

# Print the array representation of the features
print(X_count.toarray()[1:])

[[0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 ...
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]]

# CountVectorizer
count_vec = CountVectorizer()
# Convert the Pandas Series to a list of strings
X_count = count_vec.fit_transform(clean_text.tolist())
print('CountVectorizer:')
print(count_vec.get_feature_names_out()[:10], '\n')

CountVectorizer:
['aa' 'aaa' 'aamazon' 'aand' 'abandon' 'abandoned' 'abandoning'
 'abattery'
 'abc' 'abcmouse']

```

- We extracted the first 10 feature names

Next is the TF-IDF Vectorizer

```

from sklearn.feature_extraction.text import TfidfVectorizer

#Initialize the TfidfVectorizer
vectorizer = TfidfVectorizer()

# Fit the vectorizer to the corpus and transform the corpus into a TF-
IDF matrix
X_tfidf = vectorizer.fit_transform(clean_text)

# Print the TF-IDF matrix as a dense array
print(X_tfidf.toarray(), "\n")

# Print the feature names
print("Feature names:")
print(vectorizer.get_feature_names_out())

[[0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]

```

```
...  
[0. 0. 0. ... 0. 0. 0.]  
[0. 0. 0. ... 0. 0. 0.]  
[0. 0. 0. ... 0. 0. 0.]
```

Feature names:

```
['aa' 'aaa' 'amazon' ... 'zoomed' 'zooming' 'zwave']
```

Word Embedding Techniques (Word2Vec and FastText):

We used advanced word embedding techniques to capture the semantic meaning of words in the reviews.

Word2Vec: This technique uses a neural network model to learn vector representations of words based on their context in the corpus. We trained a Word2Vec model on our tokenized text data to obtain word vectors.

FastText: Similar to Word2Vec, but it also considers subword information, making it better at handling rare and out-of-vocabulary words. We trained a FastText model to generate word vectors that include subword information.

```
from gensim.models import Word2Vec  
from nltk.tokenize import word_tokenize  
  
# Tokenize the text  
sentences = [word_tokenize(doc.lower()) for doc in data['clean_text']]  
  
# Train Word2Vec model  
model = Word2Vec(sentences, vector_size=100, window=5, min_count=1,  
workers=4)  
  
# Get word vectors  
word_vectors = model.wv  
  
# Get the combined matrix of word vectors  
wordvec_matrix = word_vectors.vectors  
print(wordvec_matrix)  
  
[[-5.8863032e-01  1.2726338e+00  3.3599135e-02 ... -7.9840511e-01  
  7.1657944e-01 -1.2033071e-01]  
 [-5.3589469e-01  7.5749975e-01  3.2585797e-01 ... -2.2898255e-01  
 -2.4769144e-01 -3.5513815e-01]  
 [-1.4529744e+00  9.3826878e-01  7.3117822e-02 ... -9.9506333e-02  
  9.8324746e-01 -7.8555740e-02]  
 ...  
 [-6.3266018e-03 -1.8572644e-04  3.8351500e-03 ... -8.7167341e-03  
 -6.1650858e-03  5.2219550e-03]  
 [-9.6880654e-03  2.1857876e-02  2.5954554e-04 ... -2.0557031e-02  
 -2.2558632e-04 -8.7383231e-03]
```

```

[ 5.1067531e-04  5.0789891e-03  7.6249884e-03 ...  8.5238554e-03
 -9.4137760e-03 -8.5921250e-03]]

from gensim.models import FastText
from nltk.tokenize import word_tokenize

# Tokenize the text
sentences = [word_tokenize(doc.lower()) for doc in data['clean_text']]

# Train FastText model
model = FastText(sentences, vector_size=100, window=5, min_count=1,
workers=4)

# Get word vectors
word_vectors = model.wv

# Get the combined matrix of word vectors
fasttext_matrix = word_vectors.vectors
print(fasttext_matrix)

[[-0.573167  -0.97843057 -0.4178658  ...  0.828875  1.5272886
  0.8130299 ]
 [-1.277909   0.07784399 -1.0970418  ... -0.10225663  0.48865247
 -0.28904665]
 [-0.29094222 -0.4713004   0.45327857 ...  0.0015309  0.57042307
  0.28447413]
 ...
 [-0.46185014  0.03203178 -0.33686206 ... -0.17250736  0.34526932
  0.23940298]
 [-0.11214183 -0.10568529 -0.36992288 ... -0.1551006  0.09471703
  0.1914912 ]
 [-0.22303814 -0.16792396 -0.17555399 ... -0.24452431  0.02106808
  0.1862509  ]]
```

- Both Word2Vec and FastText are models used to create word embeddings from text data. Word2Vec focuses on capturing word meanings based on their context in sentences, while FastText adds the ability to understand word structure by considering subword information like prefixes and suffixes.

##Train test split

1. Count vectorizer

```

from sklearn.model_selection import train_test_split

# Separate features and target for each matrix
X = X_count
y = data['labeled']

# Split data into train and test sets
X_train_countvec, X_test_countvec, y_train_countvec, y_test_countvec =
```

```

train_test_split(X, y, test_size=0.2, random_state=42)

# Print the shapes of the training and test sets
print("X_train_countvec shape:", X_train_countvec.shape)
print("y_train_countvec shape:", y_train_countvec.shape)
print("X_test_countvec shape:", X_test_countvec.shape)
print("y_test_countvec shape:", y_test_countvec.shape)

X_train_countvec shape: (27243, 15517)
y_train_countvec shape: (27243,)
X_test_countvec shape: (6811, 15517)
y_test_countvec shape: (6811,)

```

1. TF-IDF VECTORIZER

```

from sklearn.model_selection import train_test_split

X = X_tfidf
y = data['labeled']

# Split data into train and test sets
X_train_tfidf, X_test_tfidf, y_train_tfidf, y_test_tfidf =
train_test_split(X, y, test_size=0.2, random_state=42)

# Print the shapes of the training and test sets
print("X_train_tfidf shape:", X_train_tfidf.shape)
print("y_train_tfidf shape:", y_train_tfidf.shape)
print("X_test_tfidf shape:", X_test_tfidf.shape)
print("y_test_tfidf shape:", y_test_tfidf.shape)

X_train_tfidf shape: (27243, 15517)
y_train_tfidf shape: (27243,)
X_test_tfidf shape: (6811, 15517)
y_test_tfidf shape: (6811,)

import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns # Import seaborn for countplot

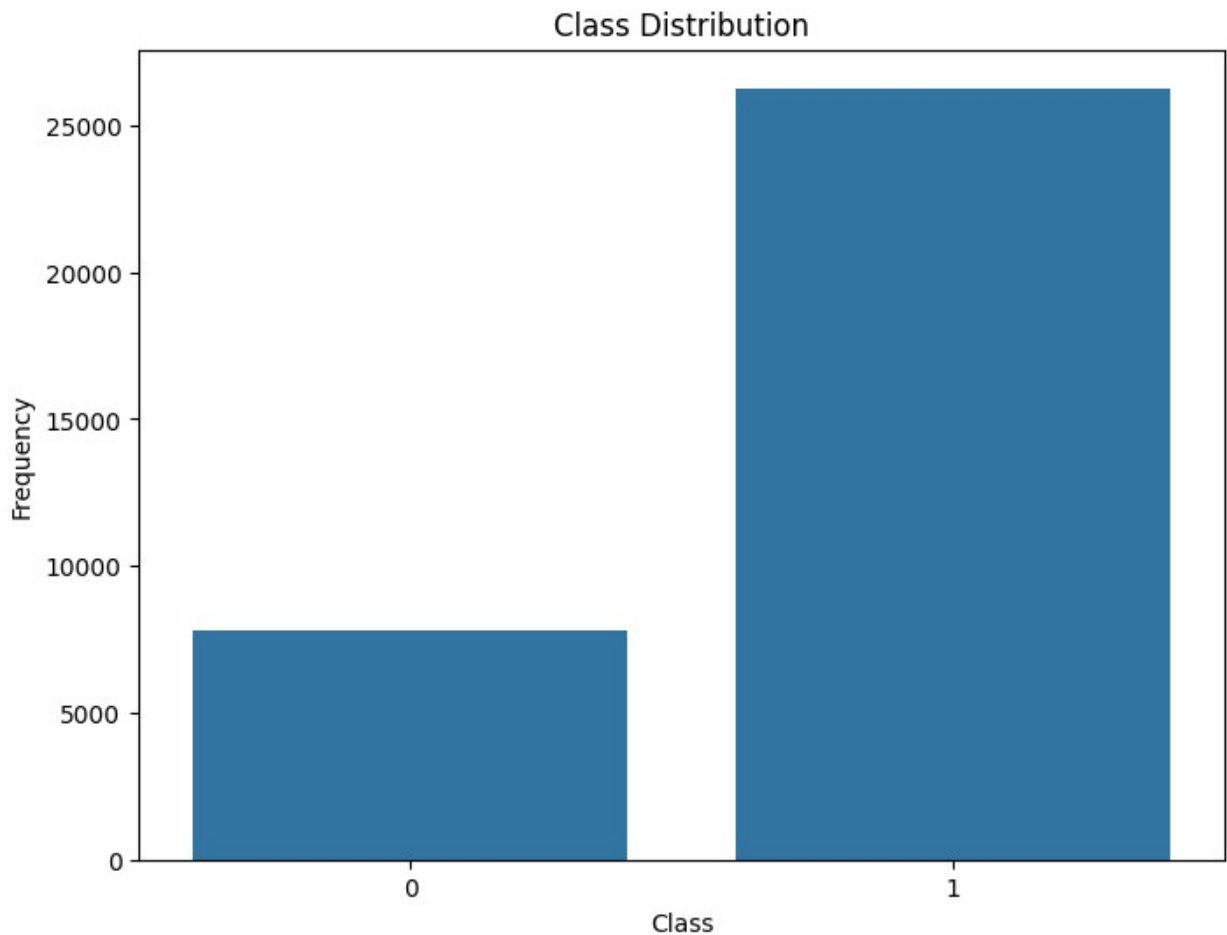
y = data['labeled']
# Create a DataFrame with the target variable
df = pd.DataFrame({'labeled': y})

# Plot the distribution of the target classes using seaborn
plt.figure(figsize=(8, 6))
sns.countplot(data=df, x='labeled')
plt.title('Class Distribution')
plt.xlabel('Class')
plt.ylabel('Frequency')
plt.show()

```



```
# Print the value counts for each class
class_counts = df['labeled'].value_counts()
print(class_counts)
```



```
labeled
1    26271
0     7783
Name: count, dtype: int64

!pip install imbalanced-learn

from imblearn.over_sampling import SMOTE
from collections import Counter

# Assuming X and y are your features and labels
smote = SMOTE(random_state=42)
X_resampled, y_resampled = smote.fit_resample(X, y)

print('Original dataset shape:', Counter(y))
print('Resampled dataset shape:', Counter(y_resampled))
```

```
# Check the new class distribution
df_resampled = pd.DataFrame({'label': y_resampled})
plt.figure(figsize=(8, 6))
sns.countplot(data=df_resampled, x='label')
plt.title('Class Distribution After SMOTE')
plt.xlabel('Class')
plt.ylabel('Frequency')
plt.show()
```

Collecting imbalanced-learn

Downloading imbalanced_learn-0.12.3-py3-none-any.whl (258 kB)
258.3/258.3 kB 2.2 MB/s eta

0:00:00

Requirement already satisfied: numpy>=1.17.3 in
/usr/local/lib/python3.10/dist-packages (from imbalanced-learn)
(1.25.2)

Requirement already satisfied: scipy>=1.5.0 in
/usr/local/lib/python3.10/dist-packages (from imbalanced-learn)
(1.11.4)

Requirement already satisfied: scikit-learn>=1.0.2 in
/usr/local/lib/python3.10/dist-packages (from imbalanced-learn)
(1.2.2)

Requirement already satisfied: joblib>=1.1.1 in
/usr/local/lib/python3.10/dist-packages (from imbalanced-learn)
(1.4.2)

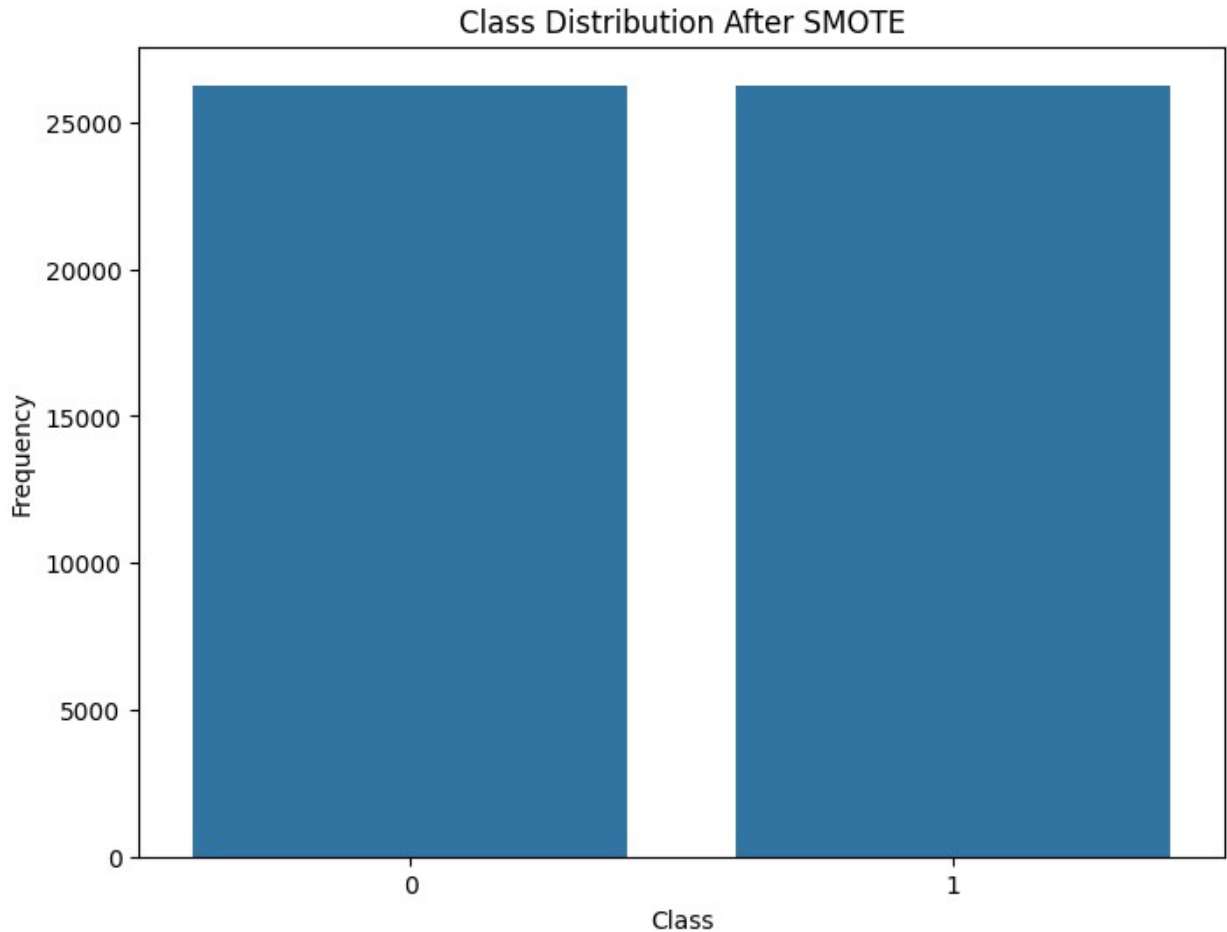
Requirement already satisfied: threadpoolctl>=2.0.0 in
/usr/local/lib/python3.10/dist-packages (from imbalanced-learn)
(3.5.0)

Installing collected packages: imbalanced-learn

Successfully installed imbalanced-learn-0.12.3

Original dataset shape: Counter({1: 26271, 0: 7783})

Resampled dataset shape: Counter({1: 26271, 0: 26271})



MODELLING

BASLINE MODEL

```
from keras.models import Sequential
from keras.layers import Embedding, SimpleRNN, Dense
from keras.callbacks import EarlyStopping

# Define the variables
MAX_NB_WORDS = 1000 # Maximum number of words to consider
EMBEDDING_DIM = 100 # Dimension of the embedding vector
MAX_SEQUENCE_LENGTH = 1000 # Maximum length of the input sequences
epochs = 10
batch_size = 32

#import Libraries
from tensorflow.keras.layers import Embedding, SimpleRNN, Dense
from tensorflow.keras.preprocessing.sequence import pad_sequences #
Import pad_sequences
from tensorflow.keras.models import Sequential
```