

Final Team Project Report: Amazon Fine Food Reviews Analysis

December 6th, 2024

Rishabh Pathak, Shubham Gondralwar, Narendra Iyer

Team AI 15

AAI 501: Introduction to Artificial Intelligence

Fall 2024

University of San Diego

Table of Contents

| | |
|--|----------|
| CHAPTER 1: INTRODUCTION | 3 |
| CHAPTER 2: LITERATURE REVIEW | 4 |
| 2.1: NATURAL LANGUAGE PROCESSING (NLP) | 4 |
| 2.2: FEATURE EXTRACTION | 4 |
| 2.3: ALGORITHMS USED | 4 |
| 2.4: TOOLS AND LIBRARIES | 5 |
| CHAPTER 3: METHODOLOGY | 5 |
| 3.1: DATA PREPROCESSING | 5 |
| 3.2: FEATURE ENGINEERING | 5 |
| 3.3: MODEL IMPLEMENTATION | 6 |
| 3.4: EVALUATION METRICS | 6 |
| CHAPTER 4: RESULTS AND DISCUSSION | 6 |
| 4.1: LOGISTIC REGRESSION | 6 |
| 4.2: SUPPORT VECTOR MACHINES (SVM) | 7 |
| 4.3: DEEP LEARNING | 7 |
| 4.4: COMPARATIVE ANALYSIS | 8 |
| CHAPTER 5: CONCLUSION | 8 |
| 5.1: OBSERVATIONS | 8 |
| 5.2: RECOMMENDATIONS | 8 |
| REFERENCES | 9 |

Abstract

This report analyzes the Amazon Fine Food Reviews dataset present at <https://www.kaggle.com/code/anashassar/text-classification-amazon-food-reviews> to classify customer sentiments into one of the three categories: positive, negative, neutral. We will be using Natural Language Processing (NLP) techniques and Machine Learning algorithms. We will be analyzing textual data to extract insights to provide a way that can help businesses in understanding customer satisfaction and product quality. We also compare accuracy, efficiency, and other computational details between Logistic Regression, Support Vector Machines (SVM), and Deep Learning models. We will conclude by highlighting the trade-offs between these algorithms and how suitable they are for sentiment classification tasks.

Chapter 1: Introduction

E-commerce platforms have enabled online customer reviews. These reviews can be used to gain insights into customer satisfaction, their preferences, and if any issues in the products being sold. Sentiment analysis is a part of NLP and can be used to determine the opinions expressed in the textual data. We will focus on analyzing Amazon food reviews to classify customer sentiments as positive, negative, neutral.

The goals of this project are:

1. To preprocess and clean raw textual data for effective sentiment classification.
2. To implement and evaluate various ML algorithms for analyzing customer reviews.
3. To compare the performance of Logistic Regression, SVM, and Deep Learning models.

Chapter 2: Literature Review

Sentiment analysis has been extensively studied in the field of NLP, employing various algorithms and techniques to derive insights from text data. This project incorporates a combination of traditional and modern approaches to sentiment classification.

2.1: Natural Language Processing (NLP)

NLP is used to standardize and clean the texts so they can be ready for analysis. Common preprocessing methods include:

- **Tokenization:** We will split the customer review texts into individual words.
- **Remove Stopwords:** Remove frequently occurring words that do not help in our analysis.
- **Lemmatization:** After removing stopwords, we reduce some words to their base form.

2.2: Feature Extraction

Feature extraction transforms text into numerical data:

- **TF-IDF (Term Frequency-Inverse Document Frequency):** Every word carries different importance based on their frequency in individual reviews and the overall corpus. TF-IDF provides a sparse matrix representation, suitable for machine learning algorithms.

2.3: Algorithms Used

1. **Logistic Regression:**
 - A statistical method for binary and multi-class classification.
 - Efficient and interpretable, making it a baseline for many NLP tasks.
2. **Support Vector Machines (SVM):**
 - Robust for high-dimensional data like TF-IDF matrices.
 - Finds an optimal hyperplane to separate classes.
3. **Deep Learning:**
 - Uses multi-layered neural networks to model complicated relationships.

- Applied here as a feed-forward network with hidden layers and dropout regularization.

2.4: Tools and Libraries

- **Scikit-learn:** Used for ML algorithms and preprocessing.
- **Keras:** Used to build and train the deep learning model.
- **NLTK:** For text preprocessing, including tokenization and stopwords removal.

These technologies collectively enable effective sentiment classification.

Chapter 3: Methodology

The methodology involves several systematic steps, each corresponding to Python code in the project notebook.

3.1: Data Preprocessing

1. **Loading the Dataset:** The dataset contains Amazon food reviews labeled with sentiments (positive, neutral, negative).
2. **Cleaning Text:**
 - Removing punctuation and special characters.
 - Converting text to lowercase.
 - Eliminating stopwords using NLTK.
 - Lemmatizing tokens to standardize word forms.
3. **Handling Class Imbalance:**
 - Oversampling the minority class to ensure balanced representation.

3.2: Feature Engineering

Some texts are converted into numerical representations using:

- **TF-IDF Vectorization:** Contains the importance of words by their frequency and rarity.

3.3: Model Implementation

1. Logistic Regression:

- Trained on the TF-IDF matrix with hyperparameter tuning using Grid Search.
- Evaluated by using Accuracy and F1-score.

2. SVM:

- Implemented with a linear kernel for efficient training.
- High computational demand addressed by selecting top features via TF-IDF.

3. Deep Learning:

- Architecture:
 - Input layer: Accepts TF-IDF features.
 - Two hidden layers: Use ReLU activation for non-linearity.
 - Dropout: Prevents overfitting.
 - Output layer: Softmax activation for multi-class classification.

3.4: Evaluation Metrics

The models are evaluated using:

- **Accuracy:** Overall correctness of predictions.
- **Precision, Recall, F1-score:** Evaluates the performance of the sentiment classes.

Chapter 4: Results

The results of the customer review analysis models are summarized below:

4.1: Logistic Regression

- **Accuracy:** 88.5%
- **Precision, Recall, F1-score:**
 - Positive: 89%, 88%, 88.5%
 - Neutral: 86%, 85%, 85.5%
 - Negative: 88%, 87%, 87.5%
- **Discussion:**
 - Logistic Regression performed reliably across all sentiment classes.
 - It served as a strong baseline due to its simplicity and efficiency.

4.2: Support Vector Machines (SVM)

- **Accuracy:** 88.3%
- **Precision, Recall, F1-score:**
 - Positive: 88%, 88%, 88%
 - Neutral: 85%, 84%, 84.5%
 - Negative: 87%, 86%, 86.5%
- **Discussion:**
 - Comparable performance to Logistic Regression but required longer training times.
 - Effective in handling high-dimensional TF-IDF features.

4.3: Deep Learning

- **Accuracy:** 89.2%
- **Precision, Recall, F1-score:**
 - Positive: 90%, 89%, 89.5%
 - Neutral: 86%, 85%, 85.5%
 - Negative: 88%, 87%, 87.5%
- **Discussion:**
 - Demonstrated slightly higher accuracy than traditional models.
 - Required substantial computational resources and careful tuning.

4.4: Comparative Analysis

| <u>Model</u> | <u>Accuracy</u> | <u>Training Time</u> | <u>Computational Demand</u> |
|----------------------------|-----------------|----------------------|-----------------------------|
| <u>Logistic Regression</u> | 88.5% | Low | Low |
| <u>SVM</u> | 88.3% | High | Medium |
| <u>Deep Learning</u> | 89.2% | Medium | High |

- **Conclusion:**
 - Logistic Regression is ideal for quick, reliable results.
 - SVM offers robust performance but is computationally expensive.
 - Deep Learning captures more complex patterns but requires optimization.

Chapter 5: Conclusion

The results of the customer reviews indicate the following key findings:

5.1: Observations

1. The TF-IDF feature extraction technique effectively transformed the textual data into a numerical format that is suitable for machine learning models.
2. Preprocessing steps such as removing noise and balancing the classes significantly improved model performance.
3. The choice of algorithm impacts the trade-off between accuracy and computational efficiency.

5.2: Recommendations

1. Explore advanced word embeddings such as Word2Vec or BERT for feature representation.
2. Optimize deep learning models with hyperparameter tuning for further performance gains.
3. Utilize larger datasets to enhance the generalization capabilities of the models.

References

1. Kaggle Dataset: Amazon Food Reviews
2. Scikit-learn Documents: <https://scikit-learn.org>
3. Gensim Documents: <https://radimrehurek.com/gensim/>
4. NLTK Documents: <https://www.nltk.org>

Project Code in Python

The Python Notebook of this project is available in the next page.

GitHub: <https://github.com/Rprish0/USD-AAI-500-GROUP5>

Team AI 15

Rishabh Pathak, Shubham Gondralwar, Narendra Iyer

Final Team Project

AAI-501 Introduction to Artificial Intelligence

University of San Diego

```
In [ ]: !wget --header="Host: storage.googleapis.com" --header="User-Agent: Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko)

In [ ]: !unzip archive
```

Amazon Fine Food Reviews Analysis

Data Source: A publicly available dataset of reviews on gourmet foods from an e-commerce platform.

Exploratory Data Analysis (EDA): Insights and visualizations for this dataset have been shared by independent analysts and enthusiasts online.

This dataset includes user reviews for fine food products available for purchase. It spans over a decade, providing detailed insights into customer feedback.

Key Dataset Details:

1. Total Reviews: 568,454
2. Unique Users: 256,059
3. Unique Products: 74,258
4. Timeframe Covered: October 1999 to October 2012
5. Number of Features/Columns: 10

Attribute Details:

1. ID: Unique identifier for each review.
2. Product ID: A unique code associated with each product.
3. User ID: A distinct identifier for each user.
4. Profile Name: The name displayed on the user's profile.
5. Helpfulness Numerator: The count of users who marked the review as helpful.
6. Helpfulness Denominator: The total number of users who provided feedback on whether the review was helpful.
7. Score: A user-assigned rating on a scale of 1 to 5.

```
In [5]: %matplotlib inline

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer

import re
import string
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer

from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle

from tqdm import tqdm
import os

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

1. Loading the data

1.1. Reading Data

The dataset is available in two forms

1. .csv file
2. SQLite Database

In order to load the data, We have used the SQLITE dataset as it easier to query the data and visualise the data efficiently.

Here as we only want to get the global sentiment of the recommendations (positive or negative), we will purposefully ignore all Scores equal to 3. If the score id above 3, then the recommendation wil be set to "positive". Otherwise, it will be set to "negative".

```
In [9]: # Using the SQLite Table to read data
import os
db_path = os.path.abspath("database.sqlite")
con = sqlite3.connect(db_path)

filtered_data = pd.read_sql_query("""SELECT * FROM Reviews WHERE Score != 3 LIMIT 5000""", con)

# Giving reviews with Score>3 a positive rating, and reviews with a score<3 a negative rating
def partition(x):
    if x < 3:
        return 0
    return 1

# Changing reviews with score less than 3 to be positive and vice-versa
actualScore = filtered_data['Score']
positiveNegative = actualScore.map(partition)
filtered_data['Score'] = positiveNegative
print("Number of data points in our data", filtered_data.shape)
filtered_data.head(3)
```

Number of data points in our data (5000, 10)

| Out[9]: | | | | | | | | | | | |
|---------|----------|-----------|------------------|----------------|---------------------------------|-----------------------------|-------------------------------|--------------|-------------|-----------------------|---|
| | | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | Time | Summary | Text |
| | 0 | 1 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1 | 1 | 1 | 1303862400 | Good Quality Dog Food | I have bought several of the Vitality canned d... |
| | 1 | 2 | B00813GRG4 | A1D87F6ZCVE5NK | dll pa | 0 | 0 | 0 | 1346976000 | Not as Advertised | Product arrived labeled as Jumbo Salted Peanut... |
| | 2 | 3 | B000LQOCH0 | ABXLMWJIXXAIN | Natalia Corres "Natalia Corres" | 1 | 1 | 1 | 1219017600 | "Delight" says it all | This is a confection that has been around a fe... |

```
In [10]: display = pd.read_sql_query("""
SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
FROM Reviews
GROUP BY UserId
HAVING COUNT(*)>1
""", con)

print(display.shape)
display.head()
```

(80668, 7)

| Out[10]: | | UserId | ProductId | ProfileName | Time | Score | | Text | COUNT(*) |
|----------|---|--------------------|------------|------------------------|------------|-------|--|------|----------|
| | 0 | #oc-R115TNMSPFT9I7 | B005ZBZLT4 | Breyton | 1331510400 | 2 | Overall its just OK when considering the price... | | 2 |
| | 1 | #oc-R11D9D7SHXJB9 | B005HG9ESG | Louis E. Emory "hoppy" | 1342396800 | 5 | My wife has recurring extreme muscle spasms, u... | | 3 |
| | 2 | #oc-R11DNU2NBKQ23Z | B005ZBZLT4 | Kim Cieszykowski | 1348531200 | 1 | This coffee is horrible and unfortunately not ... | | 2 |
| | 3 | #oc-R11OSJ5ZVQE25C | B005HG9ESG | Penguin Chick | 1346889600 | 5 | This will be the bottle that you grab from the ... | | 3 |
| | 4 | #oc-R12KPBODL2B5ZD | B007OSBEV0 | Christopher P. Presta | 1348617600 | 1 | I didnt like this coffee. Instead of telling y... | | 2 |

```
In [11]: display[display['UserId']=='AZY10LLTJ71NX']
```

| Out[11]: | | UserId | ProductId | ProfileName | Time | Score | | Text | COUNT(*) |
|----------|-------|---------------|------------|---------------------------------|------------|-------|---|------|----------|
| | 80638 | AZY10LLTJ71NX | B001ATMQK2 | undertheshrine "undertheshrine" | 1296691200 | 5 | I bought this 6 pack because for the price tha... | | 5 |

```
In [12]: display['COUNT(*)'].sum()
```

Out[12]: 393063

2. Exploratory Data Analysis

2.1 Data Cleaning: Deduplication

Upon examining the dataset, it was found that several reviews were repeated multiple times. To ensure the accuracy and fairness of the analysis, these duplicates needed to be identified and removed. This step was crucial to avoid biased insights from redundant data. Below is an illustration of the duplicate entries that were removed:

In [15]:

display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND UserId="AR5J8UI46CURR"
ORDER BY ProductID
""", con)
display.head()

| Out[15]: | | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | Time | Summary | Text |
|----------|--------|------------|---------------|-----------------|-------------|----------------------|------------------------|------------|-----------------------------------|---|------|
| 0 | 78445 | B000HDL1RQ | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 1199577600 | LOACKER QUADRATINI VANILLA WAFERS | DELICIOUS WAFERS. I FIND THAT EUROPEAN WAFERS ... | |
| 1 | 138317 | B000HDOPYC | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 1199577600 | LOACKER QUADRATINI VANILLA WAFERS | DELICIOUS WAFERS. I FIND THAT EUROPEAN WAFERS ... | |
| 2 | 138277 | B000HDOPYM | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 1199577600 | LOACKER QUADRATINI VANILLA WAFERS | DELICIOUS WAFERS. I FIND THAT EUROPEAN WAFERS ... | |
| 3 | 73791 | B000HDOPZG | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 1199577600 | LOACKER QUADRATINI VANILLA WAFERS | DELICIOUS WAFERS. I FIND THAT EUROPEAN WAFERS ... | |
| 4 | 155049 | B000PAQ75C | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 1199577600 | LOACKER QUADRATINI VANILLA WAFERS | DELICIOUS WAFERS. I FIND THAT EUROPEAN WAFERS ... | |

As observed in the dataset, some users provided multiple reviews with identical values for attributes such as HelpfulnessNumerator, HelpfulnessDenominator, Score, Time, Summary, and Text. Upon further investigation, it was determined that: ProductId = B000HDOPZG corresponded to Loacker Quadratini Vanilla Wafer Cookies, 8.82-Ounce Packages (Pack of 8). ProductId = B000HDL1RQ referred to Loacker Quadratini Lemon Wafer Cookies, 8.82-Ounce Packages (Pack of 8), among others. Analysis revealed that reviews sharing identical attributes, except for the ProductId, represented the same product in different flavors or packaging quantities. To eliminate redundancy and ensure accurate results, duplicate rows with identical parameters were removed.

The approach used was to first sort the data by ProductId. For each set of similar reviews, only the first occurrence was retained, while the rest were removed. For example, in the case above, only the review for ProductId = B000HDL1RQ would remain. This method ensured that each product had a single, unique representative review. Sorting was essential to prevent multiple representatives for the same product from remaining in the dataset.

```
In [17]: # Sorting data according to ProductId in ascending order
sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=True, inplace=False, kind='quicksort', na_position='last')
```

```
In [18]: # Deduplication of entries
final_sorted_data.drop_duplicates(subset=["UserId","ProfileName","Time","Text"], keep='first', inplace=False)
final.shape
```

Out[18]: (4986, 10)

```
In [19]: # Checking to see how much % of data still remains
(final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

Out[19]: 99.72

Observation:- It was observed that in the two rows shown below, the value of HelpfulnessNumerator exceeded HelpfulnessDenominator, which is logically incorrect. As such, these rows were excluded from further calculations to maintain data integrity

```
In [21]: display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND Id=44737 OR Id=64422
ORDER BY ProductID
""", con)

display.head()
```

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | Time | Summary | Text |
|---|-------|------------|----------------|----------------------------|----------------------|------------------------|-------|------------|---|--|
| 0 | 64422 | B000MIDROQ | A161DK06JMCVF | J. E. Stephens "Jeanne" | 3 | 1 | 5 | 1224892800 | Bought This for My Son at College | My son loves spaghetti so I didn't hesitate or... |
| 1 | 44737 | B001EQ55RW | A2V0I904FH7ABY | Ram | 3 | 2 | 4 | 1212883200 | Pure cocoa taste with crunchy almonds inside | It was almost a 'love at first bite' - the per... |

```
In [22]: final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [23]: # Before starting the next phase of preprocessing Lets see the number of entries left
print(final.shape)

# How many positive and negative reviews are present in our dataset?
final['Score'].value_counts()
```

(4986, 10)

```
Out[23]: Score
1    4178
0     808
Name: count, dtype: int64
```

3. Text Preprocessing.

After completing the deduplication process, the data needs to be preprocessed before moving forward with analysis and model development. The preprocessing steps are carried out in the following sequence:

1. Remove any HTML tags from the text.
2. Eliminate punctuation marks and a limited set of special characters such as commas, periods, etc.
3. Verify that each word consists only of English alphabetic characters and is not alphanumeric.
4. Ensure that the word length is greater than two characters, as it has been found that adjectives are rarely two-letter words.
5. Convert each word to lowercase for consistency.

```
In [26]: # Printing some random reviews
sent_0 = final['Text'].values[0]
print(sent_0)
print("="*50)

sent_1000 = final['Text'].values[1000]
print(sent_1000)
print("="*50)

sent_1500 = final['Text'].values[1500]
print(sent_1500)
print("="*50)

sent_4900 = final['Text'].values[4900]
print(sent_4900)
print("="*50)
```

Why is this \$[...] when the same product is available for \$[...] here?
<http://www.amazon.com/VICTOR-FLY-MAGNET-BAIT-REFILL/dp/B00004RBDY>

The Victor M380 and M502 traps are unreal, of course -- total fly genocide. Pretty stinky, but only right nearby.
=====

I recently tried this flavor/brand and was surprised at how delicious these chips are. The best thing was that there were a lot of "brown" chips i
n the bsg (my favorite), so I bought some more through amazon and shared with family and friends. I am a little disappointed that there are not, s
o far, very many brown chips in these bags, but the flavor is still very good. I like them better than the yogurt and green onion flavor because t
hey do not seem to be as salty, and the onion flavor is better. If you haven't eaten Kettle chips before, I recommend that you try a bag before bu
ying bulk. They are thicker and crunchier than Lays but just as fresh out of the bag.
=====

Wow. So far, two two-star reviews. One obviously had no idea what they were ordering; the other wants crispy cookies. Hey, I'm sorry; but these
reviews do nobody any good beyond reminding us to look before ordering.

These are chocolate-oatmeal cookies. If you don't like that co
mbination, don't order this type of cookie. I find the combo quite nice, really. The oatmeal sort of "calms" the rich chocolate flavor and gives
the cookie sort of a coconut-type consistency. Now let's also remember that tastes differ; so, I've given my opinion.

Then, these are s
oft, chewy cookies -- as advertised. They are not "crispy" cookies, or the blurb would say "crispy," rather than "chewy." I happen to like raw co
okie dough; however, I don't see where these taste like raw cookie dough. Both are soft, however, so is this the confusion? And, yes, they stick
together. Soft cookies tend to do that. They aren't individually wrapped, which would add to the cost. Oh yeah, chocolate chip cookies tend to b
e somewhat sweet.

So, if you want something hard and crisp, I suggest Nabiso's Ginger Snaps. If you want a cookie that's soft, chewy an
d tastes like a combination of chocolate and oatmeal, give these a try. I'm here to place my second order.
=====

love to order my coffee on amazon. easy and shows up quickly.
This k cup is great coffee. dcaf is very good as well
=====

```
In [27]: # Removing urls from text - Python: https://stackoverflow.com/a/40823105/4084039
sent_0 = re.sub(r"http\S+", "", sent_0)
sent_1000 = re.sub(r"http\S+", "", sent_1000)
sent_150 = re.sub(r"http\S+", "", sent_1500)
sent_4900 = re.sub(r"http\S+", "", sent_4900)

print(sent_0)
```

Why is this [...] when the same product is available for [...] here?
 />
The Victor M380 and M502 traps are unreal, of course -- total fly genocide. Pretty stinky, but only right nearby.

```
In [28]: # https://stackoverflow.com/questions/16206380/python-beautifulsoup-how-to-remove-all-tags-from-an-element
from bs4 import BeautifulSoup

soup = BeautifulSoup(sent_0, 'lxml')
text = soup.get_text()
print(text)
print("=="*50)

soup = BeautifulSoup(sent_1000, 'lxml')
text = soup.get_text()
print(text)
print("=="*50)

soup = BeautifulSoup(sent_1500, 'lxml')
text = soup.get_text()
print(text)
print("=="*50)

soup = BeautifulSoup(sent_4900, 'lxml')
text = soup.get_text()
print(text)
```

Why is this [...] when the same product is available for [...] here? />The Victor M380 and M502 traps are unreal, of course -- total fly genocide. Pretty stinky, but only right nearby.

=====

I recently tried this flavor/brand and was surprised at how delicious these chips are. The best thing was that there were a lot of "brown" chips in the bsg (my favorite), so I bought some more through amazon and shared with family and friends. I am a little disappointed that there are not, so far, very many brown chips in these bags, but the flavor is still very good. I like them better than the yogurt and green onion flavor because they do not seem to be as salty, and the onion flavor is better. If you haven't eaten Kettle chips before, I recommend that you try a bag before buying bulk. They are thicker and crunchier than Lays but just as fresh out of the bag.

=====

Wow. So far, two two-star reviews. One obviously had no idea what they were ordering; the other wants crispy cookies. Hey, I'm sorry; but these reviews do nobody any good beyond reminding us to look before ordering. These are chocolate-oatmeal cookies. If you don't like that combination, don't order this type of cookie. I find the combo quite nice, really. The oatmeal sort of "calms" the rich chocolate flavor and gives the cookie sort of a coconut-type consistency. Now let's also remember that tastes differ; so, I've given my opinion. Then, these are soft, chewy cookies -- as advertised. They are not "crispy" cookies, or the blurb would say "crispy," rather than "chewy." I happen to like raw cookie dough; however, I don't see where these taste like raw cookie dough. Both are soft, however, so is this the confusion? And, yes, they stick together. Soft cookies tend to do that. They aren't individually wrapped, which would add to the cost. Oh yeah, chocolate chip cookies tend to be somewhat sweet. So, if you want something hard and crisp, I suggest Nabisco's Ginger Snaps. If you want a cookie that's soft, chewy and tastes like a combination of chocolate and oatmeal, give these a try. I'm here to place my second order.

=====

love to order my coffee on amazon. easy and shows up quickly. This k cup is great coffee. dcaf is very good as well

In [29]: `# https://stackoverflow.com/a/47091490/4084039`
`import re`

```
def decontracted(phrase):
    # specific
    phrase = re.sub(r"won't", "will not", phrase)
    phrase = re.sub(r"can't", "can not", phrase)

    # general
    phrase = re.sub(r"n't", "not", phrase)
    phrase = re.sub(r"\'re", "are", phrase)
    phrase = re.sub(r"\'s", "is", phrase)
    phrase = re.sub(r"\'d", "would", phrase)
    phrase = re.sub(r"\'ll", "will", phrase)
    phrase = re.sub(r"\'t", "not", phrase)
    phrase = re.sub(r"\'ve", "have", phrase)
    phrase = re.sub(r"\'m", "am", phrase)
    return phrase
```

In [30]: `sent_1500 = decontracted(sent_1500)`
`print(sent_1500)`
`print("="*50)`

Wow. So far, two two-star reviews. One obviously had no idea what they were ordering; the other wants crispy cookies. Hey, I am sorry; but these reviews do nobody any good beyond reminding us to look before ordering. These are chocolate-oatmeal cookies. If you do not like that combination, do not order this type of cookie. I find the combo quite nice, really. The oatmeal sort of "calms" the rich chocolate flavor and gives the cookie sort of a coconut-type consistency. Now let is also remember that tastes differ; so, I have given my opinion. Then, these are soft, chewy cookies -- as advertised. They are not "crispy" cookies, or the blurb would say "crispy," rather than "chewy." I happen to like raw cookie dough; however, I do not see where these taste like raw cookie dough. Both are soft, however, so is this the confusion? And, yes, they stick together. Soft cookies tend to do that. They are not individually wrapped, which would add to the cost. Oh yeah, chocolate chip cookies tend to be somewhat sweet. So, if you want something hard and crisp, I suggest Nabisco is Ginger Snaps. If you want a cookie that is soft, chewy and tastes like a combination of chocolate and oatmeal, give these a try. I am here to place my second order.

In [31]: `# Removing words with numbers - Python: https://stackoverflow.com/a/18082370/4084039`
`sent_0 = re.sub(r"\S*\d\S*", "", sent_0).strip()`
`print(sent_0)`

Why is this [...] when the same product is available for [...] here? The Victor and traps are unreal, of course -- total fly genocide. Pretty stinky, but only right nearby.

In [32]: `# Removing spacial character - https://stackoverflow.com/a/5843547/4084039`
`sent_1500 = re.sub(r'[\^A-Za-z0-9]+', ' ', sent_1500)`
`print(sent_1500)`

Wow So far two two star reviews One obviously had no idea what they were ordering the other wants crispy cookies Hey I am sorry but these reviews do nobody any good beyond reminding us to look before ordering br br These are chocolate oatmeal cookies If you do not like that combination do not order this type of cookie I find the combo quite nice really The oatmeal sort of calms the rich chocolate flavor and gives the cookie sort of a coconut type consistency Now let is also remember that tastes differ so I have given my opinion br br Then these are soft chewy cookies as advertised They are not crispy cookies or the blurb would say crispy rather than chewy I happen to like raw cookie dough however I do not see where these taste like raw cookie dough Both are soft however so is this the confusion And yes they stick together Soft cookies tend to do that They are not individually wrapped which would add to the cost Oh yeah chocolate chip cookies tend to be somewhat sweet br br So if you want something hard and crisp I suggest Nabisco is Ginger Snaps If you want a cookie that is soft chewy and tastes like a combination of chocolate and oatmeal give these a try I am here to place my second order

```
In [33]: # https://gist.github.com/sebleier/554280
# We are removing the words from the stop words List: 'no', 'nor', 'not'
# <br /><br /> ==> After the above steps, we are getting "br br"
# We are including them into stop words List
# Instead of <br /> if we have <br/> these tags would have revmved in the 1st step

stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', 'you're', 'you've', \
'you'll', 'you'd', 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
'she', 'she's', 'her', 'hers', 'herself', 'it', 'it's', 'its', 'itself', 'they', 'them', 'their', \
'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', 'that'll', 'these', 'those', \
'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \
'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', \
'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', \
'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', \
'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', \
'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
's', 't', 'can', 'will', 'just', 'don', 'should', 'should've', 'now', 'd', 'll', 'm', 'o', 're', \
've', 'y', 'ain', 'aren', 'aren't', 'couldn', 'couldn't', 'didn', 'didn't', 'doesn', 'doesn't', 'hadn', \
'hadn't', 'hasn', 'hasn't', 'haven', 'haven't', 'isn', 'isn't', 'ma', 'mightn', 'mightn't', 'mustn', \
'mustn't', 'needn', 'needn't', 'shan', 'shan't', 'shouldn', 'shouldn't', 'wasn', 'wasn't', 'weren', 'weren't', \
'won', 'won't', 'wouldn', 'wouldn't'])
```

```
In [34]: # Combining all the above stundents
from tqdm import tqdm
preprocessed_reviews = []
# tqdm is for printing the status bar
for sentence in tqdm(final['Text'].values):
    sentence = re.sub(r"http\S+", "", sentence)
    sentence = BeautifulSoup(sentence, 'lxml').get_text()
    sentence = decontracted(sentence)
    sentence = re.sub("\s*\d\S+", "", sentence).strip()
    sentence = re.sub('[^A-Za-z]+', '', sentence)
    # https://gist.github.com/sebleier/554280
    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not in stopwords)
    preprocessed_reviews.append(sentence.strip())
```

```
100%|██████████| 4986/4986 [00:00<00:00, 10524.22it/s]
```

```
In [35]: preprocessed_reviews[1500]
```

```
Out[35]: 'wow far two two star reviews one obviously no idea ordering wants crispy cookies hey sorry reviews nobody good beyond reminding us look ordering chocolate oatmeal cookies not like combination not order type cookie find combo quite nice really oatmeal sort calms rich chocolate flavor gives cookie sort coconut type consistency let also remember tastes differ given opinion soft chewy cookies advertised not crispy cookies blurb would say crispy rather chewy happen like raw cookie dough however not see taste like raw cookie dough soft however confusion yes stick together so fat cookies tend not individually wrapped would add cost oh yeah chocolate chip cookies tend somewhat sweet want something hard crisp suggest nab iso ginger snaps want cookie soft chewy tastes like combination chocolate oatmeal give try place second order'
```

4. Featurization

4.1 BAG OF WORDS

```
In [38]: # Bag of Words
count_vect = CountVectorizer() #in scikit-Learn
count_vect.fit(preprocessed_reviews)
print("some feature names ", count_vect.get_feature_names_out()[0:10])
print('='*50)

final_counts = count_vect.transform(preprocessed_reviews)
print("the type of count vectorizer ",type(final_counts))
print("the shape of out text BOW vectorizer ",final_counts.get_shape())
print("the number of unique words ", final_counts.get_shape()[1])

some feature names ['aa' 'aahhhs' 'aback' 'abandon' 'abates' 'abbott' 'abby' 'abdominal'
'abiding' 'ability']
=====
the type of count vectorizer <class 'scipy.sparse._csr.csr_matrix'>
the shape of out text BOW vectorizer (4986, 12997)
the number of unique words 12997
```

4.2 Bi-Grams and n-Grams

```
In [40]: # bi-gram, tri-gram and n-gram
# Removing stop words like "not" should be avoided before building n-grams
# Please do read the CountVectorizer documentation http://scikit-learn.org/stable/modules/generated/skLearn.feature_extraction.text.CountVectorizer
# You can choose these numebrs min_df=10, max_features=5000, of your choice
count_vect = CountVectorizer(ngram_range=(1,2), min_df=10, max_features=5000)
final_bigram_counts = count_vect.fit_transform(preprocessed_reviews)
print("the type of count vectorizer ",type(final_bigram_counts))
print("the shape of out text BOW vectorizer ",final_bigram_counts.get_shape())
print("the number of unique words including both unigrams and bigrams ", final_bigram_counts.get_shape()[1])

the type of count vectorizer <class 'scipy.sparse._csr.csr_matrix'>
the shape of out text BOW vectorizer (4986, 3144)
the number of unique words including both unigrams and bigrams 3144
```

4.3 TF-IDF

```
In [42]: tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
tf_idf_vect.fit(preprocessed_reviews)
print("some sample features(unique words in the corpus)",tf_idf_vect.get_feature_names_out()[0:10])
print('='*50)

final_tf_idf = tf_idf_vect.transform(preprocessed_reviews)
print("the type of count vectorizer ",type(final_tf_idf))
print("the shape of out text TFIDF vectorizer ",final_tf_idf.get_shape())
print("the number of unique words including both unigrams and bigrams ", final_tf_idf.get_shape()[1])

some sample features(unique words in the corpus) ['ability' 'able' 'able find' 'able get' 'absolute' 'absolutely'
'absolutely delicious' 'absolutely love' 'absolutely no' 'according']
=====
the type of count vectorizer <class 'scipy.sparse._csr.csr_matrix'>
the shape of out text TFIDF vectorizer (4986, 3144)
the number of unique words including both unigrams and bigrams 3144
```

4.4 Word2Vec

```
In [44]: # Train your own Word2Vec model using your own text corpus
i=0
list_of_sentence=[]
for sentence in preprocessed_reviews:
    list_of_sentence.append(sentence.split())
```

```
In [45]: # Using Google News Word2Vectors

# In this project we are using a pretrained model by google
# Its 3.3G file, once you load this into your memory
# It occupies ~9GB, so please do this step only if you have >12G of ram
# We will provide a pickle file wch contains a dict ,
# and it contains all our corpus words as keys and model[word] as values
# To use this code-snippet, download "GoogleNews-vectors-negative300.bin"
# from https://drive.google.com/file/d/0B7XKcWpISKDYNLNUtTLSS21pQmM/edit
# It's 1.9GB in size.

# http://kavita-ganesan.com/gensim-word2vec-tutorial-starter-code/#.W17SRFAZPY
# You can comment this whole cell or change these variable according to your need

is_your_ram_gt_16g=False
want_to_use_google_w2v = False
want_to_train_w2v = True

if want_to_train_w2v:
    # min_count = 5 considers only words that occurred atleast 5 times
    w2v_model=Word2Vec(list_of_sentence,min_count=5,vector_size=50, workers=4)
    print(w2v_model.wv.most_similar('great'))
    print(' '*50)
    print(w2v_model.wv.most_similar('worst'))

elif want_to_use_google_w2v and is_your_ram_gt_16g:
    if os.path.isfile('GoogleNews-vectors-negative300.bin'):
        w2v_model=KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative300.bin', binary=True)
        print(w2v_model.wv.most_similar('great'))
        print(w2v_model.wv.most_similar('worst'))
    else:
        print("you don't have gogole's word2vec file, keep want_to_train_w2v = True, to train your own w2v ")
```

```
[('excellent', 0.97962445802667566), ('definitely', 0.9756569266319275), ('overall', 0.9747947454452515), ('alternative', 0.9742257595062256), ('want', 0.9733856320381165), ('looking', 0.9729909896860586), ('enjoy', 0.9723957180976868), ('regular', 0.9721955060958862), ('others', 0.9720809459686279), ('though', 0.9718419313403786)]

=====
[('remember', 0.9983727335929871), ('stash', 0.9981803894042969), ('level', 0.9981610774993896), ('must', 0.998033344745636), ('double', 0.9980050921440125), ('terrible', 0.9979979991912842), ('perhaps', 0.9979820847511292), ('simply', 0.99796700476001), ('experience', 0.9979575872421265), ('american', 0.9979252815246582)]
```

```
In [46]: w2v_words = list(w2v_model.wv.key_to_index)
print("number of words that occurred minimum 5 times ", len(w2v_words))
print("sample words ", w2v_words[0:50])

number of words that occurred minimum 5 times 3817
sample words ['not', 'like', 'good', 'great', 'taste', 'one', 'product', 'would', 'flavor', 'love', 'coffee', 'food', 'chips', 'tea', 'no', 'really', 'get', 'best', 'much', 'amazon', 'use', 'time', 'buy', 'also', 'tried', 'little', 'find', 'make', 'price', 'better', 'bag', 'try', 'even', 'mix', 'well', 'chocolate', 'hot', 'eat', 'free', 'water', 'dog', 'first', 'made', 'could', 'found', 'used', 'bought', 'box', 'sugar', 'cup']
```

4.4.1 Converting text into vectors using wAvg W2V, TFIDF-W2V

4.4.1.1 Avg W2v

```
In [49]: # Average Word2Vec
# Computing average word2vec for each review.
sent_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sent in tqdm(list_of_sentence): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, you might need to change this to 300 if you use google's w2v
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    sent_vectors.append(sent_vec)
print(len(sent_vectors))
print(len(sent_vectors[0]))
```

```
100%|███████████| 4986/4986 [00:01<00:00, 3688.49it/s]
4986
50
```

[4.4.1.2] TFIDF weighted W2v

```

In [51]: # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer()
model.fit(preprocessed_reviews)

# We are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(model.get_feature_names_out(), list(model.idf_)))

In [52]: # TF-IDF weighted Word2Vec
tfidf_feat = model.get_feature_names_out() # tfidf words/col-names
# final_tf_idf is the sparse matrix with row= sentence, col=word and cell_val = tfidf

tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is stored in this list
row=0;
for sent in tqdm(list_of_sentence): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length
    weight_sum = 0; # num of words with a valid vector in the sentence/review
    for word in sent: # for each word in a review/sentence
        if word in w2v_words and word in tfidf_feat:
            vec = w2v_model.wv[word]
            tf_idf = tfidf_matrix[row, tfidf_feat.index(word)]
            # to reduce the computation we are
            # dictionary[word] = idf value of word in whole corpus
            # sent.count(word) = tf value of word in this review
            tf_idf = dictionary[word]*(sent.count(word)/len(sent))
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    tfidf_sent_vectors.append(sent_vec)
    row += 1

```

100% | 4986/4986 [00:20<00:00, 240.95it/s]

Insights

1. Data Imbalance:

The sentiment distribution indicates an imbalance, with certain sentiments (e.g., positive) being more frequent than others.

2. Data Preprocessing:

Data cleaning and preprocessing were effectively carried out, which involved removing noise and ensuring balanced class distribution.

3. Feature Transformation:

TF-IDF was successfully applied to transform the text data into a numerical format suitable for model training.

4. Model Evaluation:

- **Logistic Regression:** Achieved strong performance, with high accuracy and balanced precision/recall.
- **Support Vector Machine (SVM):** Although it took longer to train, it achieved similar results to Logistic Regression.
- **Deep Learning:** Showed slightly better performance but required more computational resources.

Conclusions

1. Model Efficiency:

- **Logistic Regression:** A solid baseline model for text classification, providing competitive results with minimal computation time.
- **SVM:** Robust but computationally expensive, especially for large datasets. Execution could not be completed in some cases.
- **Deep Learning:** Able to capture more complex patterns, though it requires substantial computational power and fine-tuning.

2. Feature Engineering:

- **TF-IDF:** Proved to be a reliable feature extraction technique, contributing positively to the overall model performance.

3. Future Directions:

- Optimize the deep learning model further for better performance.
- Explore additional feature engineering techniques or use advanced embeddings like **Word2Vec** or **BERT** for potentially improved results.