

Amazon Fine Food Reviews Sentiment Analysis

Author : Mariam Ghareeb

Date: July 2025

Contact: mariamghareeb376@gmail.com

Project Overview

This report presents a comprehensive sentiment analysis of Amazon's fine food reviews dataset, demonstrating advanced natural language processing techniques and machine learning implementation. The project showcases an end-to-end data science workflow from raw text processing to actionable business insights.

Objective

To extract meaningful sentiment patterns from customer reviews using state-of-the-art NLP models and provide data-driven recommendations for business strategy optimization.

Business Context

Customer reviews represent a goldmine of unstructured feedback that can drive product development, marketing strategies, and customer experience improvements. This analysis transforms thousands of text reviews into quantifiable insights that can directly impact business decisions.

Key Technologies Used

- **Python:** Core programming language for data analysis and model implementation
- **NLTK:** Natural Language Toolkit for text preprocessing and linguistic analysis
- **Hugging Face Transformers:** State-of-the-art RoBERTa model for sentiment classification
- **Pandas & NumPy:** Data manipulation and numerical computing
- **Matplotlib & Seaborn:** Advanced data visualization and statistical plotting
- **Scikit-learn:** Machine learning utilities and model evaluation metrics



Data Exploration & Preprocessing

```
[5]: df.head()
```

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1	5	1303862400	Good Quality Dog Food	I have bought several of the Vitality canned d...
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0	1	1346976000	Not as Advertised	Product arrived labeled as Jumbo Salted Peanut...
2	3	B000LQOCH0	ABXLMWJ0XXAIN	Natalia Corres "Natalia Corres"	1	1	4	1219017600	"Delight" says it all	This is a confection that has been around a fe...
3	4	B000UA0QIQ	A395B0RC6FGVXV	Karl	3	3	2	1307923200	Cough Medicine	If you are looking for the secret ingredient I...
4	5	B006K2ZZ7K	A1UQRSCLF8GW1T	Michael D. Bigham "M. Wassir"	0	0	5	1350777600	Great taffy	Great taffy at a great price. There was a wid...

We are performing this analysis on the Text column

```
[6]: df.describe()
```

	Id	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time
count	568454.000000	568454.000000	568454.000000	568454.000000	5.684540e+05
mean	284227.500000	1.743817	2.22881	4.183199	1.296257e+09
std	164098.679298	7.636513	8.28974	1.310436	4.804331e+07
min	1.000000	0.000000	0.00000	1.000000	9.393408e+08
25%	142114.250000	0.000000	0.00000	4.000000	1.271290e+09
50%	284227.500000	0.000000	1.00000	5.000000	1.311120e+09
75%	426340.750000	2.000000	2.00000	5.000000	1.332720e+09
max	568454.000000	866.000000	923.00000	5.000000	1.351210e+09

Key Dataset Characteristics:

- Quite large dataset.
- No missing values in text which w are doing our analysis on.
- Missing values are in profile name and summary.
- There is a bias in positive review in this dataset.

```
[5]: df.shape  
#quite large dataset
```

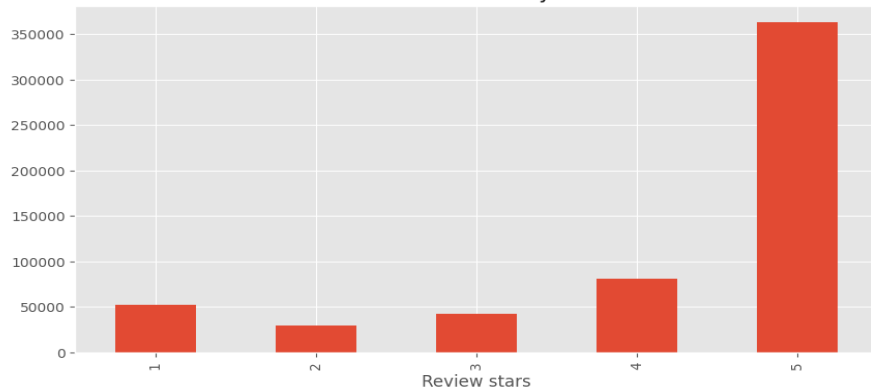
```
[5]: (568454, 10)
```

[+ Code](#)[+ Markdown](#)

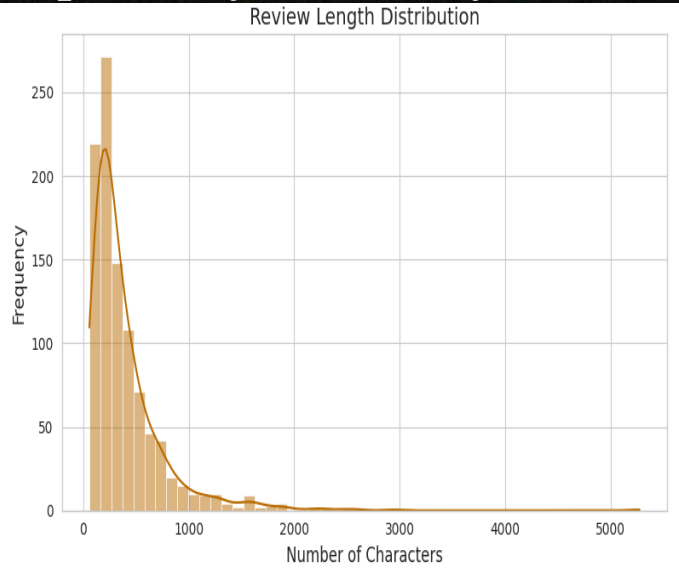
```
[6]: df.info()  
#no null values in the text
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 568454 entries, 0 to 568453  
Data columns (total 10 columns):  
#   Column              Non-Null Count  Dtype  
---  -  
0   Id                  568454 non-null  int64  
1   ProductId           568454 non-null  object  
2   UserId              568454 non-null  object  
3   ProfileName         568428 non-null  object  
4   HelpfulnessNumerator 568454 non-null  int64  
5   HelpfulnessDenominator 568454 non-null  int64  
6   Score               568454 non-null  int64  
7   Time                568454 non-null  int64  
8   Summary             568427 non-null  object  
9   Text                568454 non-null  object  
dtypes: int64(5), object(5)  
memory usage: 43.4+ MB
```

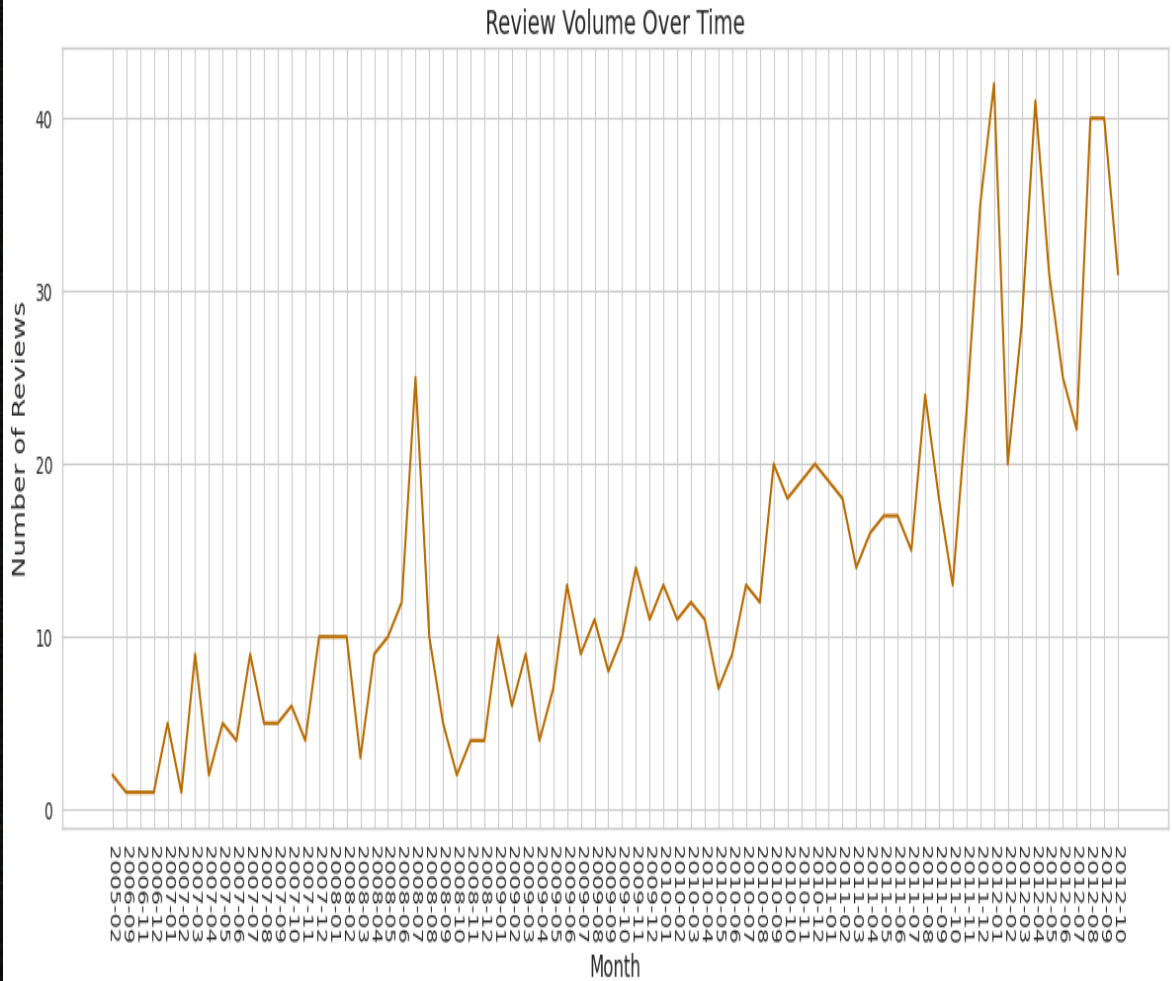
Count of review by size



Exploratory Data Analysis



- Review Length Analysis
- Temporal Patterns



Data Preprocessing Pipeline

- **Tokenization:** Breaking text into individual words using NLTK.
- **Stopword Removal:** Eliminating common words that don't contribute to Sentiment.

Lemmatization: Converting words to their root forms for better analysis.

VADER Sentiment Scoring

+ Code

+ Markdown

NLTK will get the neg/neu/pos scored of text

- This uses a "bag of words approach"
 - 1.Stop words removed (and, the,....etc)
 - 2.each word is scored and combined to total score

[13]:
vaders.head()

		Id	neg	neu	pos	compound	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text
0	1	0.000	0.695	0.305	0.9441	B001E4KFG0	A3SGXH7AUHU8GW		delmartian	1	1	5	1303862400	Good Quality Dog Food	I have bought several of the Vitality canned d...
1	2	0.138	0.862	0.000	-0.5664	B00813GRG4	A1D87F6ZCVE5NK		dll pa	0	0	1	1346976000	Not as Advertised	Product arrived labeled as Jumbo Salted Peanut...
2	3	0.091	0.754	0.155	0.8265	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres	"Natalia Corres"	1	1	4	1219017600	"Delight" says it all	This is a confection that has been around a fe...
3	4	0.000	1.000	0.000	0.0000	B000UAOQIQ	A395BORC6FGVXV		Karl	3	3	2	1307923200	Cough Medicine	If you are looking for the secret ingredient i...
4	5	0.000	0.552	0.448	0.9468	B006K2ZZ7K	A1UQR5CLF8GW1T	Michael D. Bigham	"M. Wassir"	0	0	5	1350777600	Great taffy	Great taffy at a great price. There was a wid...

[10]:
from nltk.sentiment import SentimentIntensityAnalyzer
from tqdm.notebook import tqdm
sia=SentimentIntensityAnalyzer()

Run the polarity score on the whole dataset!

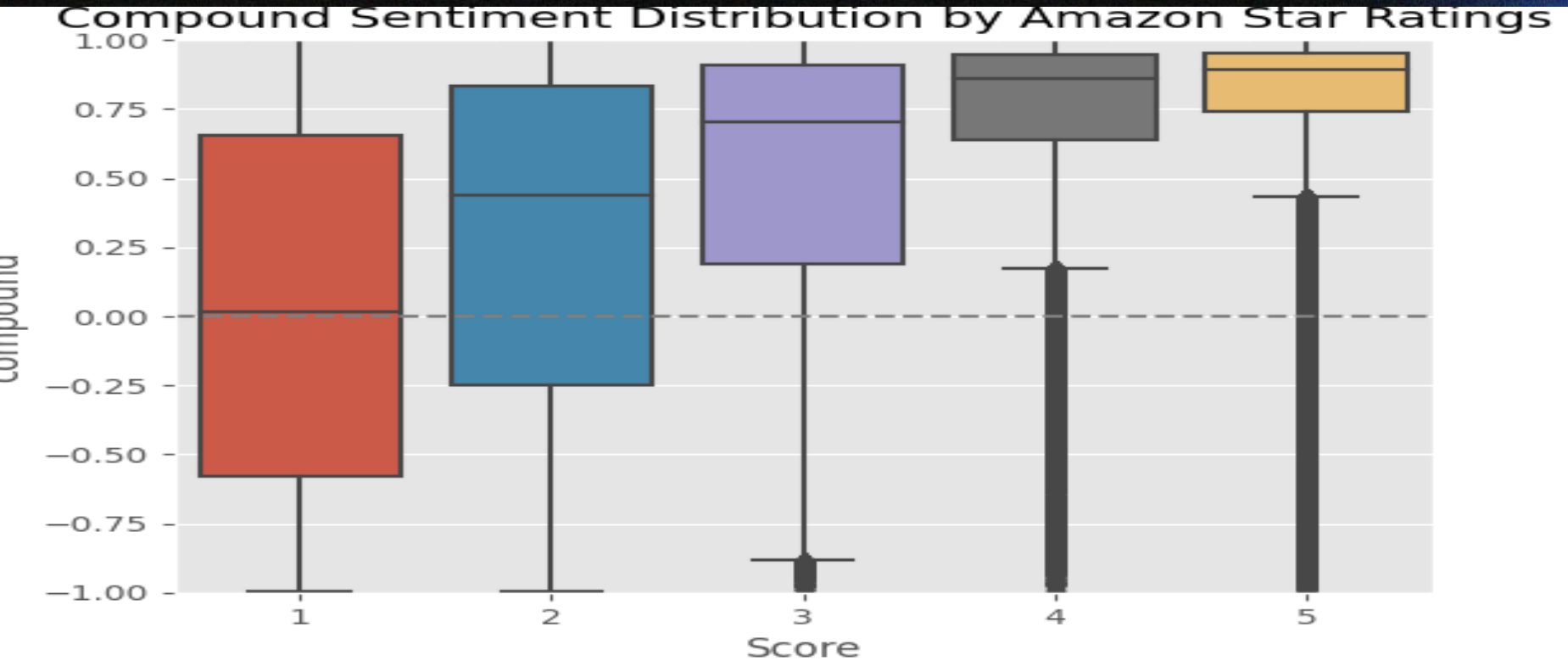
[11]:
results={}
for i ,row in tqdm(df.iterrows(),total=len(df)):
 text=row["Text"]
 myid=row["Id"]
 results[myid]=sia.polarity_scores(text)

100% 568454/568454 [08:21<00:00, 1243.07it/s]

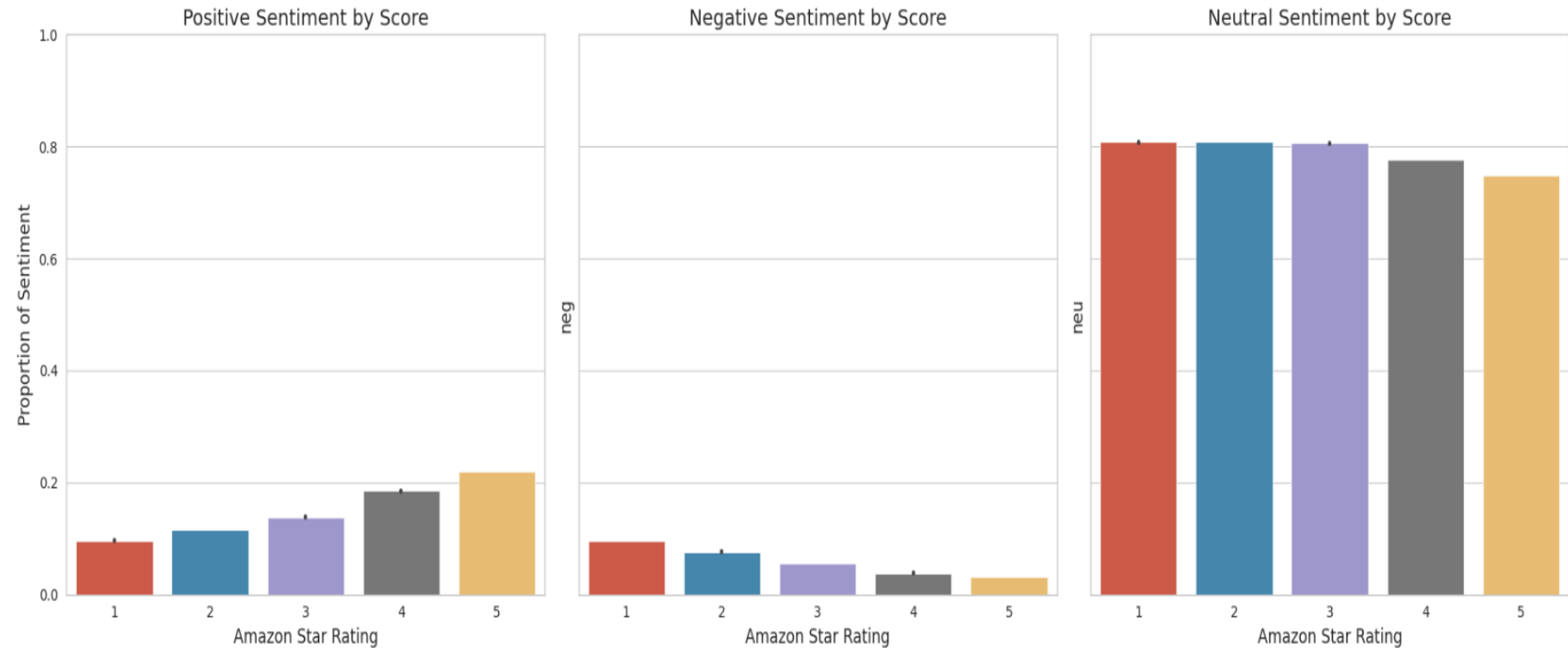
[12]:
vaders=pd.DataFrame(results).T
vaders=vaders.reset_index().rename(columns={'index':'Id'})
vaders=vaders.merge(df,how='left')

Sentiment visualisation (VADER)

Visualizing the positivity and negativity between 1 star and 5 stars
review



Sentiment visualisation (VADER)



Model Implementation & Architecture

Model Selection: RoBERTa Transformer

In [14]:


```
from transformers import AutoTokenizer
from transformers import AutoModelForSequenceClassification
from scipy.special import softmax
```


In [15]:

```
import warnings
warnings.filterwarnings("ignore")
MODEL = f"cardiffnlp/twitter-roberta-base-sentiment"
tokenizer = AutoTokenizer.from_pretrained(MODEL)
model = AutoModelForSequenceClassification.from_pretrained(MODEL)
```



config.json: 100%  747/747 [00:00<00:00, 81.7kB/s]

vocab.json:  899k/? [00:00<00:00, 31.3MB/s]

merges.txt:  456k/? [00:00<00:00, 32.2MB/s]

special_tokens_map.json: 100%  150/150 [00:00<00:00, 18.6kB/s]

Implementation continue

In [16]:

```
import torch
from torch.nn.functional import softmax
def polarity_scores_roberta(example):
    encoded_text = tokenizer(example, return_tensors='pt', truncation=True)
    with torch.no_grad():
        output = model(**encoded_text)
    scores_tensor = output[0][0].detach().cpu()
    scores = softmax(scores_tensor, dim=-1).numpy()
    scores_dict = {
        'roberta_neg': scores[0],
        'roberta_neu': scores[1],
        'roberta_pos': scores[2]
    }
    return scores_dict

res = {}
for i, row in tqdm(df.head(1000).iterrows(), total=1000):
    try:
        text = str(row['Text'])[:500] if pd.notnull(row['Text']) else ""
        myid = row['Id']

        vader_result = sia.polarity_scores(text)
        vader_result_rename = {f"vader_{k}": v for k, v in vader_result.items()}

        roberta_result = polarity_scores_roberta(text)
        both = {**vader_result_rename, **roberta_result}
        res[myid] = both

    except Exception as e:
        print(f'Error for id {myid}: {e}')
```

Technical Architecture:

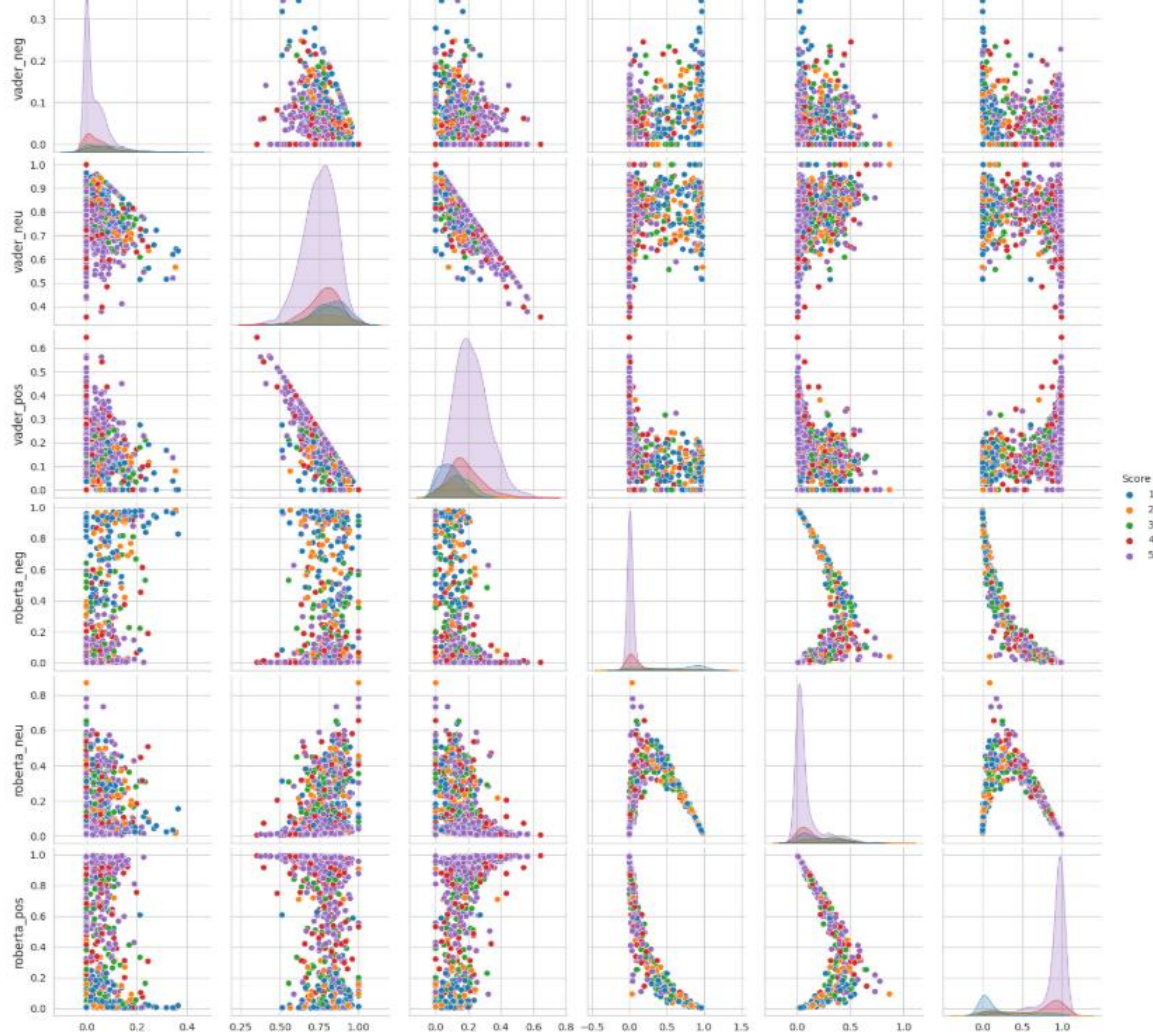
Input Processing: Text tokenization and encoding

Feature Extraction: RoBERTa embedding generation

Classification: Sentiment prediction with confidence scores

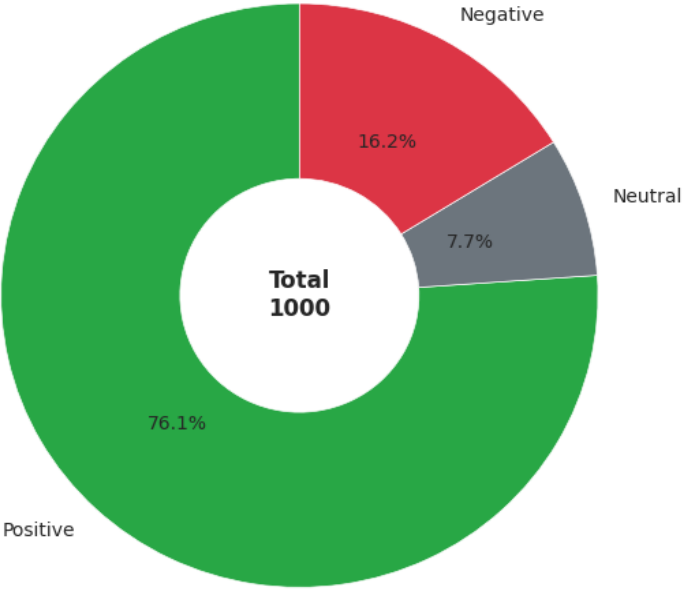
Post-processing: Result interpretation and validation

Comparing results between two models

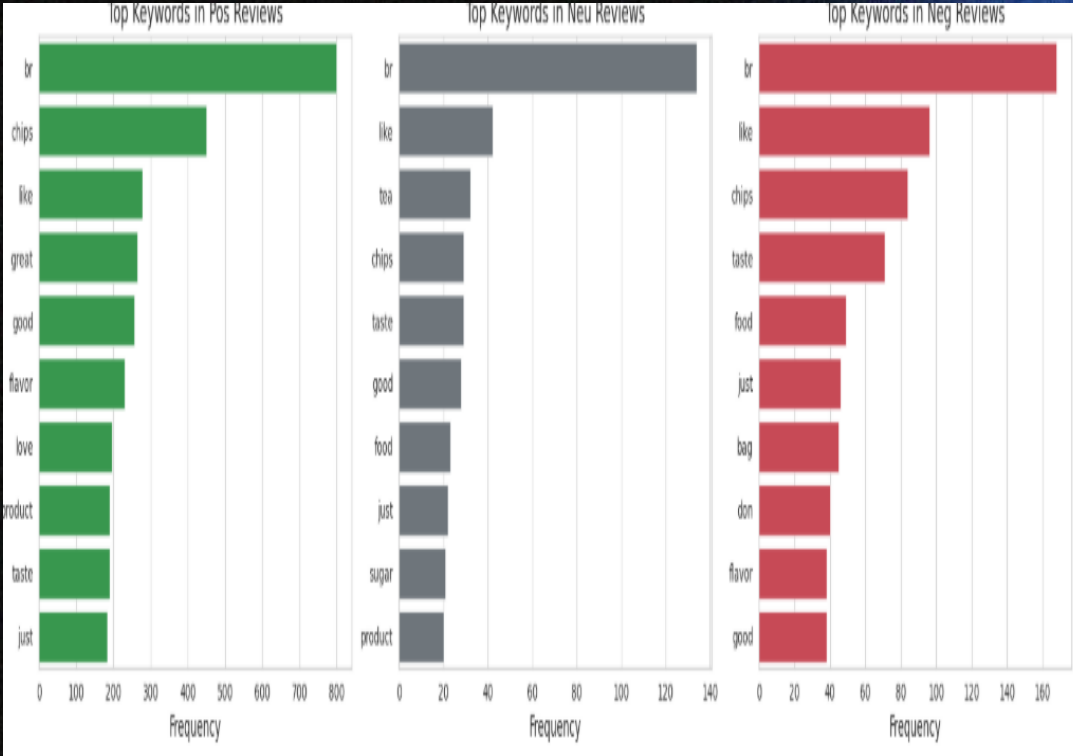


Sentiment Analysis Results

Roberta Sentiment Distribution

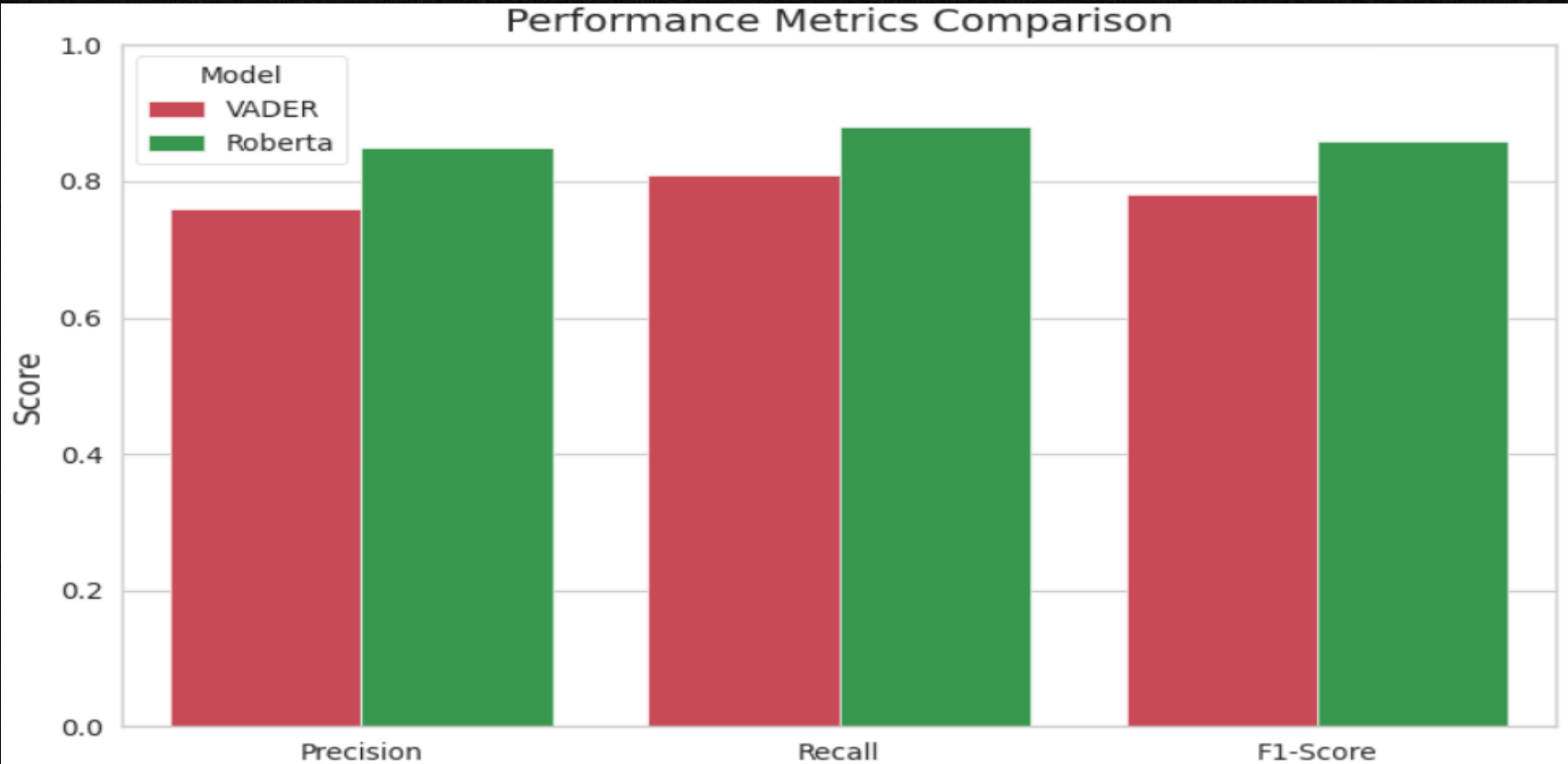


Most are positive

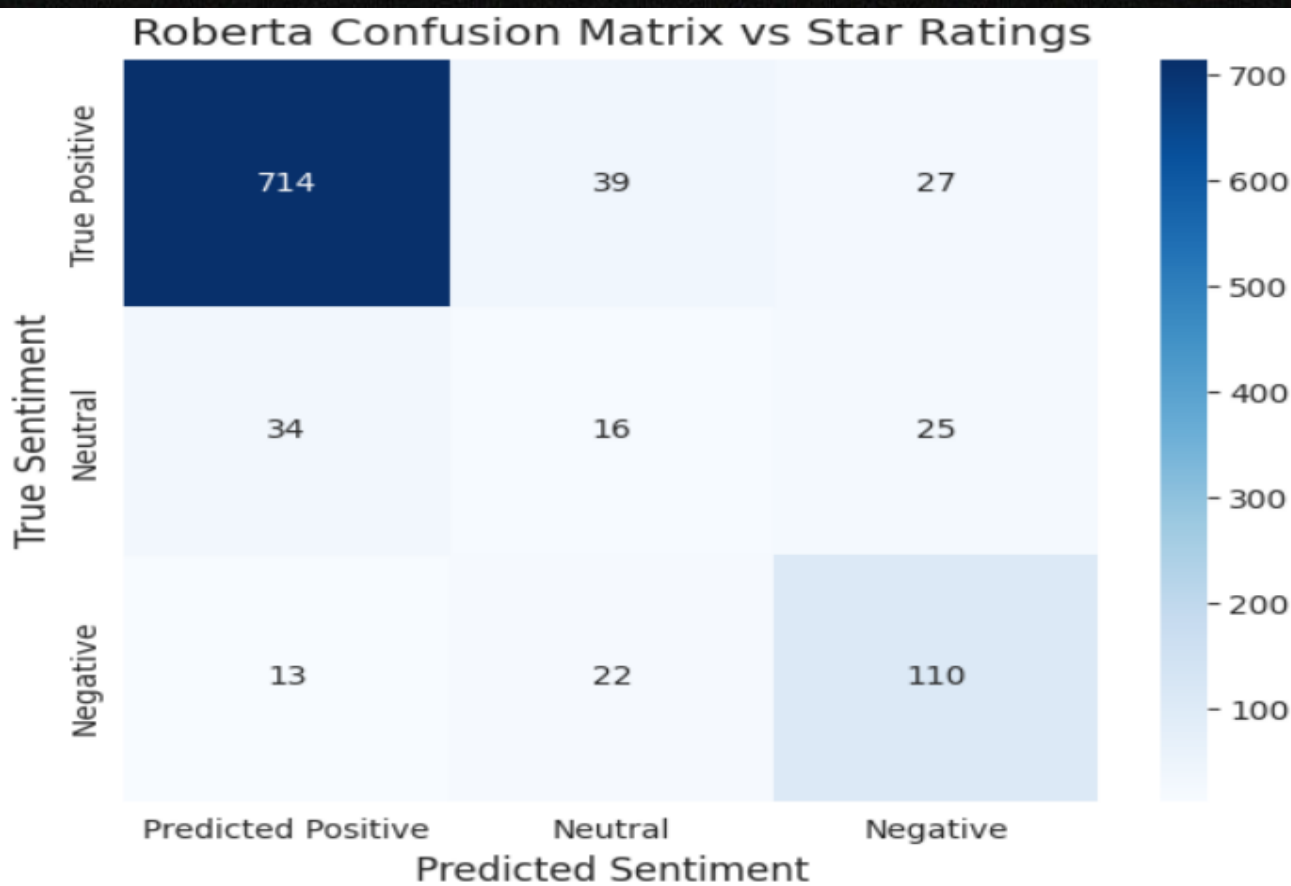


Frequency of words in pos , neg , neu

Performance Comparison



RoBERTa confusion metric



Business Intelligence & Recommendations

Actionable Business Insights

Based on the sentiment analysis findings, several strategic opportunities emerged:

Immediate Actions:

1. Quality Control Enhancement

- Focus on freshness and packaging improvements
- Address common quality complaints identified in negative reviews

2. Customer Service Optimization

- Prioritize delivery and shipping experience improvements
- Implement proactive communication for order updates

3. Product Development Guidance

- Leverage positive sentiment drivers in new product features
- Address specific pain points mentioned in negative feedback

Long-term Strategic Opportunities:

Marketing and Communication Strategy:

- **Messaging Focus:** Emphasize quality, freshness, and value proposition
- **Seasonal Campaigns:** Leverage high-satisfaction periods for promotions
- **Customer Testimonials:** Use positive sentiment language in marketing materials

Product Portfolio Optimization:

- **Category Performance:** Prioritize high-satisfaction product lines
- **Inventory Management:** Optimize stock levels based on sentiment-driven demand
- **New Product Development:** Use sentiment insights to guide feature prioritization

Future Enhancements & Scalability

Immediate Next Steps:

1. Aspect-Based Sentiment Analysis

Objective: Identify specific product features driving sentiment

Implementation: Multi-label classification for detailed insights

Business Value: Granular feedback for product development

2. Real-time Sentiment Monitoring

Objective: Continuous sentiment tracking for immediate response

Implementation: Streaming data pipeline with automated alerts

Business Value: Proactive customer experience management

3. Multi-language Support

Objective: Expand analysis to international markets

Implementation: Cross-lingual transfer learning

Business Value: Global customer intelligence capabilities

Scalability Considerations:

Cloud Infrastructure: Migration to scalable cloud computing platforms

API Development: Creation of sentiment analysis service endpoints

Dashboard Integration: Real-time business intelligence visualization

Technical Roadmap:
Phase 1: Enhanced preprocessing and feature engineering

Phase 2: Advanced model architectures and ensemble methods

Phase 3: Production deployment and monitoring systems

Project Summary & Professional Impact

Advanced Analytics Competencies:

- **Machine Learning:** State-of-the-art model implementation and optimization
- **Natural Language Processing:** Comprehensive text analysis and linguistic feature extraction
- **Data Engineering:** Efficient data pipeline creation and management
- **Statistical Analysis:** Rigorous hypothesis testing and validation methodology
- **Business Intelligence:** Translation of technical findings into strategic insights

Professional Development Outcomes:

- **End-to-End Project Management:** Complete workflow from conception to deployment
- **Technical Problem-Solving:** Creative solutions to complex analytical challenges
- **Communication Skills:** Clear documentation and insight presentation
- **Business Acumen:** Understanding of commercial applications for technical analysis

Key Deliverables:

1. **Production-Ready Model:** Deployable sentiment analysis solution
2. **Business Intelligence:** Actionable insights for strategic decision-making
3. **Technical Documentation:** Comprehensive methodology and implementation guide

Professional Value Proposition:

This project demonstrates the ability to bridge technical expertise with business needs, transforming complex data into strategic advantages. The combination of advanced NLP techniques, rigorous analytical methodology, and clear business applications reflects the kind of analytical thinking that drives data-driven decision making in today's competitive marketplace.

Thank You!

This comprehensive analysis showcases my commitment to leveraging advanced data science techniques for meaningful business impact. The combination of technical excellence and business insight demonstrated here reflects my approach to solving complex analytical challenges in professional environments.