# Amazon Product Reviews Analysis Project

# 2. Project Overview

## Dataset Description

The project utilizes the Amazon Product Reviews dataset from Kaggle, which contains customer reviews for various Amazon products. The dataset includes key information such as: - Product names and categories - Customer reviews and ratings (1-5 stars) - Review text and metadata - Product categories and descriptions

## Initial Hypothesis

1. Product reviews can be effectively categorized into distinct sentiment classes (positive, negative, neutral) based on both star ratings and review text.
2. Products can be meaningfully clustered into 4-6 meta-categories based on their descriptions and reviews.
3. Review sentiment patterns within product categories can reveal valuable insights about product performance and customer satisfaction.

## Analysis Structure and Process

1. **Data Loading and Initial Assessment**
   - Loaded the dataset using pandas
   - Performed initial data exploration and statistics
   - Identified key columns and data types

```python
df = pd.read_csv("dataset.csv")
df.columns
df['reviews.rating'].value_counts()
```

2. **Feature Engineering**
   - Mapped star ratings to sentiment classes:
     - 1-2 stars → Negative
     - 3 stars → Neutral
     - 4-5 stars → Positive
   - Created text features for clustering
   - Generated product embeddings using SentenceTransformer
3. **Model Development**
   - Implemented sentiment analysis using DistilBERT
   - Developed product clustering using K-means
   - Created summarization pipeline using BART

# 3. Data Wrangling and Cleaning

## Major Challenges

1. **Missing Data**

```python
df.isnull().sum()
```

   - Handled missing review texts

- o Dealt with missing product categories
- o Managed incomplete ratings

2. **Duplicate Reviews**

```python
unique_count = df['name'].nunique()
duplicate_count = df['name'].duplicated().sum()
```

- o Identified and removed duplicate reviews
- o Preserved unique product information
- o Maintained data integrity

3. **Text Preprocessing**

```python
def preprocess_text(text):
    # Lowercase
    text = text.lower()
    # Remove HTML tags
    text = re.sub(r'<.*?>', '', text)
    # Remove URLs
    text = re.sub(r"http\S+|www\S+|https\S+", '', text)
    # Remove user mentions and hashtags
    text = re.sub(r'\@\w+|\#', '', text)
    # Remove punctuation
    text = text.translate(str.maketrans('', '', string.punctuation))
    # Remove numbers
    text = re.sub(r'\d+', '', text)
    # Remove extra whitespace
    text = re.sub(r'\s+', ' ', text).strip()
    return text
```

## Data Enrichment Methods

1. **Text Feature Extraction**

- o Used SentenceTransformer for generating embeddings
- o Created TF-IDF features for text analysis

```python
model = SentenceTransformer('all-MiniLM-L6-v2')
embeddings = model.encode(product_texts, show_progress_bar=True)
```

2. **Category Clustering**

```python
kmeans = KMeans(n_clusters=num_clusters, random_state=42)
clusters = kmeans.fit_predict(embeddings)
```

- o Implemented K-means clustering
- o Generated cluster labels using most common words
- o Created meaningful category names

3. **Sentiment Analysis Enhancement**

- o Used TextBlob for additional sentiment features
- o Extracted pros and cons from reviews

3

```python
def analyze_sentiment(reviews):
    pros = defaultdict(int)
    cons = defaultdict(int)
    for review in reviews:
        blob = TextBlob(review)
        for sentence in blob.sentences:
            polarity = sentence.sentiment.polarity
            # ... sentiment analysis logic ...
```

## Challenge Resolution

1. **Data Quality Improvements**
   - Implemented robust text preprocessing
   - Handled multilingual reviews
   - Standardized product categories
2. **Feature Engineering Solutions**
   - Created meaningful product clusters
   - Generated sentiment scores
   - Extracted key product attributes
3. **Performance Optimization**
   - Used efficient data structures
   - Implemented batch processing
   - Optimized memory usage

The data cleaning and preprocessing steps were crucial for ensuring high-quality input for our models and generating meaningful insights from the reviews.

# 4. Exploratory Data Analysis

## Analysis Methods

1. **Distribution Analysis**

```python
# Rating Distribution
sns.set_theme(style='darkgrid', font_scale=1.15, palette="Set3")
ax = sns.countplot(x='reviews.rating', data=df)
for p in ax.patches:
    ax.annotate('{}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))
```

2. **Product Category Analysis**

```python
# Product Category Distribution
plt.figure(figsize=(10, 6))
scatter = plt.scatter(reduced_embeddings[:, 0], reduced_embeddings[:, 1],
            c=clusters, cmap='viridis', alpha=0.6)
plt.title('Product Category Clustering')
plt.colorbar(scatter)
```

3. **Sentiment Patterns**

- o Analyzed sentiment distribution across categories
- o Investigated rating patterns over time
- o Examined correlation between review length and sentiment

## Key Insights

1. **Rating Distribution**
   - o Most reviews are positively skewed (4-5 stars)
   - o Small percentage of extremely negative reviews (1-2 stars)
   - o Neutral reviews (3 stars) are relatively uncommon
2. **Category Patterns**
   - o Identified 4 distinct product clusters
   - o Electronics and accessories form the largest category
   - o Some categories show more polarized reviews than others
3. **Review Characteristics**
   - o Longer reviews tend to be more detailed and balanced
   - o Negative reviews often contain more specific details
   - o Positive reviews frequently use common positive phrases

# 5. Teamwork & Project Management

## Workflow Execution

1. **Original Plan vs. Reality**
   - o Successfully implemented core functionality
   - o Added additional features:
     - Enhanced visualization
     - More sophisticated clustering
     - Improved sentiment analysis
2. **Timeline Management**
   - o Met most deadlines
   - o Some delays in model optimization
   - o Successfully adapted to challenges

## Team Collaboration

1. **Strengths**
   - o Clear communication
   - o Regular code reviews
   - o Effective task distribution
2. **Areas for Improvement**
   - o Better documentation practices
   - o More frequent progress updates
   - o Enhanced version control practices

## Risk Management

1. **Identified Risks**
   - o Data quality issues
   - o Model performance
   - o Technical dependencies
2. **Mitigation Strategies**
   - o Regular data validation
   - o Model performance monitoring
   - o Backup solutions for critical components

# 6. Major Obstacles

## Primary Challenges

1. **Model Performance Issues**
   - o Initial sentiment analysis accuracy was lower than expected
   - o Clustering results needed refinement
   - o Resource constraints with large models
2. **Resolution Steps**
   - o Implemented more sophisticated preprocessing
   - o Used lighter, more efficient models
   - o Enhanced feature engineering

## Lessons Learned

1. **Technical Insights**
   - o Importance of thorough data preprocessing
   - o Value of efficient model selection
   - o Need for robust error handling
2. **Process Improvements**
   - o Better initial planning
   - o More comprehensive testing
   - o Regular performance monitoring

## Future Improvements

1. **Potential Enhancements**
   - o Implement more sophisticated clustering
   - o Add multi-language support
   - o Enhance visualization capabilities
2. **Alternative Approaches**
   - o Consider different model architectures
   - o Explore additional data sources
   - o Implement more automated testing

# 7. Conclusion and Insights

## Hypothesis Evaluation

1. **Sentiment Classification**
   - Successfully implemented three-class sentiment analysis
   - Achieved good accuracy in rating prediction
   - Validated correlation between ratings and review text
2. **Category Clustering**
   - Successfully identified meaningful product categories
   - Clustering provided useful insights
   - Some categories showed distinct review patterns

## Key Findings

1. **Product Insights**
   - Identified key factors in product satisfaction
   - Discovered common complaint patterns
   - Found correlation between price and review sentiment
2. **Review Patterns**
   - Longer reviews tend to be more informative
   - Sentiment patterns vary by category
   - Time-based trends in review behavior

## Implications

1. **Business Impact**
   - Improved product categorization
   - Better understanding of customer satisfaction
   - Identified areas for product improvement
2. **Technical Implications**
   - Demonstrated effectiveness of transformer models
   - Validated clustering approach
   - Identified areas for model improvement

## Unanswered Questions

1. **Technical Limitations**
   - Impact of review age on relevance
   - Cross-category product relationships
   - Long-term sentiment trends
2. **Future Research**
   - Multi-language review analysis
   - Deep dive into specific categories
   - Time-series analysis of ratings

The project successfully demonstrated the value of NLP in analyzing product reviews, while also highlighting areas for future improvement and research.