

# NLP Analysis of Meesho Product Reviews

Group 47

## **Team Members:**

Rohit Raju Kamble (2411AI61)

Satyam Kumar (2411AI44)

Abhisar Anand (2411AI43)

## **Guided by:**

Dr. Sourav Kumar Dandapat

April 21, 2025

# Agenda

Introduction

Data Collection & Preprocessing

Feature Engineering

Visualizations

Model Training

Model Evaluation

Conclusion

# Problem Statement

- ▶ Analysis of Meesho product reviews using NLP techniques
- ▶ Objectives:
  - ▶ Perform sentiment analysis on customer reviews
  - ▶ Extract key aspects and opinions
  - ▶ Build classification models to predict sentiment
  - ▶ Visualize insights from the data
- ▶ Dataset: 24,999 product reviews from Meesho

# Data Collection

- ▶ Scraped reviews using Apify client
- ▶ Combined multiple CSV files into one dataset
- ▶ Initial dataset columns:

```
1 Index(['productUrl', 'review/author/name',  
2 'review/comments', 'review/product_name',  
3 'review/rating', 'scrapedAt'], dtype='object')
```

# Text Cleaning Pipeline

```
1 def clean_text(text):  
2     text = text.lower()  
3     text = re.sub(r'http\S+|www\S+|https\S+', '', text  
4 )  
5     text = re.sub(r'\\@\\w+|\\#', '', text)  
6     text = re.sub(r'<.*?>', '', text)  
7     text = re.sub(r'[^\\w\\s]', '', text)  
8     text = re.sub(r'\\d+', '', text)  
9     text = re.sub(r'\\s+', ' ', text).strip()  
10    return text
```

Additional steps:

- ▶ Tokenization and lemmatization
- ▶ Stopword removal
- ▶ Spelling correction

# N-gram Extraction

```
1 def get_top_ngrams(text_series, n=2, top_n=10):
2     vec = CountVectorizer(ngram_range=(n, n)).fit(
3         text_series)
4     bag_of_words = vec.transform(text_series)
5     sum_words = bag_of_words.sum(axis=0)
6     words_freq = [(word, sum_words[0, idx])
7                    for word, idx in vec.vocabulary_.
8                    items()]
9     return sorted(words_freq, key=lambda x: x[1],
10                   reverse=True)[:top_n]
```

# Top N-grams

## Top Bigrams:

- ▶ good price (795)
- ▶ low quality (769)
- ▶ suit need (757)
- ▶ poor quality (568)
- ▶ daily use (566)

## Top Trigrams:

- ▶ fell apart cheap (501)
- ▶ highly functional great (341)
- ▶ excellent daily use (317)
- ▶ daily use love (317)
- ▶ poor performance worth (310)

# Top Product Aspects

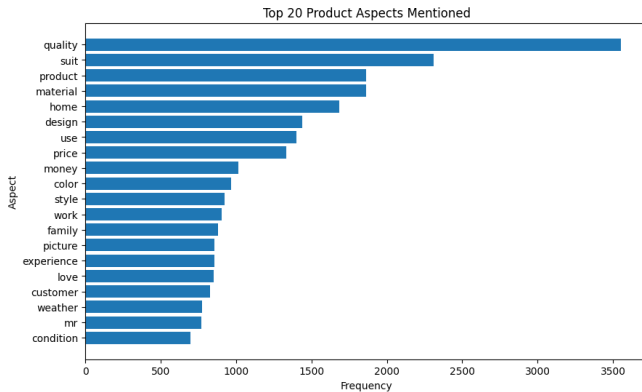


Figure: Top Product Aspects



## Word Cloud



Figure: Word cloud of most frequent terms in reviews

# Feature Extraction Methods

```
1 # Bag of Words
2 bow_vectorizer = CountVectorizer(max_features=5000)
3 X_train_bow = bow_vectorizer.fit_transform(X_train)
4
5 # TF-IDF
6 tfidf_vectorizer = TfidfVectorizer(
7     max_features=5000,
8     ngram_range=(1, 2))
9 X_train_tfidf = tfidf_vectorizer.fit_transform(X_train
10 )
11
12 # Train-test split
13 X_train, X_test, y_train, y_test = train_test_split(
14     X, y, test_size=0.2, random_state=42)
```

# Model Implementation

```
1 def train_and_evaluate(model, X_train, X_test,
2                           y_train, y_test, model_name):
3     model.fit(X_train, y_train)
4     y_pred = model.predict(X_test)
5     print(classification_report(y_test, y_pred))
6     print(f"Accuracy: {accuracy_score(y_test, y_pred)
7           :.4f}")
8     return model
9
10 # Example: SVM with TF-IDF
11 svm_tfidf = SVC(kernel='linear',
12                  probability=True,
13                  random_state=42)
14 train_and_evaluate(svm_tfidf, X_train_tfidf,
15                     X_test_tfidf, y_train, y_test,
16                     "SVM with TF-IDF")
```

# Model Performance Summary

Model	Features	Accuracy
Naive Bayes	BoW	0.6956
Naive Bayes	TF-IDF	0.6976
Random Forest	TF-IDF	0.6956
XGBoost	TF-IDF	0.6956
SVM	TF-IDF	0.7016

# Detailed Model Performance

## Best Model: SVM with TF-IDF

- ▶ Precision: 0.74 (class 0), 0.66 (class 1)
- ▶ Recall: 0.73 (class 0), 0.66 (class 1)
- ▶ F1-score: 0.73 (class 0), 0.66 (class 1)
- ▶ Accuracy: 0.7016

## Model Parameters

- ▶ Kernel: linear
- ▶ Probability: True
- ▶ Random state: 42
- ▶ Features: TF-IDF (5000)
- ▶ N-gram range: (1,2)

# Classification Report

```
1 SVM with TF-IDF Performance:
2           precision    recall  f1-score   support
3
4      0       0.74       0.73       0.73       2817
5      1       0.66       0.66       0.66       2183
6
7   accuracy                0.70       5000
8   macro avg              0.70       0.70       0.70       5000
9   weighted avg           0.70       0.70       0.70       5000
```

# Key Findings

- ▶ SVM with TF-IDF performed best (70.16% accuracy)
- ▶ All models showed similar performance ( 70% accuracy)
- ▶ Key phrases reveal product quality and value concerns
- ▶ "Good price" and "low quality" most common bigrams
- ▶ Model performance consistent across metrics

# Future Work

- ▶ Analyze temporal trends in reviews
- ▶ Build recommendation system based on aspects
- ▶ Hyperparameter tuning for better performance



# Contributions

- ▶ Abhisar - Automatic Data Collection
- ▶ Satyam - Text prepossessing and visualization
- ▶ Rohit - feature engineering and Model Building

Thank You!