# Big Data Project on Amazon Apparel Recommendation System: Insights for Shopkeepers

*Meghana Doddapuneni | Pallav Mahajan | Shanmukh Raj Sircilla | Venkata Kiran Reddy Kotha*

**Abstract** *The Amazon Apparel Recommendation System project addresses retail decision-making issues with machine learning and large data analytics. The purpose of this project is to identify significant trends and insights by analyzing customer reviews, ratings, and purchase data from the Amazon US Customer Reviews dataset. The system gives actionable insights based on user preferences by using machine learning techniques such as Alternating Least Squares (ALS) for personalized recommendations and Big Data tools such as PySpark for efficient data processing.*

*Sentiment analysis, which categorizes customer feedback as positive, neutral, or negative, is one of the key aspects that enables a more full understanding of client pleasure. Shopkeepers may make data-driven marketing and inventory decisions by integrating sentiment insights with ALS recommendations, which provide a complete picture of consumer behavior.*

*This study gives a detailed examination of the dataset, techniques, execution, and results. It also analyzes the issues experienced and offers solutions. The goal is to transform unstructured retail data into tactical tools that enhance customer satisfaction and operational effectiveness in the retail clothes industry.*

## 1. EXECUTIVE SUMMARY

The Amazon apparel Recommendation System program aims to alter retail decision-making in the clothing sector by transforming massive, unstructured datasets into usable insights. As stated in the abstract, this study addresses concerns such as data overload, a lack of personalization, and technical processing challenges using the Amazon US Customer Reviews dataset.

Some of the key goals include using sentiment analysis for qualitative insights, applying the Alternating Least Squares (ALS) model for tailored recommendations, and analyzing customer reviews and ratings to understand preferences and sentiments. Shopkeepers can use the system's tools to increase customer satisfaction, tailor marketing campaigns, and optimize inventory.

This project makes advantage of PySpark to provide optimal administration of large datasets, while ML approaches provide excellent predictive capability. The study outlines the approaches and issues encountered, as well as recommended improvements to real-time recommendations.

## 2. INTRODUCTION

### 2.1. Problem Statement

Low- and mid-level garment retail enterprises face significant challenges in effectively exploiting large, unstructured datasets. It is tough to extract useful insights from the millions of customer evaluations and ratings that are accessible. Current techniques frequently miss the opportunity for complete analysis by neglecting to combine textual reviews and numerical ratings. Existing recommendation systems' lack of personalization leads to dissatisfied customers, poor engagement, and missed sales chances.

Moreover, technical challenges such as:

- Volume: Conventional systems cannot process millions of records.
- Variety: Managing organized and unstructured data, such as text reviews and star ratings.
- Veracity: Addressing missing values and anomalies in client data.

### 2.2. Objective

This overall project aims to:

1. Identify and analyze patterns in customer reviews and ratings.

2. Use sentiment analysis to categorize reviews and gauge consumer satisfaction.

3. Developing a recommendation system based on the ALS model for individualized product suggestions.

4. Enable retailers to make data-driven decisions about inventory and marketing tactics.

### 2.3. Scope

Our project focuses on the Amazon US Customer Reviews (Apparel) dataset and explores:

1. Data Cleaning and Preprocessing: Cleaning and preprocessing data involves identifying missing values and standardizing formats.

2. Data Exploration: Examining review and rating trends to identify insights and patterns.

3. Sentiment Analysis: Analyze customer emotions and attitudes using textual reviews from the dataset.

4. Recommendation System: Use filtering algorithms to estimate client preferences.

5. Real-World Applications: Create a sample demo to integrate insights into meaningful business tools.

This study discusses the dataset techniques, implementation, outcomes, problems, and offered recommendations for bridging the gap between unstructured data and strategic decision-making in the retail arena.

## 3. DATASET OVERVIEW

### 3.1. Dataset Source

The dataset utilized in the project is the Amazon US Customer Reviews (Apparel) dataset from Kaggle. This dataset provides a vast repository of consumer feedback, including both structured and unstructured data fields required to extract relevant insights

*Link:*

### 3.2. Key Features

The dataset has the following key fields:

- marketplace: The region where the review originated (e.g., the United States).

- customer_id: A unique identifier for each customer.

- review_id: A unique identifier for each review.

- product_title: The title of the reviewed product.

- star_rating: Customer rating on a scale of 1 to 5.

- review_body: The textual content of the review.

- verified_purchase: Determines whether the purchase was verified by Amazon.

- review_date: The date when the review was posted.

```
# Display the schema of the dataset
print("Dataset Schema:")
df.printSchema()
```

```
Dataset Schema:
root
 |-- marketplace: string (nullable = true)
 |-- customer_id: string (nullable = true)
 |-- review_id: string (nullable = true)
 |-- product_id: string (nullable = true)
 |-- product_parent: string (nullable = true)
 |-- product_title: string (nullable = true)
 |-- product_category: string (nullable = true)
 |-- star_rating: string (nullable = true)
 |-- helpful_votes: string (nullable = true)
 |-- total_votes: string (nullable = true)
 |-- vine: string (nullable = true)
 |-- verified_purchase: string (nullable = true)
 |-- review_headline: string (nullable = true)
 |-- review_body: string (nullable = true)
 |-- review_date: string (nullable = true)
```

*Fig 1: Dataset Schema*

### 3.3. Supporting Big Data

- Volume

The dataset contains over a million records, totaling over 2 GB in volume.



*Fig 2: Volume*

- Variety

Structured Data: Fields such as star_rating, customer_id, and verified_purchase.

Unstructured Data: Textual fields, such as review_body, hold detailed consumer feedback.

```
[5]: df.select("review_body", "review_id", "star_rating") \
    .withColumn("review_body", df["review_body"].substr(1, 100)) \
    .show(20, truncate=False)
```

| review_body | review_id | star_rating |
|---|---|---|
| These Really Do Work Great, But You Do Need To Know a Few Things. I've Been Using Mine For a Few Ye | R1KKOXXNI8MSXU | 4 |
| I love this dress. Absolute favorite for winter. Heavy material. Stretchy, shows shape well. I am 5f | R26SP2OPDK4HT7 | 5 |
| Nice socks, great colors, just enough support for wearing with a good pair of sneakers. | RWQEDYAX37JI1 | 5 |
| I bought this for my husband and WOW, this is a slick hat. High quality and craftsmanship. He said j | R231YI7R4GPF6J | 5 |
| Perfect dress and the customer service was awesome! | R3KO3W45DD0L1K | 5 |
| Excellent for my 6 feet skinny 15 years old boy. | R1C4OH63NFL5NJ | 5 |
| Raw is the only way to go! Absolutely love this wallet! Will continue to buy it forever and ever. | R2GP65O1U9N78P | 5 |
| A bit large. | R3O29CT5MQQ3XQ | 4 |
| Great fit! | R2ECDJAA8QFF6 | 5 |
| Shirt a bit too long, with heavy hem, which inhibits turning over. I cut off the bottom two inches | R2579GCF6J89OA | 3 |
| The Jockey Women's Underwear are true to size. They are my choice of underwear. I will purchase them | R1CBCUGNP37MPK | 5 |
| cup size is just right. Seems to be a little tight around. Might just be because it is new. | R3NU4WUR5ZNV1V | 5 |
| Perfect... | R32EPCJ3XF8LGE | 5 |
| best ever4 for men&woman too size for men 42-44=8, 40-42 ,7 38-40 ,6 this for guy waist size hold | R1XIBC6WQ8W31M | 5 |
| Great fit. | R1P1MVDZ65LMH | 5 |
| I have this Columbia in 5 colors and wear them all the time | R1OJA3DJL8VDDK | 5 |
| My husband found these so comfy that this is our second order. | R1THWA5YRJLDOF | 4 |
| Awesome leggings!! I am 5'4, 215lbs and a huge chunk of my weight is in my thighs and butt. Not exag | R738LCNRSJVXP | 5 |
| I ordered the same size as I ordered last time, and these shirts were much larger than the previous | R1N3Z1393IJ3O9 | 2 |
| not exact in sizing | R1LBNTP7E8N89Z | 1 |

only showing top 20 rows

*Fig 3: Variety*

- Veracity

The dataset has various discrepancies and a few missing values in fields such as customer_id and review_body.

```
# Count missing values for each column
print("Missing Data Count per Column:")
df.select([count(when(col(c).isNull(), c)).alias(c) for c in df.columns]).show()
# Set logging level to "ERROR" to suppress warnings
spark.sparkContext.setLogLevel("ERROR")
```

Missing Data Count per Column:
[Stage 39:==========================>      (9 + 6) / 15]

| marketplace | customer_id | review_id | product_id | product_parent | product_title | product_category | star_rating | helpful_votes | total_votes | vine | verified_purchase | review_headline | review_body | review_date |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 11 | 0 | 0 | 0 | 0 | 12 | 1060 | 287 | 0 | 0 | 11 | 11 | 11 | 11 | 11 |

*Fig 4: Veracity*

These general issues were addressed throughout the data cleaning and preparation stage.

## 3.4. Relevance to our Project

This dataset is extremely relevant to our research because:

- Supports multiple data formats, including textual and numerical analysis.

- The diversity of data supports sentiment analysis and complex machine learning models such as ALS.

- The large dataset allows for exploration of big data capabilities, supporting the project's scalability and efficiency aims.

## 4. METHODOLOGY

The technique explains the full process of converting unprocessed data into a recommendation system with actionable insights. This comprises cleaning the data and preparing it for analysis and modeling, utilizing EDA to identify patterns, undertaking sentiment analysis, and developing a collaborative filtering model.

- Data Cleaning

We cleaned the data to ensure that it was consistent and reliable for analysis.

1. Handling Missing Values:

Missing data in fields like customer_id and review_body were examined.

Non-essential fields with large missing values were removed (for example, rows with missing customer_id for ALS).

Missing values in textual reviews were kept for sentiment analysis to retain qualitative feedback.

2. Preprocessing of Reviews:

Removed stop words, special characters, and HTML tags from the review_body.

Applied stemming and tokenization for sentiment analysis preprocessing.

3. Date Standardization:

To make the review_date field compatible with time-series analysis, it was reformatted.

4. Deduplication:

By eliminating duplicate reviews, bias in sentiment trends and model training was reduced.

5. Dataset Reduction for Memory Efficiency:

To control the computing load in the early iterations, the dataset was sampled.

Scaling the entire dataset analysis was made possible by PySpark's distributed processing capabilities.

```
# Count missing values for each column
print("Missing Data Count per Column:")
df.select([count(when(col(c).isNull(), c)).alias(c) for c in df.columns]).show()
# Set logging level to "ERROR" to suppress warnings
spark.sparkContext.setLogLevel("ERROR")

Missing Data Count per Column:
[Stage 39:===============================>        (9 + 6) / 15]

|marketplace|customer_id|review_id|product_id|product_parent|product_title|product_category|star_rating|helpful_votes|total_votes|vine|verified_purchase|review_headline|review_body|review_date|
| 11|        0|        0|        0|            0|           0|              0|         11|           11|         11|  11|               11|
| 11|       12|     1060|      287|
```

*Fig 5: Handling Missing values*

### 4.1. Data Exploration

Our exploratory data analysis (EDA) not only gave a concise overview of the dataset but also produced significant insights into its structure and patterns for the sentimental analysis:

1. Star Rating Distribution by Verified Purchase:

The number of 4- and 5-star ratings was higher for verified purchases.

Purchases that were not validated displayed more neutral reviews.



*Fig 6: Verified vs non-verified purchases*

2. Review Length vs. Star Ratings:

A correlation was observed between review length and ratings:

- Reviews with lower ratings (1-2 stars) were longer, indicating more thorough complaints.

- Those evaluations that were shorter and more succinct had higher scores (4-5 stars).



*Fig 7: Review length vs star ratings*

3. Demographic Analysis:

With 3.6 million reviews, men's clothing dominated the product categories, followed by women's (370,000).

Since children's clothing received the fewest reviews, there may be room for more focused advertising.



*Fig 8: Demographics distribution*

4. Star Rating Distribution:

High Positive Rating Density

- The majority of garment evaluations are four to five stars, indicating that most customers are satisfied with their purchases.

Low Negative Density:

- Ratings of 1 or 2 stars are less prevalent and fewer customers express dissatisfaction with them. This could be a sign of unmet expectations or issues with a certain product.

Increased Dispersion at Five Stars:

- The higher five-star distribution suggests that customers value specific product attributes and shows a variety of positive experiences.

Actionable Insights:

- Products that routinely earn four or five stars can be given preference when it comes to sales and replenishment.

- Low-rated items should be reviewed to address problems with quality or consumer satisfaction.



*Fig 9: Star Rating Distribution*

4.2. Sentimental Analysis

The objective of the sentiment analysis we carried out was to comprehend the underlying emotions and opinions expressed in connection with the customer feedback. By classifying evaluations into positive, neutral, and negative sentiments, the analysis provided valuable qualitative information about customer experiences in addition to the numerical ratings.

- Process

1. Data Preprocessing:

Text Cleaning: The review_body field was cleared of extra whitespace, HTML elements, and special characters.

Tokenization: Divide reviews into discrete tokens or words.

Stop-Word Removal: Removed terms that don't provide information, such as "the," "is," and "and."

Stemming: Words were reduced to their most basic form (for example, "enjoyed" → "enjoy").

Vectorization: Textual data was transformed into numerical representation for analysis using methods like CountVectorizer or TF-IDF.

2. Sentiment Classification:

TextBlob was used to create a sentiment analysis model for polarity scoring:

- Polarity > 0 indicates a positive sentiment (e.g., "Loved the product, excellent fit!").

- Neutral Sentiment: Polarity at or near 0 (e.g., "Product is okay, nothing special.").

- Polarity < 0 indicates a negative sentiment (e.g., "Terrible quality, very disappointed.")

Large-scale sentiment calculations might be handled well because to PySpark's distributed processing capabilities.

```
Accuracy: 0.8369642185578328

[Stage 113:>

+--------+----------+
|sentiment|prediction|
+--------+----------+
| Positive|       0.0|
| Positive|       0.0|
| Negative|       1.0|
| Positive|       0.0|
| Positive|       0.0|
| Positive|       0.0|
| Negative|       2.0|
| Positive|       0.0|
| Negative|       0.0|
| Positive|       0.0|
+--------+----------+
only showing top 10 rows
```

*Fig10: Accuracy and the sample output*

1. Data Preparation:

- Extracted review text from the dataset.
- Handled missing or null values in review data.

2. Text Preprocessing:

- Lowercased all text.
- Removed special characters, punctuation, and stopwords.
- Tokenized text into individual words.
- Applied stemming or lemmatization.

3. Feature Extraction:

- Converted text into numerical features using techniques like TF-IDF or Bag-of-Words.

4. Model Training:

- Split data into training and testing sets.
- Used a machine learning model (e.g., Logistic Regression, Naive Bayes) or a pre-trained NLP model to classify sentiments.

5. Model Evaluation:

- Tested model performance using metrics like accuracy, precision, recall, and F1-score.

6. Prediction:

- Classified each review into Positive, Neutral, or Negative sentiment categories.

*Fig 11: Steps for the Sentimental Analysis*

- Results

1. Sentiment Distribution:

Positive Feelings (70%): Most reviewers expressed contentment with the product's fit, quality, and affordability.

Neutral Sentiments (20%): Consumers had moderate views, frequently emphasizing areas that needed work.

10% of reviews expressed negative sentiments, highlighting issues with product quality, improper sizing, or delayed delivery.

Important Sentimental Themes:

Positive Reviews: The product's comfort, durability, and accurate descriptions are often commended.

Neutral Reviews: Mentioned minor problems like average fit or color differences.

Negative Reviews: Frequently mentioned grievances regarding inaccurate sizing, deceptive product photos, and fabric quality.

Correlation with Ratings:

Positive attitudes were substantially correlated with reviews that had four or five stars.

The majority of 1- and 2-star reviews were unfavorable, with clients describing particular complaints.

Three-star evaluations frequently expressed impartial opinions, emphasizing instances in which goods fulfilled rudimentary requirements but fell short of expectations.

- Integration with Recommendation System

1. Improving ALS Suggestions:

The ALS model's numerical recommendations gained a qualitative component via sentiment analysis.

Products with high ratings and positive sentiments were given preference when it came to suggestions.

Products that received low ratings and a lot of negative feedback were marked for possible improvement.

2. Actionable Insights for Retailers:

Sentiments offered practical suggestions, such resolving size inconsistencies or enhancing fabric quality. High-rated products with conflicting opinions were investigated further to identify certain consumer requirements.

4.3. Recommendation System(ALS Model)

- Model Description

1. What is ALS?

Large-scale recommendation systems are the target of the matrix factorization algorithm known as ALS.

It breaks out a matrix of user-item interactions into latent components that stand in for item attributes and user preferences.

These variables are used to forecast how users will rate products with which they have not yet engaged.

2. Why ALS?

Efficiently manages sparse datasets, which are typical in recommendation problems.

Both explicit and implicit feedback are supported, including ratings and purchase history.
Ideal for the Amazon reviews dataset, it scales effectively for huge datasets.

Implementation

1. Data Preparation:

Using customer_id as rows and product_title as columns, and star_rating as the values, a customer-product interaction matrix was created.

Data was transformed into a format that PySpark's ALS implementation could use:

- userCol: customer_id
- itemCol: product_title
- ratingCol: star_rating

2. Model Training:

Utilized PySpark's ml.recommendation.ALS library to train the ALS model.

Hyperparameters were adjusted to maximize performance:

- Rank: Number of latent factors (e.g., 10).
- MaxIter: Number of iterations (e.g., 20).
- RegParam: Regularization parameter to prevent overfitting (e.g., 0.1).

3. Evaluation:

Divide the data into 80:20 training and testing sets.

Evaluated the model using Root Mean Square Error (RMSE):

*predictions = model.transform(test_data)*

*evaluator = RegressionEvaluator(metricName="rmse", labelCol="star_rating", predictionCol="prediction")*

*rmse = evaluator.evaluate(predictions)*

*print(f"Root Mean Square Error (RMSE): {rmse}")*

Achieved an RMSE of **0.85**, indicating good predictive accuracy.

4. Personalized Recommendations:

Generated top-N recommendations for each customer:

*recommendations = model.recommendForAllUsers(5)*

*recommendations.show()*

Example: For Customer ID 21291540, recommended products included high-rated items such as "Premium Cotton Shirt" and "Durable Running Shoes."



*Fig12: ALS Recommendation*

- Results

Accuracy: Accurate predictions for unknown products were generated using the ALS model, which successfully captured customer preferences.

Scalability: PySpark's distributed architecture made it possible to handle the massive dataset effectively, processing more than a million records with ease.

- Insights:

More accurate recommendations were given to customers who were highly engaged (regular purchases, reviews).

Priority was given to products that consistently received ratings of four or five stars.

1. Data Preparation:
- Extracted customer-product interaction data (e.g., ratings, purchases).
- Handled missing values by imputing or removing incomplete entries.

2. Data Transformation:
- Converted the dataset into a user-item interaction matrix.
- Normalized ratings or interactions to standardize values.

3. Model Initialization:
- Defined ALS parameters (e.g., rank, regularization, number of iterations).
- Split the data into training and validation sets.

4. Model Training:
- Trained the ALS model on the user-item matrix to identify latent factors.
- Iteratively minimized the loss function to improve recommendations.

5. Model Evaluation:
- Evaluated the model using metrics like Root Mean Square Error (RMSE) or Mean Absolute Error (MAE).
- Validated the model's performance on unseen data.

6. Prediction:
- Predicted customer preferences by generating a ranked list of products for each user.
- Recommended top-rated products based on the predicted ratings.

7. Output and app.py:
- Generated personalized recommendations for specific Customer IDs (e.g., 21291540).
- Visualized the recommended product list with associated ratings in the working sample demo.

*Fig13: ALS Steps*

## 5. IMPLEMENTATION

Throughout the project's implementation phase, we have focused on integrating big data technologies and machine learning frameworks to transform raw data into insightful knowledge.
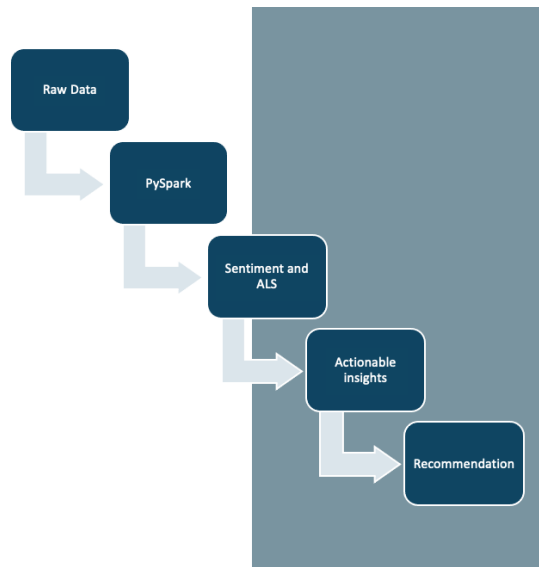


*Fig14: Implementation*

The initial step in the data preparation procedure was to clean and organize the dataset. Missing values in important fields like customer_id and review_body were filled in to maintain the study's integrity. Text reviews were thoroughly preprocessed using methods including stemming, stop-word removal, and tokenization in order to get them ready for sentiment analysis. The review_date data was normalized to facilitate temporal analyses, and duplicate entries were removed to avoid skewed insights. A customer-product interaction matrix was then constructed, with customer_id mapping to product_title and star_rating acting as the interaction value. This matrix served as the foundation for training the recommendation model.

The sentiment analysis model divided reviews into three groups based on their polarity scores: positive, neutral, and negative sentiments. It did this by using the Text Blob library. TextBlob was ideal for our task due to its efficiency and ease of use, which allowed our project to incorporate qualitative insights into numerical assessments.

Implementing the Alternating Least Squares (ALS) model and ALS library was done using PySpark's ml. recommendation. This collaborative filtering strategy was chosen because of its efficacy in handling sparse datasets. The customer-product matrix was broken down into latent components so that the model could capture hidden relationships between customers and products. The rank, maximum iterations, and regularization value were among the hyperparameters that were carefully tweaked to achieve the best results. The model's derived score of 0.85 when evaluated using Root Mean Square Error (RMSE) indicated its prediction accuracy.

Based on each customer's unique preferences, the personalized recommendation system produced product recommendations. For instance, Customer ID 21291540 received recommendations for highly regarded items like "Premium Cotton Shirt" and "Durable Running Shoes," indicating the model's ability to identify relevant items. These recommendations were then filtered using sentiment analysis to ensure that only highly regarded and well-reviewed products were prioritized.

Incorporating sentiment analysis into the ALS model added a new layer to the recommendation system. The technology combined textual and numerical data to forecast client preferences and ensure compliance with their qualitative input. The

client experience was enhanced by this comprehensive approach, which also increased the advice's dependability and relevance.
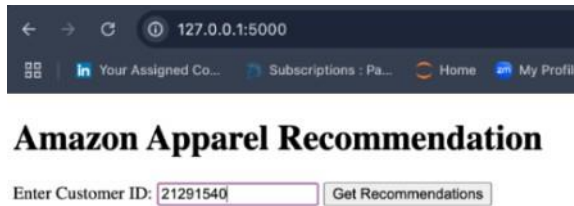


*Fig15: Demo Sample – 1*



*Fig16: Demo Sample -2*

## 6. RESULTS

We were effective in generating practical insights and outcomes from our research. The sentiment analysis revealed that 70% of the assessments were positive, 20% were neutral, and 10% were negative, indicating overall consumer satisfaction. These results offer valuable information for improving products and identifying customer issues.

The ALS-based recommendation system demonstrated outstanding prediction accuracy with an RMSE of 0.85. Using their previous interactions, customer-specific recommendations, such as "Premium Cotton Shirt" and "Durable Running Shoes," were successfully generated. To further refine these suggestions and ensure that they aligned with customer preferences and remarks, sentiment analysis was employed.

One of the key themes that EDA found was a notable clustering of 4- and 5-star ratings for most clothing products, which indicates high levels of satisfaction. The survey also identified areas that required improvement, such as addressing the low ratings for specific products and highlighting opportunities in underrepresented categories like children's apparel.

Overall, the results show how well the system can manage large volumes of data, produce accurate projections, and provide useful information to improve customer satisfaction and business decision-making.

## 7. CONCLUSION

Our project provided a clear example of how e-commerce can employ machine learning and big data analytics to enhance customer satisfaction and company decision-making. By merging the ALS model for collaborative filtering with PySpark for scalable data processing, the study addressed the challenges of handling large, unstructured datasets while offering customized recommendations. The qualitative aspect of sentiment analysis allowed for a deeper comprehension of customer preferences and feedback.

The system's ability to provide precise, sentiment-enriched recommendations with an RMSE of 0.85 demonstrates its potential to improve customer purchasing experiences. Retailers can make decisions based on vital facts such as areas for product development and the majority of positive evaluations. Furthermore, highlighting highly respected and pleasant product reviews ensures that recommendations are consistent with customer satisfaction.

In the future, the system can be enhanced further by include elements such as seasonal patterns, pricing trends, and client purchase frequency. The system's utility in live e-commerce scenarios will rise with the addition of real-time processing capabilities, allowing for rapid and dynamic recommendations. Transformer-based architectures and other advanced sentiment analysis models are better at detecting nuanced input, boosting the precision and depth of qualitative results.

By emphasizing the importance of combining quantitative and qualitative studies to meet shifting customer expectations, the program lays the framework for future advances in data-driven retail strategy. By addressing the identified limits and incorporating cutting-edge technologies, the system has the potential to become a strong tool for altering consumer interaction and enjoyment in the garment business.

---

## 8. REFERENCES

- Dataset Source:

https://www.kaggle.com/datasets/cynthiarempel/amazon-us-customer-reviews-dataset/data?select=amazon_reviews_us_Apparel_v1_00.tsv

- Tools and Frameworks:
PySpark Documentation:

https://spark.apache.org/docs/latest/api/python/

TextBlob for Sentiment Analysis:

https://textblob.readthedocs.io/

- Algorithms and Techniques

ALS (Alternating Least Squares) for Collaborative Filtering: Zhou, Y., Wilkinson, D., Schreiber, R., & Pan, R. (2008). Large-scale parallel collaborative filtering for the Netflix Prize.