

Danmarks
Tekniske
Universitet



Recommendation System for Skincare Products

02807 COMPUTATIONAL TOOLS FOR DATA SCIENCE

AUTHORS

Ana Marija Pavičić - s232468
Anna Sky Kastl Jensen - s194824
Raquel Moleiro Marques - s243636
Sree Keerthi Desu - s243933

December 5, 2024

Contents

Table of Contents

1	Introduction	1
1.1	Motivation	1
1.2	Approach	1
2	Datasets	1
2.1	Product Dataset	1
2.2	Review Dataset	2
3	Methods	2
3.1	Sentiment Analysis	2
3.2	Frequent Itemsets	2
3.3	Similar Items	3
3.3.1	Similarity of Highlights	3
3.3.2	Ingredient Similarity	3
3.3.3	Reciprocal Rank Fusion (RRF)	4
4	Recommendation System	4
5	Concluding remarks	4
A	Appendix	5
A.1	Group Member’s Contribution to the Project	5
A.2	Code	5
	References	18

1 Introduction

1.1 Motivation

In today's consumer-driven world, where corporations frequently flood the market with new products and aggressive promotions, choosing the right products can be overwhelming. Our recommendation system offers a simple solution by tailoring suggestions to a user's preferences and predicting potential responses to specific inputs. By presenting users with options they are likely to enjoy, these systems can reduce cognitive overload, making the decision-making process more efficient. As a result, users are more likely to purchase products they genuinely value, ultimately reducing wasteful consumption.

1.2 Approach

Collaborative filtering is a specific type of recommendation system that focuses on the relationships between users and products. Recommendations are generated by identifying similarities between user profiles, which are typically based on their product purchase history [1]. In essence, these recommendations suggest products that other users frequently buy alongside the input product. However, traditional collaborative filtering struggles with the cold start problem, where new users or products lack sufficient data for accurate recommendations.

For our skincare recommendation system, the goal is to provide recommendations based on a singular input skincare product, which simulates the situation where there is a lack of a solid user profile. To achieve this, we first use the A-priori algorithm to analyze product reviews and identify frequent itemsets — groups of products often bought together. This step narrows the search space, focusing on products that are contextually relevant to the user's input.

Next, sentiment analysis is applied to user reviews to gauge overall satisfaction and identify potential adverse reactions. Additionally, we refine the recommendations by calculating similarity scores based on two critical aspects: product highlights (e.g., claims like "brightening" or "hydrating") and ingredients (e.g., hyaluronic acid or salicylic acid). Combining these factors would ensure that the recommended products align with user preferences. By doing so, we aim to answer the following research question:

“How can we deliver product recommendations for customers with no prior user data?”

2 Datasets

The chosen datasets are available on Kaggle under the name "**Sephora Products and Skincare Reviews**" [2]. It consists of a product information file and several files containing different reviews. Please note that all our datasets and code are available in the GitHub repository [3].

2.1 Product Dataset

There are approximately 8,000 products available on Sephora, and their corresponding information includes price, highlights, ingredient list, average rating, size, etc. However, not all of these products are related to skincare. Using the primary category, only skincare products were filtered resulting in 2,420 unique product IDs. Among the remaining rows, there are also mini versions and gift sets made of products that already exist in the dataset and therefore were also removed to ensure that there were no duplicates for the association rules to avoid recommending the same product. Rows that have only either highlights or ingredients were removed, as both columns are necessary for obtaining similar items. This action resulted in 1,583 unique products.

The ingredient list for each product was preprocessed to make sure that ingredients like water (that were listed in multiple different ways such as “aqua”) were standardized, and special characters towards the beginning and the end of the list were removed. The uniformity of the column would help the similarity score further in the report.

2.2 Review Dataset

The multiple files containing reviews were merged into one. Most of the products seem to have at least 65 reviews; hence, the review dataset size is considerably large and has over a million reviews. The dataset was reduced to include the reviews corresponding to only the 1583 products obtained from data cleaning of the product dataset. Additionally, the dataset was filtered to keep only the reviews that had 3, 4, or 5 stars as we do not want to consider any negative reviews before finding association rules among them. After all this preprocessing, we end up with 700,000 reviews in the dataset.

3 Methods

3.1 Sentiment Analysis

To ensure that the products we recommend are those that users truly enjoy, we first had to narrow down the data used in our analysis and models. With that in mind, sentiment analysis was implemented as it is a Natural Language Processing (NLP) technique that evaluates the emotional tone of text, classifying it as positive, negative, or sometimes neutral. In the context of product reviews, sentiment analysis helps identify whether users had a positive or negative experience with a product, which is essential when making product recommendations.

Originally, we considered using the *is_recommended* column already present in the reviews dataset, to more easily determine which products users spoke well of. However, we found inconsistencies in this parameter, including a significant amount of missing values and 5-star reviews marked as "not recommended". This made the *is_recommended* column unreliable, making sentiment analysis a necessary alternative.

In our analysis, we focused on reviews rated 3 stars and above, excluding 1- and 2-star reviews as they were considered too negative. Since 4- and 5-star reviews are clearly positive, they were included directly in our study, while sentiment analysis was applied to only 3-star reviews, as their sentiments are ambiguous. To perform the analysis, we used two *Hugging Face* models: *xlnet-base-cased-product-review-sentiment-analysis* [4], specialized for product reviews, and *cardiffnlp/twitter-xlm-roberta-base-sentiment* [5], which supports multiple languages. This allowed us to capture both English and the few non-English reviews. While these models are somewhat of a black box, they are verified, widely used, and well-suited to our problem. To ensure accuracy, we intersected the results of both models, keeping only reviews marked as positive by both. This decreased false positives and therefore resulted in 14,454 positive reviews out of 60,703 3-star reviews for further analysis.

Note that, although there are challenges when taking the sentiment of product reviews such as ambiguous or softened feedback, the implementation of sentiment analysis still plays a crucial role in improving the quality of our recommendations.

3.2 Frequent Itemsets

The primary objective of using frequent itemsets is to uncover associations between products based on customer reviews. Through mining frequent itemsets and generating association rules, we aim to predict which products a customer might like based on their interest in another product. This insight allows us to build a recommendation system that enhances the shopping experience by suggesting items that other similar users have enjoyed. We chose to use the A-priori algorithm because its systematic approach to pruning the search space of itemsets ensures computational efficiency, making it suitable for our dataset size and structure.

Choosing the appropriate minimum support threshold and deciding how many reviews to exclude were critical steps in our analysis. Excluding more users resulted in a larger number of association rules, providing richer data for recommendations. However, this approach risks overfitting to a small subset of users, making the rules less generalizable to the broader population. Conversely, including too many users with fewer reviews requires significantly lowering the minimum support threshold in order to have enough association rules to work with, which can lead to less reliable rules. After exploring various configurations, as can be seen in Table 1, we determined that restricting the dataset to users with 10 or more reviews struck an

effective balance. This approach ensured that the generated rules supported meaningful and actionable recommendations.

A limitation of association rule mining is that not all products will have associated rules, especially those with low representation in the dataset. To address this, we incorporated a similar items approach. When a product lacks direct recommendations, we identify the most similar product that is present in the association rules. This ensures every product can act as an entry point for personalized recommendations. A more detailed accounting of this is provided later in the paper.

By mining frequent itemsets and generating association rules, we enable a data-driven recommendation system that suggests products based on user preferences. Furthermore, the integration of similar items enhances its usability across diverse customer journeys.

User Threshold	Association Rules Size	Unique Products Included	Comments
5+ reviews	6	3	Although we are including more users in our dataset, the unique products the resulting rules cover are nowhere near enough to make assumption about the rest of the dataset.
10+ reviews	1704	82	This seems to be a decent middle-ground, where enough unique products are included in the association rules, but we are not cutting out too much of the population.
30+ reviews	526,853	277	Here, we are approaching a super small subset of users. They account for only 0.001% of the total users, which is not enough to make generalizations.

Table 1: Finding a user threshold based on resulting association rules

3.3 Similar Items

Similar items can be found in different ways. Two of the most used measures are Cosine similarity and Jaccard similarity scores. Cosine similarity measures the distance between two vectors in a vector space, whereas Jaccard similarity measures the proportion of shared to the union of the items [1]. As mentioned briefly in our approach, the highlights and ingredients of each product will be used to determine a score. This score is obtained by calculating the similarity scores for highlights and ingredients individually, and then combining them using an approach known as Reciprocal Rank Fusion (RRF).

3.3.1 Similarity of Highlights

The product highlights are sets of the most emphasized features of the product. For example, a certain product targeted at users with dry skin would have a "dry skin" highlight to convey this attribute of the product. To find the similarity of preprocessed highlights between a random product and the input, we computed Jaccard Similarity. It measures how many features are in common within the sets of highlights. The more similar the sets, the more alike the products are in how they affect the skin and address skin issues.

3.3.2 Ingredient Similarity

Before calculating the similarity scores for ingredients, we apply a TF-IDF vectorizer to transform the ingredient lists into vectors. It prevents common ingredients, like water which is a common ingredient in almost all products, from disproportionately influencing similarity, focusing instead on more unique ingredients. The TF-IDF vectorizer accounts for both the frequency of an ingredient within a single product and its frequency across all products. [1]. The resulting vectors are then used to compute cosine similarity scores. The higher the similarity score, the more alike the ingredient lists are in terms of truly relevant ingredients.

3.3.3 Reciprocal Rank Fusion (RRF)

As the similarity scores are generated using different measures, their scales vary. Cosine similarity ranges from -1 to 1, whereas Jaccard similarity ranges from 0 to 1. Therefore, simply taking the average of the scores or any weighted sum would not be recommended.

To avoid this, we use an algorithm known as Reciprocal Rank Fusion (shown in Equation 1),

$$\text{RRF} = \sum_{i=1}^N \frac{1}{k + r_i(d)} \quad (1)$$

where $r_i(d)$ is the rank of product d in the i^{th} ranked list, k is a constant (typically $k = 60$) used to prevent division by zero and to adjust the impact of higher ranks, and N is the number of ranked lists (in our case 2) [6].

RRF uses the ranks of the product within the individual lists and then computes a singular score. The rank is simply the position of the product on a list sorted based on the similarity score in descending order. By doing so, we ensure that the higher the score, the more the product is similar to the input in terms of both highlights and ingredients. In the end, the resulting output of RRF is a ranked list of products.

4 Recommendation System

As mentioned briefly in our approach, the recommendation system operates with a logic similar to Collaborative Filtering, but is specifically designed to address the challenge of limited user history. Instead of relying on prior user data, it uses a single product as input, simulating scenarios where no purchase history or profile information is available. To recommend only products that other users have positively reviewed, we first filter reviews with ratings of 4 and 5 stars, as these are inherently positive. For 3-star reviews, we perform sentiment analysis, retaining only those with high sentiment scores, to ensure they reflect a positive experience. The A-priori algorithm is then applied to this filtered data to identify frequent itemsets, revealing which products are commonly bought together. This assumes that users with similar skin conditions tend to purchase and like related items. To handle cases where a product lacks any association rules — such as newly introduced items — the system calculates similarity scores against antecedents in all the found association rules. It identifies the most similar antecedent and utilizes its association rules to generate recommendations. This approach maintains consistency in the type of recommendations provided. Similarly, if there are insufficient association rules, the system uses the similar items method to find more products related to the consequents to ensure a consistent and diverse output.

5 Concluding remarks

One challenge we faced was in determining how many reviews to include in the analysis. Using a smaller subset of reviews often results in more association rules, which provides more granular information for recommendations. However, this approach risks bias, as it may not accurately reflect the broader population's preferences. To address this, two strategies were considered. By plotting user reviews (2+ review, 3+ reviews, etc.) against average confidence of the generated rules, the elbow point in the curve could help us identify a trade-off between rule confidence and dataset size. Another approach would be to use Algorithm 1 as described in the article "Mining association rules with multiple minimum supports" [7], where all the association rules are tested, and only the most influential are used to make assumptions about user preferences. Although we tried to incorporate these methods, we couldn't make it work with the time constraints we had, so we chose arbitrary values for the minimum threshold and user cutoffs. They were chosen based off of what would give decent results without giving over-fitted results. If given the opportunity to improve our recommendation system, we would like to implement these methods.

Another improvement that can be made is to increase the size of the dataset by scraping the data to include more products along with more information about the products themselves such as the product descriptions. The latter would have proved useful in performing similarity score analysis on as it has more textual information. An attempt was made to scrape data, but due to time constraints, we had to pivot.

A Appendix

A.1 Group Member's Contribution to the Project

	Ana Marija Pavičić	Anna Sky Kastl Jensen	Raquel Moleiro Marques	Sree Keerthi Desu
Motivation	40%	20%	0%	40%
Approach	35%	30%	0%	35%
Product Dataset	30%	0%	40%	30%
Review Dataset	20%	0%	60%	20%
Sentiment Analysis	0%	20%	80%	0%
Frequent Itemsets	0%	70%	30%	0%
Similar Items	50%	0%	0%	50%
Recommendation System	35%	20%	10%	35%
Concluding Remarks	15%	70%	0%	15%

Table 2: Contributions Table

A.2 Code

```

1  # -*- coding: utf-8 -*-
2  """final_notebook.ipynb
3
4  Automatically generated by Colab.
5
6  Original file is located at
7  https://colab.research.google.com/drive/1w5MPyzzZFuzAlRrybLQoq_WMXtN2hRW
8
9  # 02807 Computational Tools for Data Science – Final Project
10
11  Group members:
12  - Ana Marija Pavičić (s232468)
13  - Anna Sky Kastl Jensen (s194824)
14  - Raquel Moleiro Marques (s243636)
15  - Sree Keerthi Desu (s243933)
16
17  #### Imports
18  """
19
20  # Commented out IPython magic to ensure Python compatibility.
21  # %capture
22  # !pip install nbimporter
23  # from imports import *
24
25  """# 1. Data pre-processing"""
26
27  # Data paths
28  raw_data_directory = 'data'
29  raw_data_path = 'data/product_info.csv'
30  processed_data_directory = 'processed_data'
31  processed_data_path = processed_data_directory + '/skincare.csv'
32
33  """## Data Cleaning"""
34
35  # Read all products info
36  df = pd.read_csv(raw_data_path)
37
38  # Filter df for rows where 'primary_category' is 'Skincare'
39  skincare_df = df[df['primary_category'] == 'Skincare']
40  print("Skin care data size:", len(skincare_df))
41
42  # Remove rows where highlight are non existent
43  skincare_df = skincare_df[skincare_df['highlights'].notna() & (skincare_df['highlights'] !=
44  '')]
45  print("Skin care data size after removing empty highlights:", len(skincare_df))
46
47  # Remove rows where ingredients are non existent
48  skincare_df = skincare_df[skincare_df['ingredients'].notna() & (skincare_df['ingredients']
49  != '')]
50  print("Skin care data size after removing empty ingredients:", len(skincare_df))

```

```

50 # Remove 'Mini' size products from data as it could be seen as a duplicate
51 skincare_df = skincare_df[~skincare_df['product_name'].str.contains('mini', case=False, na=
52 False)]
53 # Also removing additional products with 'Mini Size' as their secondary category
54 skincare_df = skincare_df[skincare_df['secondary_category'] != 'Mini Size']
55 print(f"Size of skincare data after removing 'mini' from product names (case-insensitive): {
    len(skincare_df)}")
56 # Remove 'Limited edition' from data
57 skincare_df = skincare_df[~skincare_df['product_name'].str.contains('limited edition', case=
58 False, na=False)]
59 print(f"Size of skincare data after removing 'limited edition' from product names (case-
    insensitive): {len(skincare_df)}")
60 # Remove 'Value & Gift Sets' from data to focus on individual products
61 skincare_df = skincare_df[skincare_df['secondary_category'] != 'Value & Gift Sets']
62 print("Skin care data size after removing 'Value & Gift Sets' secondary_category:", len(
    skincare_df))
63 # Clean highlight column
64 highlights = skincare_df['highlights']
65 highlights = [h.replace("[", "").replace("]", "").replace("Best for ", "").
66               replace("Good for: ", "").replace(" Skin", "").replace("/", " ").replace(" for h in highlights)]
67 skincare_df['highlights'] = highlights
68 # Function to clean each ingredient row
69 def clean_ingredients(row):
70     # Replace unwanted characters
71     row = row.replace("[", "").replace("]", "").replace(" ", "").replace(" (Vegan)*", "").
72     replace(".", "")
73     # Remove text inside parentheses
74     row = re.sub(r'\([^)]*\)', '', row)
75     # Replace ", " with a single comma (in case extra spaces after commas)
76     row = row.replace(", ", ",")
77     row = row.replace(" ", "")
78     # Check for "water", "agua", or "eau" and replace first occurrence
79     if "water" in row.lower() or "agua" in row.lower() or "eau" in row.lower():
80         row_list = row.split(", ")
81         # Find the first occurrence of "water", "agua", or "eau"
82         index = next((i for i, s in enumerate(row_list) if 'water' in s.lower() or "agua" in s
83                     .lower() or "eau" in s.lower()), -1)
84         if index != -1:
85             # Replace the identified word with "Water"
86             row_list[index] = "Water"
87             row = ", ".join(row_list)
88     return row
89 # Apply the clean_ingredients function
90 skincare_df['ingredients'] = skincare_df['ingredients'].apply(clean_ingredients)
91 # Ensure the directory exists
92 output_dir = os.path.dirname(processed_data_path)
93 if not os.path.exists(output_dir):
94     os.makedirs(output_dir)
95 # Save cleaned file
96 skincare_df.to_csv(processed_data_path, index=False)
97 """## Clean and Filter Product Reviews"""
98 # Convert product_id column to a set for efficient filtering
99 product_ids = set(skincare_df['product_id'])
100 def clean_reviews(raw_data_path, processed_data_path):
101     """
102     Clean the reviews in the given CSV file by filtering based on the conditions.
103     - Filters rows where 'product_id' is in the skincare dataset.
104     - Keeps reviews with rating 5, 4, and only 3-star reviews if they have review text.
105     """
106     try:
107         # Read the data
108         df = pd.read_csv(raw_data_path)
109         # Keep only the filtered skincare data
110         df = df[df['product_id'].isin(product_ids)]
111         # Check for 3-star reviews with no text or empty text
112         three_star_reviews = df[df['rating'] == 3]
113         no_text_reviews = three_star_reviews[three_star_reviews['review_text'].isna() | (
114             three_star_reviews['review_text'].str.strip() == '')]

```



```

126     # Keep reviews with rating 5, 4, or 3 (only keep 3-star reviews if they have review
127     text)
128     df = df[
129         (df['rating'].isin((5, 4))) |
130         ((df['rating'] == 3) & df['review_text'].notna() & (df['review_text'].str.strip
131         () != "")) # Include 3-star reviews only if they have text
132     ]
133     # Ensure the output directory exists
134     output_dir = os.path.dirname(processed_data_path)
135     if not os.path.exists(output_dir):
136         os.makedirs(output_dir)
137     # Save the cleaned file
138     df.to_csv(processed_data_path, index=False)
139     print(f"Processed and saved: {processed_data_path}")
140
141 except Exception as e:
142     print(f"Error processing {raw_data_path}: {e}")
143
144 # Apply the function to each review data file
145 clean_reviews('data/reviews_0-250.csv', 'processed_data/reviews_0-250.csv')
146 clean_reviews('data/reviews_250-500.csv', 'processed_data/reviews_250-500.csv')
147 clean_reviews('data/reviews_500-750.csv', 'processed_data/reviews_500-750.csv')
148 clean_reviews('data/reviews_750-1250.csv', 'processed_data/reviews_750-1250.csv')
149 clean_reviews('data/reviews_1250-end.csv', 'processed_data/reviews_1250-end.csv')
150
151 print("All files have been processed and saved.")
152
153 """## Combining All Reviews Files Into a Single Dataframe"""
154 # List of file paths for the processed reviews
155 file_paths = [
156     'processed_data/reviews_0-250.csv',
157     'processed_data/reviews_250-500.csv',
158     'processed_data/reviews_500-750.csv',
159     'processed_data/reviews_750-1250.csv',
160     'processed_data/reviews_1250-end.csv'
161 ]
162
163 # Combine all files into one DataFrame
164 combined_reviews_df = pd.DataFrame()
165
166 for file_path in file_paths:
167     try:
168         # Read each processed file with updated argument for handling bad lines
169         df = pd.read_csv(file_path,
170                         on_bad_lines='skip', # Skip bad lines instead of using deprecated
171                         arguments
172                         encoding='utf-8', # Ensure proper encoding
173                         engine='python') # Use Python engine for more robust handling
174
175         # Append it to the combined DataFrame
176         combined_reviews_df = pd.concat([combined_reviews_df, df], ignore_index=True)
177         print(f"Successfully added {file_path}")
178     except Exception as e:
179         print(f"Error reading {file_path}: {e}")
180
181 # Save the combined DataFrame to a single CSV file
182 output_combined_path = 'processed_data/combined_reviews.csv'
183 combined_reviews_df.to_csv(output_combined_path, index=False)
184 print(f"All files have been combined and saved to: {output_combined_path}")
185
186 """## 2. Exploratory Data Analysis (EDA)
187
188 ### Skincare data
189
190 # Display the first few rows to get a sense of the data
191 print("Skincare Data Overview:")
192 print(skincare_df.head())
193
194 # Get general info
195 print("\nSkincare Data Info:")
196 print(skincare_df.info())
197
198 # Describe numerical columns for basic statistics
199 print("\nSkincare Data Statistics:")
200 print(skincare_df.describe())
201
202 """### Product Reviews Data"""
203
204 # Display the first few rows to get a sense of the data
205 print("\nReview Data Overview:")
206

```

```

208 print(combined_reviews_df.head())
209
210 # Get general info
211 print("\nReview Data Info:")
212 print(combined_reviews_df.info())
213
214 # Describe numerical columns for basic statistics
215 print("\nReview Data Statistics:")
216 print(combined_reviews_df.describe())
217
218 """### Distribution of categorical and numerical features
219
220 ### Skincare data
221 """
222
223 # Check distribution of secondary categories
224 print("\nSkincare Secondary Categories Distribution:")
225 secondary_category_counts = skincare_df['secondary_category'].value_counts().sort_values(
    ascending=False)
226 print(secondary_category_counts)
227
228 # Plot the distribution of secondary categories
229 plt.figure(figsize=(10,6))
230 sns.countplot(
231     y='secondary_category',
232     data=skincare_df,
233     order=secondary_category_counts.index # Sort the plot by frequency
234 )
235 plt.title('Distribution of Secondary Categories')
236 plt.show()
237
238 # Get the count and percentage of top brands
239 brand_count = skincare_df['brand_name'].value_counts()
240 brand_percentage = skincare_df['brand_name'].value_counts(normalize=True) * 100
241
242 # Create a DataFrame for the top brands
243 brand_df = pd.DataFrame({'Brand': brand_count.index, 'Counts': brand_count.values, 'Percent'
    : np.round(brand_percentage.values, 2)})
244
245 # Display the top 10 brands
246 print(brand_df.head(10))
247
248 # Plot the top 10 brands horizontally
249 plt.figure(figsize=(8, 6))
250 sns.barplot(x='Counts', y='Brand', data=brand_df.head(10))
251 plt.xlabel('Count')
252 plt.title('Top 10 Occurring Brands')
253 plt.show()
254
255 """#### Product Reviews Data"""
256
257 # Review rating distribution
258 print("\nReview Ratings Distribution:")
259 print(combined_reviews_df['rating'].value_counts())
260
261 # Plot the rating distribution
262 plt.figure(figsize=(8,6))
263 sns.countplot(x='rating', data=combined_reviews_df)
264 plt.title('Distribution of Review Ratings')
265 plt.show()
266
267 # Review length (characters in review text)
268 combined_reviews_df['review_length'] = combined_reviews_df['review_text'].apply(lambda x:
    len(str(x)))
269
270 # Plot distribution of review lengths
271 plt.figure(figsize=(10,6))
272 sns.histplot(combined_reviews_df['review_length'], bins=50, kde=True)
273 plt.title('Distribution of Review Lengths')
274 plt.show()
275
276 """
277
278 ## Filtering 3-Star Reviews
279 """
280
281 # Load the combined reviews data from the saved file
282 combined_reviews_df = pd.read_csv('processed_data/combined_reviews.csv', encoding='utf-8',
    engine='python')
283
284 # Filter for reviews with rating = 3 to perform sentiment analysis on the same
285 three_stars = combined_reviews_df[combined_reviews_df['rating'] == 3]
286
287 print(f"Total reviews with rating of 3: {len(three_stars)}")
288
289 """## *is_recommended* column in the dataset

```

```

290 To ensure we would recommend products well-regarded by users, we explored the *
291 is_recommended* column in our reviews dataset and performed a series of checks to
    evaluate its reliability. As per the following:
292 """
293 # Count the number of 1s, 0s, and NAs in the 'is_recommended' column for 3-5 stars reviews
294 is_recommended_counts = combined_reviews_df['is_recommended'].value_counts(dropna=False)
295 is_recommended_counts
296
297 """We can already tell that there are a lot of missing values in the *is_recommended* column
    . Either way, before proceeding, we'll check if any 5-star reviews are incorrectly
    marked as not recommended, as we would expect these to generally be recommended. This
    will help ensure data quality."""
298
299 # Filter 5-star reviews
300 five_star_reviews = combined_reviews_df[combined_reviews_df['rating'] == 5]
301
302 # Check how many of the 5-star reviews are not recommended
303 not_recommended_5_stars = five_star_reviews[five_star_reviews['is_recommended'] == 0]
304
305 # Count the number of 5-star reviews that are not recommended
306 not_recommended_count = not_recommended_5_stars.shape[0]
307
308 print(f"Number of 5-star reviews marked as not recommended: {not_recommended_count}")
309
310 # Show a random review_text from the not recommended 5-star reviews, random seed set for
    reproducibility
311 random_example = not_recommended_5_stars['review_text'].sample(1, random_state=42).iloc[0]
312
313 print("\nRandom example of 5-star review marked as not recommended:")
314 print(random_example)
315
316 """Since there are 610 5-star reviews marked as "not recommended," and the random review
    test shows a very positive review from a happy user, this *is_recommended* column can't
    be trusted. Therefore, we'll rely on sentiment analysis to accurately identify products
    users truly recommend and feel good about.
317
318 # 3. Sentiment Analysis on Product Reviews
319
320 ## xlnet-base-cased-product-review-sentiment-analysis
321
322 Below we use [xlnet-base-cased-product-review-sentiment-analysis](https://huggingface.co/
    dipawidia/xlnet-base-cased-product-review-sentiment-analysis) model from Hugging face
    to perform sentiment analysis.
323 """
324
325 # Initialize the tokenizer and model
326 tokenizer = XLNetTokenizer.from_pretrained("dipawidia/xlnet-base-cased-product-review-
    sentiment-analysis")
327
328 model = TFXLNetForSequenceClassification.from_pretrained("dipawidia/xlnet-base-cased-product-
    review-sentiment-analysis")
329
330 # Function to analyze sentiment for a batch of reviews
331 def get_sentiment_batch(reviews_batch):
332     tokenize_text = tokenizer(reviews_batch.tolist(), padding=True, truncation=True,
        return_tensors='tf', max_length=512)
333     preds = model.predict(dict(tokenize_text))['logits']
334     class_preds = np.argmax(tf.keras.layers.Softmax()(preds), axis=-1)
335     labels = ['Positive' if pred == 1 else 'Negative' for pred in class_preds]
336     return labels
337
338 # Batch size for processing reviews
339 batch_size = 64
340
341 # Initialize the progress bar for batching
342 num_batches = len(three_stars) // batch_size + (1 if len(three_stars) % batch_size != 0 else
    0)
343
344 # List to hold sentiment results
345 sentiment_results = []
346
347 # Process reviews in batches
348 for i in tqdm(range(0, len(three_stars), batch_size), total=num_batches, desc="Processing
    reviews", unit="batch"):
349     batch = three_stars['review_text'][i:i + batch_size] # Get the current batch of reviews
350     batch_sentiments = get_sentiment_batch(batch) # Get sentiments for the current batch
351     sentiment_results.extend(batch_sentiments) # Append results
352
353 # Assign the sentiment results back to the dataframe
354 three_stars['sentiment'] = sentiment_results
355
356 # Count positive and negative sentiments
357 sentiment_counts = three_stars['sentiment'].value_counts()
358 print(sentiment_counts)
359
360 # Filter the reviews with positive sentiment

```

```

361 positive_reviews_xlnet = three_stars[three_stars['sentiment'] == 'Positive']
362
363 # Display the count of positive sentiment reviews
364 print(f"Number of positive sentiment reviews: {len(positive_reviews_xlnet)}")
365
366 print("Positive Sentiment Examples:")
367 print(positive_reviews_xlnet[['review_text']].head())
368
369 """
370 ----
371 ## cardiffnlp/twitter-xlm-roberta-base-sentiment
372
373 Below we use [cardiffnlp/twitter-xlm-roberta-base-sentiment](https://huggingface.co/
cardiffnlp/twitter-xlm-roberta-base-sentiment) model from Hugging face to perform again
sentiment analysis.
"""
374
375 # Model path for the new model
376 model_path = "cardiffnlp/twitter-xlm-roberta-base-sentiment"
377
378 # Initialize the tokenizer and model
379 tokenizer = AutoTokenizer.from_pretrained(model_path)
380 model = TFAutoModelForSequenceClassification.from_pretrained(model_path)
381
382 # Function to analyze sentiment for a batch of reviews
383 def get_sentiment_batch(reviews_batch):
384     # Tokenize the batch of reviews
385     tokenize_text = tokenizer(reviews_batch.tolist(), padding=True, truncation=True,
386                             return_tensors='tf', max_length=512)
387
388     # Get predictions from the model
389     preds = model(**tokenize_text)['logits']
390
391     # Apply softmax to get class probabilities
392     scores = softmax(preds.numpy(), axis=-1)
393
394     # Get the class with the highest probability (0 = negative, 2 = positive)
395     class_preds = np.argmax(scores, axis=-1)
396
397     # Assign sentiment labels based on the class prediction
398     labels = ['Negative' if pred == 0 else 'Neutral' if pred == 1 else 'Positive' for pred
399              in class_preds]
400
401     return labels
402
403 # Batch size for processing reviews
404 batch_size = 64
405
406 # Ensure the reviews are valid strings and filter out NaN values
407 three_stars['review_text'] = three_stars['review_text'].fillna('').astype(str)
408
409 # Initialize the progress bar for batching
410 num_batches = len(three_stars) // batch_size + (1 if len(three_stars) % batch_size != 0 else
411 0)
412
413 # List to hold sentiment results
414 sentiment_results = []
415
416 # Process reviews in batches
417 for i in tqdm(range(0, len(three_stars), batch_size), total=num_batches, desc="Processing
418 reviews", unit="batch"):
419     batch = three_stars['review_text'][i:i + batch_size] # Get the current batch of reviews
420     batch_sentiments = get_sentiment_batch(batch) # Get sentiments for the current batch
421     sentiment_results.extend(batch_sentiments) # Append results
422
423 # Assign the sentiment results back to the dataframe
424 three_stars['sentiment'] = sentiment_results
425
426 # Count positive, negative, and neutral sentiments
427 sentiment_counts = three_stars['sentiment'].value_counts()
428 print(sentiment_counts)
429
430 # Filter the positive sentiment reviews
431 positive_reviews_cardiffnlp = three_stars[three_stars['sentiment'] == 'Positive']
432
433 # Display the count of positive sentiment reviews
434 print(f"Number of positive sentiment reviews: {len(positive_reviews_cardiffnlp)}")
435
436 print("Positive Sentiment Examples:")
437 print(positive_reviews_cardiffnlp[['review_text', 'product_id']].head())
438
439 """
440 ----
441 ## Intersect all positive reviews found
442
443 # Merge the two positive review datasets based on the 'Unnamed: 0' column
444 intersected_positive_reviews = pd.merge(

```

```

443     positive_reviews_xlnet[['Unnamed: 0', 'review_text', 'product_id', 'author_id', '
444         review_title',
445         'product_name', 'brand_name', 'price_usd']],
446     positive_reviews_cardiffnlp[['Unnamed: 0', 'product_id', 'author_id', 'review_title',
447         'product_name', 'brand_name', 'price_usd']],
448     on='Unnamed: 0'
449     how='inner' # 'inner' ensures only matching rows are returned
450 )
451 # Rename columns to avoid conflicts (the '_x' and '_y' suffixes)
452 intersected_positive_reviews.rename(columns={
453     'review_text_x': 'review_text',
454     'product_id_x': 'product_id',
455     'author_id_x': 'author_id',
456     'review_title_x': 'review_title',
457     'product_name_x': 'product_name',
458     'brand_name_x': 'brand_name',
459     'price_usd_x': 'price_usd'
460 }, inplace=True)
461
462 # Now drop the columns from the second dataframe that are not needed (with the '_y' suffix)
463 intersected_positive_reviews.drop(columns=[col for col in intersected_positive_reviews.
464     columns if col.endswith('_y')], inplace=True)
465
466 # Select only desired columns
467 final_columns = ['Unnamed: 0', 'author_id', 'review_text', 'review_title', 'product_id', '
468     product_name', 'brand_name', 'price_usd']
469 intersected_positive_reviews = intersected_positive_reviews[final_columns]
470
471 # Display the count of intersected positive sentiment reviews
472 print(f"Number of intersected positive sentiment reviews: {len(intersected_positive_reviews)}")
473
474 # Check the first few rows of the final dataframe
475 print(intersected_positive_reviews.head())
476
477 # Save the intersected positive reviews to a CSV file
478 intersected_positive_reviews.to_csv('processed_data/3_star_positive_reviews.csv', index=
479     False)
480
481 """
482 -----
483 """
484 ## Filtering and Combining Desired Reviews
485
486 # Load the combined reviews data from the saved file
487 combined_reviews_df = pd.read_csv('processed_data/combined_reviews.csv', encoding='utf-8',
488     engine='python')
489
490 # Filter the reviews with 4 and 5 stars
491 frequent_items_reviews = combined_reviews_df[combined_reviews_df['rating'].isin([4, 5])]
492
493 # Load the positive reviews from the sentiment analysis on the 3-star reviews
494 positive_reviews_df = pd.read_csv('processed_data/3_star_positive_reviews.csv', encoding='
495     utf-8', engine='python')
496
497 # Combine the filtered 4 and 5-star reviews with the 3-star positive reviews
498 frequent_items_reviews = pd.concat([frequent_items_reviews, positive_reviews_df],
499     ignore_index=True)
500
501 print(f"Total reviews to be used in frequent items: {len(frequent_items_reviews)}")
502
503 # Check how many unique products there are based on 'product_id'
504 unique_products_count = frequent_items_reviews['product_id'].nunique()
505
506 # Pretty printing the result
507 print(f"Total number of unique products in the reviews dataset: {unique_products_count}")
508
509 """# 4. Frequent Items"""
510
511 # Step 1: Aggregate reviews by user
512 user_review_counts = frequent_items_reviews.groupby('author_id')['product_id'].nunique()
513
514 # Step 2: Define thresholds for grouping
515 thresholds = list(range(2, 11)) + [20, 30] + list(range(40, user_review_counts.max() + 10,
516     10))
517
518 # Step 3: Count users per group
519 users_per_group = []
520 for t in thresholds:
521     count = (user_review_counts >= t).sum()
522     users_per_group.append(count)
523
524 # Step 4: Display results
525 print("Users per group:")

```

```

520 for t, count in zip(thresholds, users_per_group):
521     print(f"{t}+ reviews: {count} users")
522
523 plt.figure(figsize=(12, 6))
524 plt.bar([str(t) + "+" for t in thresholds], users_per_group, color="skyblue", alpha=0.8)
525 plt.title("Number of Users per Group Based on Minimum Reviews")
526 plt.xlabel("Minimum Reviews per User")
527 plt.ylabel("Number of Users")
528 plt.xticks(rotation=45)
529 plt.grid(axis='y', linestyle='—', alpha=0.7)
530 plt.show()
531
532 """## Threshold on users with 10+ reviews"""
533
534 # Selecting users with 10+ reviews
535 filtered_df = frequent_items_reviews.groupby('author_id').filter(lambda x: x['product_id'].
    nunique() > 9)
536
537 # Display the filtered DataFrame
538 print(f"Number of rows in the filtered dataframe: {len(filtered_df)}")
539 print(f"Number of unique authors in the filtered dataframe: {len(set(filtered_df['author_id']
    '))}")
540 #set(filtered_df['author_id'])
541
542 # Combining all the products reviewed for each person in the dataset
543 selected_columns = filtered_df[['author_id', 'product_id']]
544
545 # Convert product_id to string before applying 'join'
546 combined_reviews = selected_columns.groupby('author_id')['product_id'].apply(lambda x: ' '.
    join(x.astype(str))).reset_index()
547
548 print(combined_reviews)
549
550 # Convert the 'reviews' into a list of transactions
551 transactions = combined_reviews['product_id'].str.split().tolist()
552
553 # Create a DataFrame for one-hot encoding
554 # Flatten all unique items (reviews) and create a unique item list
555 te = TransactionEncoder()
556 te_array = te.fit(transactions).transform(transactions)
557 df_encoded = pd.DataFrame(te_array, columns=te.columns_)
558
559 # Apply the Apriori algorithm
560 frequent_itemsets = apriori(df_encoded, min_support=0.03, use_colnames=True)
561
562 # Focus on frequent itemsets containing only a single product (not pairs or larger sets)
563 frequent_pairs = frequent_itemsets[frequent_itemsets['itemsets'].apply(len) == 1]
564
565 # Generate association rules, with support as minimum of 3%
566 rules = association_rules(frequent_itemsets, metric="support", min_threshold=0.03,
    num_itemsets=len(frequent_itemsets))
567
568 print(frequent_pairs)
569
570 # Filter rules to include only those where antecedents have a single item
571 rules_filtered = rules[rules['antecedents'].apply(len) == 1].copy()
572
573 rules_display = rules_filtered[["antecedents", "consequents", "support", "confidence", "lift
    "]].copy()
574
575 # Convert frozensets to readable strings
576 rules_display["antecedents"] = rules_display["antecedents"].apply(lambda x: ', '.join(list(x
    )))
577 rules_display["consequents"] = rules_display["consequents"].apply(lambda x: ', '.join(list(x
    )))
578
579 # Sorting by confidence
580 rules_display = rules_display.sort_values(by="support", ascending=False)
581 print(rules_display.to_string(index=False))
582
583 """### Testing"""
584
585 filtered_rules = rules_display[rules_display['antecedents'].apply(lambda x: 'P500633' in x)]
586 print(filtered_rules)
587
588 # Preprocess the rules DataFrame to ensure 'antecedents' and 'consequents' are sets
589 def preprocess_rules(rules):
590     """
591     Preprocess the rules DataFrame to ensure 'antecedents' and 'consequents' are sets.
592     """
593     # Create a copy of the DataFrame to avoid modifying the original
594     rules = rules.copy()
595     rules['antecedents'] = rules['antecedents'].apply(lambda x: x.split(", ") if isinstance(x
        , str) else x)

```



```

596 rules['consequents'] = rules['consequents'].apply(lambda x: x.split(", ") if isinstance(x,
597 str) else x)
598 return rules
599 # Function to find items associated with an input item
600 def find_associated_items(input_item, rules):
601     """
602     Find all items associated with an input item based on association rules.
603
604     Parameters:
605         input_item (str): The item to find associations for.
606         rules (pd.DataFrame): A DataFrame of association rules with columns 'antecedents' and
        'consequents'.
607
608     Returns:
609         list: A list of items associated with the input item, preserving the order of rules.
610     """
611     # Sort rules by confidence in descending order
612     rules = rules.sort_values(by='confidence', ascending=False).reset_index(drop=True)
613
614     input_item_antecedents = [input_item]
615     associated_items = []
616
617     # Iterate through the rules
618     for _, rule in rules.iterrows():
619         antecedents = rule['antecedents']
620         consequents = rule['consequents']
621
622         # Check if the input item is in the antecedents
623         if input_item_antecedents == antecedents:
624             # Add consequents to the associated items if not already present
625             for item in consequents:
626                 if item not in associated_items:
627                     associated_items.append(item)
628
629     return associated_items
630
631 rules_df = pd.DataFrame(rules_display)
632 # Preprocess the rules DataFrame
633 rules_df = preprocess_rules(rules_df)
634
635 # Input item
636 input_item = 'P500633'
637
638 # Find associated items
639 associated_items = find_associated_items(input_item, rules_df)
640 print(f"Items associated with '{input_item}': {associated_items}")
641
642 rules_display.to_csv('processed_data/association_rules.csv', index=False)
643
644 """
645 -----
646 # 5. Similar Items Based on Ingredients and Highlights
647
648 ### Calculate similar items using ingredients
649
650 To calculate similar items using ingredients we choose to embed ingredients with TF-IDF.
651 This was chosen because TF-IDF will take into account the uniqueness of ingredients. For
652 example water is a frequent ingredient in the list which does not tell us a lot about
653 the product. Where as salicylic acid does not occur as often leading to higher
654 importance. To compute the similarity between products we choose cosine similarity
655 because TF-IDF gives results in the form of embedding.
656
657 #### Display ingredients
658 """
659 df = pd.read_csv("processed_data/skincare.csv")
660 # Select only 'Name' and 'Ingredients' columns
661 df_selected = df[['product_name', 'ingredients']]
662
663 # Display the first 5 rows
664 print(df_selected.head())
665
666 """#### Define functions"""
667 def custom_tokenizer(text):
668     return text.split(", ")
669
670 def similarity_search_ingredients(df, query):
671     """
672     Perform a similarity search based on cosine similarity of TF-IDF vectors.
673
674     Parameters:
675         query (str): The input query string.
676         df (pd.DataFrame): A DataFrame containing 'id' and 'ingredients' columns.
677     """

```

```

677 Returns:
678 _results_df (pd.DataFrame): Rows from the original DataFrame sorted by similarity.
679 """
680 # Initialize TfidfVectorizer with a custom tokenizer (adjust lowercase as needed)
681 vectorizer = TfidfVectorizer(tokenizer=custom_tokenizer, lowercase=False)
682
683 # Extract the 'ingredients' column
684 ingredients = df['ingredients']
685
686 # Fit and transform the ingredients
687 tfidf_matrix = vectorizer.fit_transform(ingredients)
688
689 # Transform the query into the TF-IDF space
690 query_tfidf = vectorizer.transform([query])
691
692 # Compute cosine similarity between the query and all documents
693 similarities = cosine_similarity(query_tfidf, tfidf_matrix).flatten()
694
695 # Add similarity scores to the DataFrame
696 df['similarity_score_ingredients'] = similarities
697
698 # Sort the DataFrame by similarity scores in descending order
699 results_df = df.sort_values(by='similarity_score_ingredients', ascending=False).
700     reset_index(drop=True)
701
702 results_df['rank_ingredients'] = range(1, len(results_df) + 1)
703
704 return results_df
705
706 """### Calculate similarity using the highlights
707 In the dataset we choose we were not provided with description but rather highlights of a
708 product. Those are sets of words that describe the product, so therefore there was no
709 need to use the minhashing as the length of sets is not big to begging with. To compare
710 the sets we choose to use Jaccard Similarity because the more highlight product have in
711 common the more similar they are.
712
713 """
714
715 ##### Display a few rows of highlights
716 """
717 df_selected = df[['product_name', 'highlights']]
718
719 # Display the first 5 rows
720 print(df_selected.head())
721
722 """##### Define functions"""
723
724 def jaccard_similarity(set_a, set_b):
725     """
726     Calculate the Jaccard Similarity between two sets.
727     """
728     intersection = len(set_a.intersection(set_b))
729     union = len(set_a.union(set_b))
730     return intersection / union if union != 0 else 0.0
731
732 def similarity_search_highlights(df, token_list):
733     """
734     Perform similarity search based on Jaccard similarity between df and a token list.
735     :param df: A pandas DataFrame with columns ['product_id', 'product_name', 'highlights'].
736     :param token_list: A list of tokens to compare against (highlights).
737     :return: A DataFrame with IDs and their Jaccard similarity scores, sorted by similarity
738             score.
739     """
740     # Convert the token list to a set
741     token_set = set(token_list)
742
743     # List to store the similarity scores
744     similarity_scores = []
745
746     # Iterate over the rows of the DataFrame
747     for index, row in df.iterrows():
748         product_id = row['product_id']
749         title = row['product_name']
750         # Convert the tokens for this ID to a set
751         id_token_set = set(row['highlights'].split(" "))
752
753         # Calculate Jaccard similarity
754         similarity_score = jaccard_similarity(id_token_set, token_set)
755
756         # Append the result as a tuple (id, score)
757         similarity_scores.append((product_id, title, similarity_score, (row['highlights'])))
758
759     # Convert the list of similarity scores to a DataFrame
760     similarity_df = pd.DataFrame(similarity_scores, columns=['product_id', 'product_name', '
761         similarity_score_highlights', 'highlights'])
762
763     # Sort the DataFrame by the 'similarity_score' column in descending order

```



```

758 similarity_df_sorted = similarity_df.sort_values(by='similarity_score_highlights',
759         ascending=False).reset_index(drop=True)
760 similarity_df_sorted['rank_highlights'] = range(1, len(similarity_df_sorted) + 1)
761
762 return similarity_df_sorted
763
764 """### Combining similarities based on highlight and ingredients
765 Jaccard similarity measure is in range [0, 1] and cosine similarity is in range [-1, 1]
    which means that we can not simply calculate the average or weighted sum. Due to that
    reason we choose to use reciprocal rank fusion algorithm. To use it we need to rank our
    similarities tables which is already done in their respective functions.
766 The **Reciprocal Rank Fusion (RRF)** score for a product _d_ is calculated as:
767 
$$RRF(d) = \sum_{i=1}^N \frac{1}{r_i(d) + k}$$

768 \text{}
769 Where:
770 -  $r_i(d)$  is the rank of product sd in the  $i^{th}$  ranked list.
771 -  $k$  is a constant (typically  $k = 60$ ), used to prevent division by zero and to adjust the
772   impact of higher ranks.
773 -  $N$  is the number of ranked lists (models or sources).
774 The reciprocal rank is typically defined as  $\frac{1}{r_i(d)}$ , where higher ranks give more
775   weight to the document.
776 """
777
778 def reciprocal_rank_fusion(df_highlights, df_ingredients, k=60):
779     """
780     Compute Reciprocal Rank Fusion (RRF) scores based on rank_highlights and rank_ingredients
781
782     Parameters:
783     - df_highlights (pd.DataFrame): DataFrame containing 'product_id', 'rank_highlights', and
784       other relevant columns.
785     - df_ingredients (pd.DataFrame): DataFrame containing 'product_id', 'rank_ingredients',
786       and other relevant columns.
787     - k (int): A constant for RRF computation (default=60).
788
789     Returns:
790     - combined_df (pd.DataFrame): A new DataFrame with overall RRF scores and combined
791       ranking.
792     """
793     # Merge the two DataFrames on 'product_id'
794     merged_df = pd.merge(
795         df_highlights, # Include all columns from df_highlights
796         df_ingredients[['product_id', 'rank_ingredients']], # Include only product_id and
797         rank_ingredients
798         on='product_id',
799         how='inner',
800     )
801     # Fill missing ranks with a large value (e.g., very low relevance)
802     merged_df['rank_highlights'] = merged_df['rank_highlights'].fillna(float('inf'))
803     merged_df['rank_ingredients'] = merged_df['rank_ingredients'].fillna(float('inf'))
804
805     # Compute the RRF score
806     merged_df['rrf_score'] = (
807         1 / (k + merged_df['rank_highlights']) +
808         1 / (k + merged_df['rank_ingredients'])
809     )
810
811     # Sort by the RRF score in descending order
812     merged_df = merged_df.sort_values(by='rrf_score', ascending=False).reset_index(drop=True)
813
814     # Add a new rank based on the RRF score
815     merged_df['overall_rank'] = range(1, len(merged_df) + 1)
816
817     return merged_df
818
819 """### Run code for specific product_id"""
820
821 def get_similar_items(product_id, df, n = 5):
822     """
823     Retrieve the top N products most similar to a given product based on highlights and
824     ingredients.
825
826     This function takes a product ID, performs similarity searches on the product's
827     highlights and ingredients,
828     and combines the results using a reciprocal rank fusion algorithm. It returns the top N
829     most similar products.
830
831     Parameters:
832     product_id : int or str
833         The ID of the product for which similar items are being searched.
834     """

```

```

830     n : int, optional
831         The number of similar products to return. Default is 5.
832
833     Returns:
834     ---merged_results (pd.DataFrame): A new DataFrame containing top n similar items
835     """
836
837     # get the selected product
838     product = df[df['product_id'] == product_id]
839
840     # get the product highlights and ingredients
841     product_highlights = list(product['highlights'])[0].split(", ")
842     product_ingredients = str(product['ingredients'])
843
844     # remove the product I am searching for
845     df = df[(df['product_id'] != product_id)]
846
847     # perform similarity searches
848     highlights_similarity_results = similarity_search_highlights(df, product_highlights)
849     ingredients_similarity_results = similarity_search_ingredients(df, product_ingredients)
850
851     # combine similarity searches with reciprocal rank fusion algorithm
852     merged_results = reciprocal_rank_fusion(highlights_similarity_results,
853                                             ingredients_similarity_results)
854
855     # Return only the top-n products
856     return merged_results[:n]
857
858 df = pd.read_csv("processed_data/skincare.csv")
859 product_id = "P442001"
860 product_name = df[df['product_id'] == product_id]['product_name'].iloc[0]
861 merged = get_similar_items(product_id, df)['product_name'].head(5)
862 print(f"The most similar products to the {product_name} are: ")
863 print(merged)
864
865 """
866 -----
867 # 6.Recommender
868
869 To handle cases where a product lacks association rules such as newly introduced
870 items the system calculates
871 similarity scores against all other products in the database. It identifies the most similar
872 product and
873 utilizes its association rules to generate recommendations. This approach maintains
874 consistency in the type
875 of recommendations provided. Similarly, if there are insufficient association rules, the
876 system leverages the
877 "similar items" method by using related products as a basis to ensure a consistent and
878 diverse output.
879 """
880
881 # Recommender function
882 def recommender(product_id, rules, df, n=5):
883     """
884     Recommend products based on association rules and similarity.
885
886     Parameters:
887     product_id (str): The product ID for which recommendations are needed.
888     rules (pd.DataFrame): A DataFrame of association rules with columns 'antecedents'
889         and 'consequents'.
890     n (int): Number of recommendations to return.
891
892     Returns:
893     --- list: Recommended products.
894     """
895
896     # Step 1: Find associated items
897     associated_items = find_associated_items(product_id, rules)
898
899     # no rule for chosen item
900     if len(associated_items) == 0:
901         # find the most similar antecedents to the product_id
902
903         # get all the single antecedents in the rules
904         antecedents = [
905             list(item)[0]
906             for item in rules["antecedents"]
907             if len(item) == 1
908         ]
909
910         # filter skincare so it searches for the similar items only among the antecedents
911         # and product_id
912         antecedents.append(product_id)
913         df_antecedents = df[df['product_id'].isin(antecedents)]

```

```

908     # the most similar antecedents to a product_id
909     the_most_similar_product_id = get_similar_items(product_id, df_antecedents, n = 1)['
        product_id'].iloc[0]
910
911     associated_items = find_associated_items(the_most_similar_product_id, rules)
912
913     # get the top n associated items
914     if len(associated_items) >= n:
915         return associated_items[:n]
916
917     # if there are not enough associated items get the similar items to the product_id as
    well
918     else:
919         number_of_similar_items = n - len(associated_items)
920
921         # for each of the consequents find similar items and combine them in one dataframe
922         similar_items_dataframes = []
923         for i in associated_items:
924             # get the top number_of_similar_items
925             similar_items = get_similar_items(product_id, df)
926             similar_items_dataframes.append(similar_items)
927
928         all_similar_items = pd.concat(similar_items_dataframes, axis=0)
929
930         # order is by similarity_score_highlights
931         all_similar_items = all_similar_items.sort_values(by='rrf_score', ascending=False).
            reset_index(drop=True)
932
933         # take only the top number_of_similar_items
934         similar_products = list(all_similar_items['product_id'])[:number_of_similar_items]
935
936         return associated_items + similar_products
937
938     """### Load data"""
939
940     association_rules = pd.read_csv('processed_data/association_rules.csv')
941     association_rules = preprocess_rules(association_rules)
942     skincare_df = pd.read_csv("processed_data/skincare.csv")
943
944     """### Example of usage"""
945
946     random_product_id = skincare_df['product_id'].sample(n=1).iloc[0]
947
948     product_name = skincare_df[skincare_df['product_id'] == random_product_id]['product_name'].
        iloc[0]
949     merged = get_similar_items(random_product_id, skincare_df)['product_name'].head(5)
950
951     print(f"The most similar products to the {product_name} are: ")
952     print(merged)

```

References

- [1] J. Leskovec, A. Rajaraman, and J. D. Ullman, *Mining of Massive Datasets*. USA: Cambridge University Press, 2nd ed., 2014.
- [2] Kaggle, “Sephora products and skincare reviews,” 2023. <https://www.kaggle.com/datasets/nadyinky/sephora-products-and-skincare-reviews>.
- [3] “Github repository for the project.” <https://github.com/raquelmdtum/CTDS-Final-Project>.
- [4] H. Face, “dipawidia/xlnet-base-cased-product-review-sentiment-analysis.” <https://huggingface.co/dipawidia/xlnet-base-cased-product-review-sentiment-analysis>.
- [5] H. Face, “twitter-xlm-roberta-base for sentiment analysis.” <https://huggingface.co/cardiffnlp/twitter-xlm-roberta-base-sentiment>.
- [6] D. Shah, “Reciprocal rank fusion (rrf) explained in 4 mins.,” 2024.
- [7] B. Liu, W. Hsu, and Y. Ma, “Mining association rules with multiple minimum supports,” in *Proceedings of the Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’99, (New York, NY, USA), p. 337–341, Association for Computing Machinery, 1999.