

Module: ITNPBD2 Autumn 2025

Assessment: Main assignment

Due Date/Time: 14/11/2025

AIAS Levels Allowed: 2

	Please tick the boxes/include appropriate information below
Student ID Number	3461226
Word Count (penalties apply for exceeding the stated limit)	957
I have read and understand the severity of academic misconduct – see link below	<input checked="" type="checkbox"/>
I give consent for my work to be used as an exemplar to future students.	<input checked="" type="checkbox"/>
I have checked my submitted document to ensure it complies with module requirements.	<input checked="" type="checkbox"/>
Link to version-controlled file (i.e on OneDrive, Google Docs, Github, or other) which contain evidence of the process I undertook to complete this assignment. Information on how to create a Microsoft 365 OneDrive folder is available HERE . *Please see notes below	https://github.com/kristina-pal/3461226_BD2
I understand that if there is a concern about potential academic misconduct, including inappropriate use of AI tools, then I could be asked to provide evidence of my drafting process during an academic integrity meeting if I have not done so using the link above. Not providing evidence of my drafting process could prejudice the outcome of academic misconduct cases.	<input checked="" type="checkbox"/>
Tailored feedback. If you would like tailored feedback on a specific aspect (or aspects) of your work (e.g., referencing, writing style, grammar), then please give details here.	I would appreciate tailored feedback on my data analysis approach and the development of more creative and effective problem-solving strategies to improve future data-driven work.
If you used AI at (or below) the level allowed, please explain briefly which AI, how you used it, and for what purpose.	I used ChatGPT at level 2 to check grammar in the report, and to assist with debugging minor Python syntax issues in Jupyter Notebook. All analytical work and coding decisions were my own.

**This may include (but is not limited to) drafts, versions of the finished document, notes, references, AI output, and AI prompts. These materials are not marked or graded, but they are simply a way to demonstrate how your work was created and to confirm that any AI use in your final submission is within the permitted AIAS scale for your assessment. Providing this helps safeguard you, showing your authentic process, and protecting you should any academic integrity questions arise.*

<https://www.stir.ac.uk/about/professional-services/student-academic-and-corporate-services/academic-registry/policy-and-procedure/>

JC Penney Data Analysis Report

Overview

This report provides an analytical overview of data describing products sold by JC Penney, along with customer demographics and reviews.

The analysis was conducted in **Python**, applying data-processing and machine-learning techniques to identify key factors that influence product popularity and customer satisfaction.

The findings indicate that the company serves a wide and diverse customer base, with no dominant age group. However, the most popular product categories across age groups tend to have relatively low customer ratings, suggesting issues with product quality.

Conversely, products with high ratings often exhibit low popularity, indicating limited promotional visibility.

An evaluation of all product ratings showed that only **6%** of products achieved a rating above **4.0**, emphasising the need for quality improvement.

The **Random Forest Regression** model confirmed that **customer engagement and satisfaction** are stronger drivers of product popularity than pricing.

Based on these insights, it is recommended that JC Penney:

- Focus on **improving product quality** to meet customer expectations;
- **Encourage more customer reviews** to strengthen engagement;
- **Promote high-rated but low-visibility products** through targeted marketing campaigns.

Body of the Report

- 1. Purpose
- 2. Analysis and Findings
- 3. Forecast of Popularity Factors (Machine Learning Model)
- 4. Conclusion
- 5. Recommendations
- **Appendix:** Python-based solution

1. Purpose

The purpose of this report is to analyse product and customer data from JC Penney, an American retail company, in order to identify potential issues and gain key insights for improvement.

The analysis focuses on understanding:

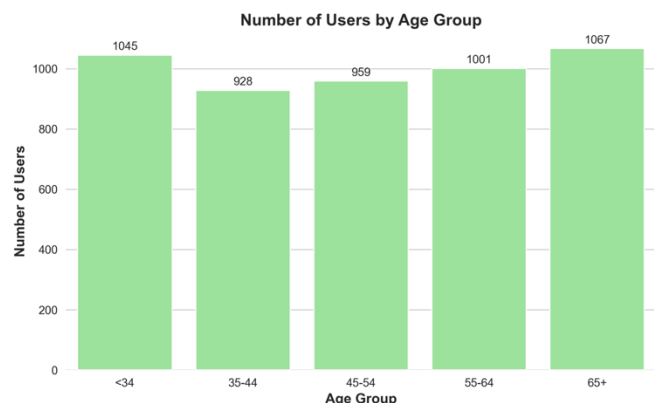
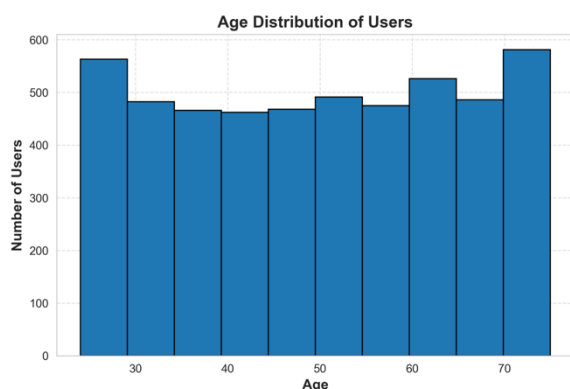
- The age distribution of customers and target audience segments;
- The relationship between product popularity and customer satisfaction;
- The identification of high-performing categories;
- The prediction of future popularity factors using a machine-learning model.

All data processing and analysis were performed in **Python**, applying statistical techniques, visualisation tools, and a **Random Forest Regression model** to explore the main drivers of product popularity.

2. Analysis and Findings

1) Age Distribution

To find out target audience customer by age, we need to see their age distribution.



Customers were grouped as follows for more detailed analysis:

Under 34 years — *Young Adults*

35–44 — *Adults 1*

45–54 — *Adults 2*

55–64 — *Older Adults*

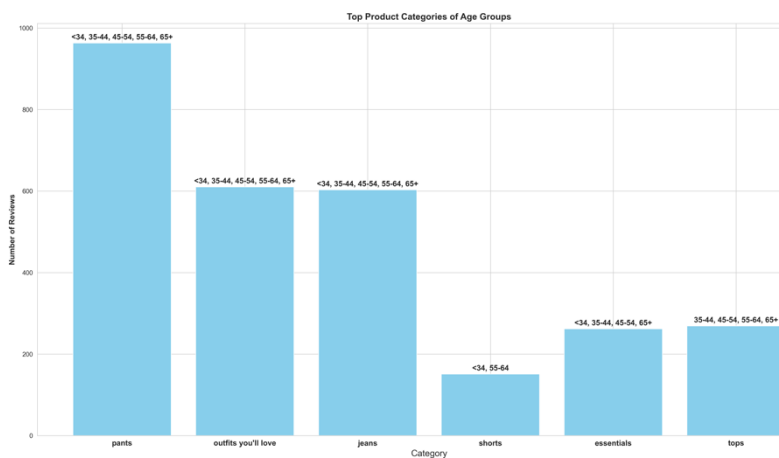
65 and above — *Elderly People*

The analysis shows that customer ages are relatively evenly distributed, with no significant peaks. Customers represent a wide range of ages, from 24 to 72 years old. The difference between the most and least common ages is relatively small, indicating that customers of all age groups are well represented and that the store appeals to a broad, diverse audience without a pronounced age bias.

2) Analysis of the 5 top categories for each age group

The analysis identified the five most frequently purchased categories for each age group.

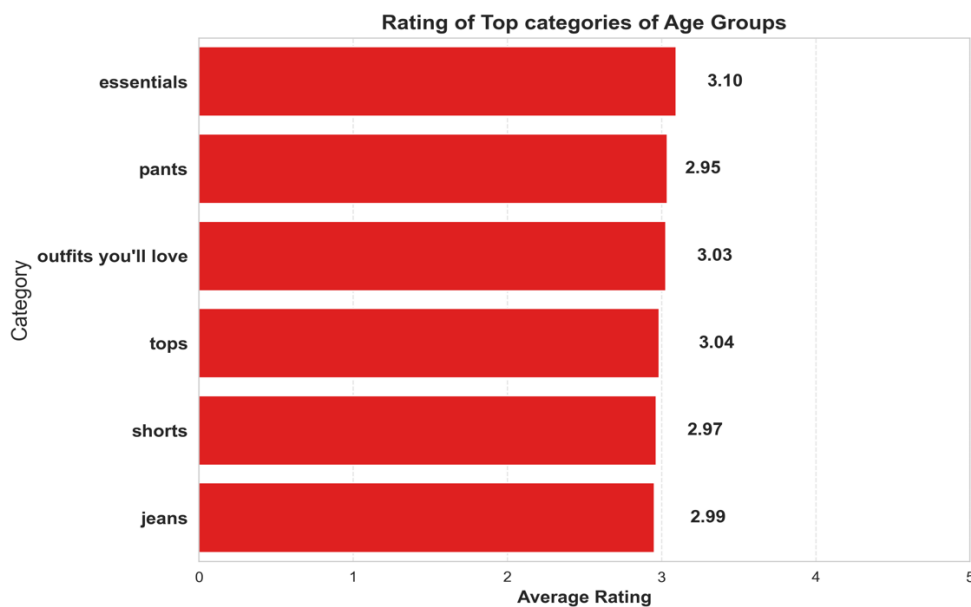
	<34	35-44	45-54	55-64	65+
0	pants (195)	pants (179)	pants (175)	pants (204)	pants (210)
1	outfits you'll love (139)	jeans (106)	outfits you'll love (115)	jeans (135)	jeans (130)
2	jeans (136)	outfits you'll love (104)	jeans (96)	outfits you'll love (125)	outfits you'll love (127)
3	shorts (76)	tops (67)	essentials (62)	shorts (75)	tops (80)
4	essentials (69)	essentials (62)	tops (58)	tops (64)	essentials (69)



There are **6 categories** that the most frequently purchased by all age groups:

- **Pants;**
- **Outfits you'll love;**
- **Jeans;**
- **Shorts;**
- **Essentials;**
- **Tops.**

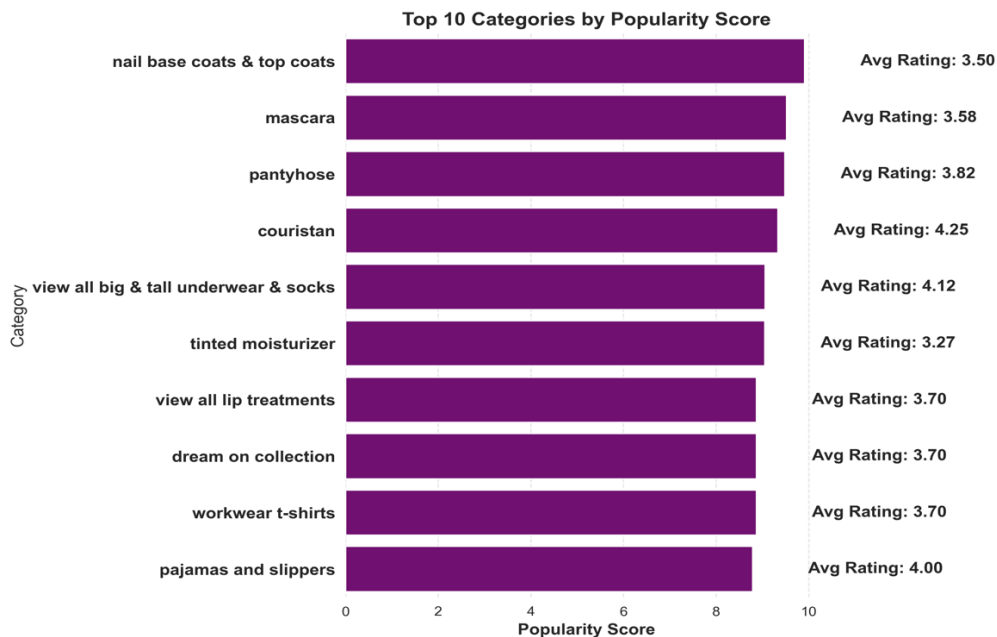
The analysis of the rating of the above categories yielded the following results:



Visualisations show that the most popular product categories across age groups receive relatively low average ratings. This suggests a problem with product quality that does not meet customers' expectations.

3) High-Popularity and Low-Rating Categories

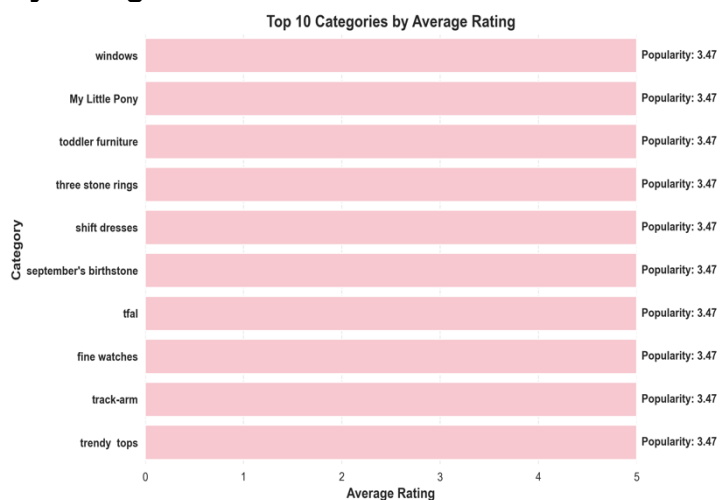
The ten most popular categories were examined to compare their popularity levels with their average ratings.



Although these categories appear statistically high-performing, their relatively low ratings confirm earlier findings that popular products often fail to meet quality expectations. *This suggests a broader issue with product quality across the range.* Such items may be damaging customer trust and require a **quality review, supplier evaluation, or product redesign.**

4) High-Rated but Low-Popularity Categories

Some categories have high average ratings but low popularity scores. This likely reflects **limited exposure or niche demand** rather than product excellence alone. These items represent **untapped marketing opportunities** — products that customers love but few people know about.



The company could consider promoting or featuring these products in marketing campaigns to increase visibility and sales.

5) High-Rated and High-Popularity Categories

The scatter plot presents all product categories.

Categories highlighted in red represent those achieving exceptionally high customer satisfaction (average rating above 4.0)

These best - performing categories should be viewed as **strategic benchmarks** for improving product quality and enhancing overall customer satisfaction.

The fact that only a few categories meet these standard underscores the challenge of maintaining consistent product quality.



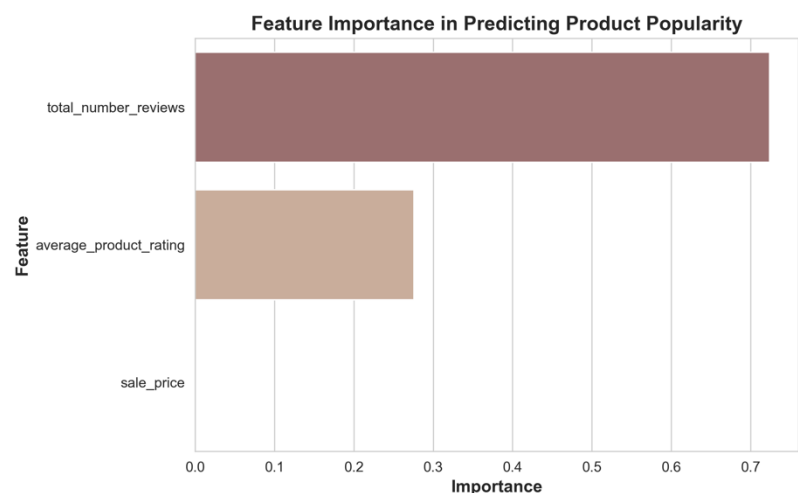
3. Forecast of Popularity Factors (Machine Learning Model)

A **Random Forest Regressor** was applied to forecast and evaluate the factors contributing to product popularity within the JC Penney dataset.

3 key predictors were used:

- **Sale price**
- **Average product rating**
- **Total number of reviews**

After training the model with an 80/20 train-test split, the feature importance analysis showed:



- **Total number of reviews (0.72)** – the strongest influence on popularity;
- **Average product rating (0.28)** – a moderate positive influence;
- **Sale price (~0.00)** – negligible influence.

These results indicate that popularity is driven primarily by **customer engagement** and **satisfaction**, rather than by pricing.

Encouraging customers to leave more reviews and maintaining high satisfaction levels will be more effective for increasing popularity than price adjustments.

4. Conclusion

The analysis shows that JC Penney serves a wide range of customers with no prioritised age group.

The most popular product categories tend to have low customer ratings, while high-rated products are less popular — suggesting insufficient promotional activity.

Only **6%** of products meet customer quality expectations (rating > 4.0), highlighting a broader challenge in maintaining consistent quality.

The machine-learning model confirmed that **customer engagement and satisfaction** are the primary drivers of product popularity, rather than pricing strategies.

5. Recommendations

- **Improve product quality** to better meet customer expectations.
- **Encourage customers to leave more reviews** to strengthen engagement.
- **Promote high-rated but low-visibility products** through marketing campaigns.
- Focus strategic investments on **product design, durability, and customer satisfaction**, as these factors have a stronger impact on popularity and sales than pricing alone.

Appendix

```
In [3]: import pandas as pd
        from datetime import datetime
        import json
```

```
In [4]: # Upload users table
        users=pd.read_csv('users.csv', index_col=0)
        display(users.head())
```

	State	Age
DOB		
1983-07-31	Oregon	42
1998-07-27	Massachusetts	27
1950-08-08	Idaho	75
1969-08-03	Florida	56
2001-07-26	Georgia	24

```
In [5]: # Upload products table
        products=pd.read_csv('products.csv', index_col=0)
        display(products[:3])
```

	SKU	Name	Description
Uniq_id			
b6c0b6bea69c722939585baeac73c13d	pp5006380337	Alfred Dunner® Essential Pull On Capri Pant	You'll return to our Alfred Dunner pull-on capr...
93e5272c51d8cce02597e3ce67b7ad0a	pp5006380337	Alfred Dunner® Essential Pull On Capri Pant	You'll return to our Alfred Dunner pull-on capr...
013e320f2f2ec0cf5b3ff5418d688528	pp5006380337	Alfred Dunner® Essential Pull On Capri Pant	You'll return to our Alfred Dunner pull-on capr...

```
In [6]: # Upload reviews table
```



```
reviews = pd.read_csv('reviews.csv', index_col=1)
display(reviews[:3])
```

	Uniq_id	Score	Review
Username			
fsdv4141	b6c0b6bea69c722939585baeac73c13d	2	You never have to worry about the fit...Alfred...
krpz1113	b6c0b6bea69c722939585baeac73c13d	1	Good quality fabric. Perfect fit. Washed very ...
mbmg3241	b6c0b6bea69c722939585baeac73c13d	2	I do not normally wear pants or capris that ha...

```
In [7]: # Upload jcpenny_products.json table
records = []
with open('jcpenny_products.json', 'r') as f:
    for line in f:
        line = line.strip()
        if line: # skip empty lines
            try:
                record = json.loads(line) # parse each JSON object
                records.append(record)
            except json.JSONDecodeError:
                # If the line is partial or malformed, try to fix/replace
                pass
jcpenny_products = pd.DataFrame(records)
if 'uniq_id' in jcpenny_products.columns:
    jcpenny_products.rename(columns={'uniq_id': 'Uniq_id'}, inplace=True)
    jcpenny_products.set_index('Uniq_id', inplace=True)
display(jcpenny_products[:3])
```

Uniq_id		sku	name_title	description
b6c0b6bea69c722939585baeac73c13d	pp5006380337		Alfred Dunner® Essential Pull On Capri Pant	You'll return to our Alfred Dunner pull-on cap...
93e5272c51d8cce02597e3ce67b7ad0a	pp5006380337		Alfred Dunner® Essential Pull On Capri Pant	You'll return to our Alfred Dunner pull-on cap...
013e320f2f2ec0cf5b3ff5418d688528	pp5006380337		Alfred Dunner® Essential Pull On Capri Pant	You'll return to our Alfred Dunner pull-on cap...

```
In [8]: # Upload jcpenny_reviewers.json table
jcpenny_reviewers = pd.read_json('jcpenny_reviewers.json', lines=
jcpenny_reviewers.set_index('Username', inplace=True)
display(jcpenny_reviewers.head())
```

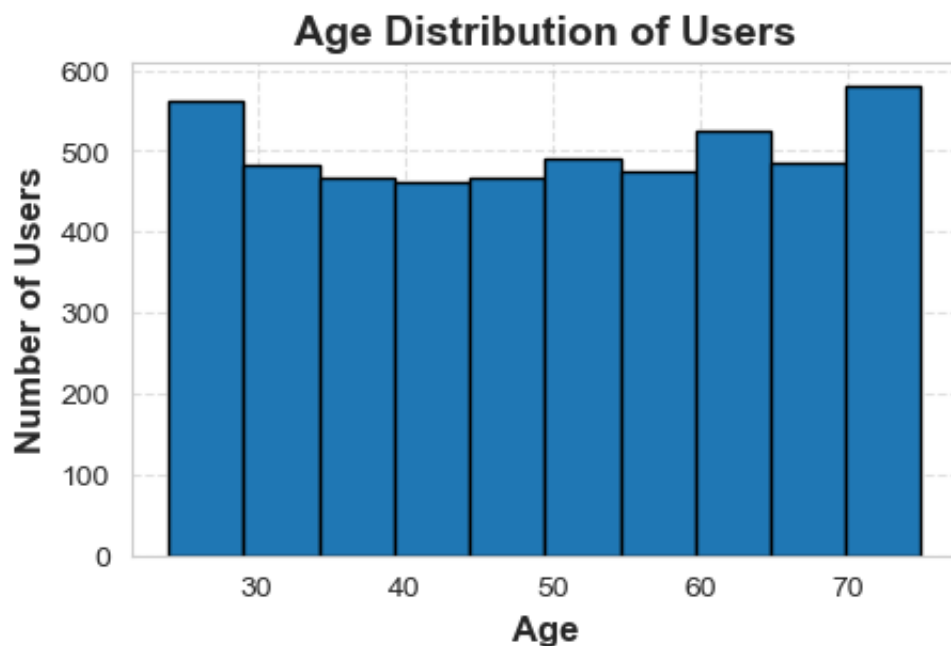
	DOB	State	Reviewed
Username			
bkpn1412	31.07.1983	Oregon	[cea76118f6a9110a893de2b7654319c0]
gqjs4414	27.07.1998	Massachusetts	[fa04fe6c0dd5189f54fe600838da43d3]
eehe1434	08.08.1950	Idaho	[]
hkxj1334	03.08.1969	Florida	[f129b1803f447c2b1ce43508fb822810, 3b0c9bc0be6...]
jjbd1412	26.07.2001	Georgia	[]

```
In [9]: # Add Ages of customers
from datetime import date
users = pd.read_csv('users.csv')
users['DOB'] = pd.to_datetime(users['DOB'], errors='coerce')
def calculate_age(born):
    if pd.isnull(born):
        return None
    today = date.today()
    return today.year - born.year - ((today.month, today.day) < (bo
users['Age'] = users['DOB'].apply(calculate_age)
users = users[['DOB', 'State', 'Age']]
users.to_csv('users.csv', index=False)
```

```
display(users[:3])
```

	DOB	State	Age
0	1983-07-31	Oregon	42
1	1998-07-27	Massachusetts	27
2	1950-08-08	Idaho	75

```
In [27]: # Visualisation of Ages
import matplotlib.pyplot as plt
plt.figure(figsize=(5, 3))
plt.hist(users['Age'], bins=10, edgecolor='black')
plt.title('Age Distribution of Users', fontsize=14, fontweight='bold')
plt.xlabel('Age', fontsize=12, fontweight='bold')
plt.ylabel('Number of Users', fontsize=12, fontweight='bold')
plt.grid(True, linestyle='--', alpha=0.6)
plt.savefig('AgeDistribution.png', dpi = 300)
plt.show()
```



Observation

The age distribution of customers is relatively even, with no significant peaks. Customers represent a wide range of age groups, from 24 to 72 years old. The difference between the most and least common ages is relatively small, indicating that customers of all age groups are well represented and that the store appeals to a broad, diverse audience without a pronounced age bias.

```
In [11]: # Ages groups with numbers of customers
bins = [0, 34, 44, 54, 64, 100]
labels = ['<34', '35-44', '45-54', '55-64', '65+']
users['Age_Group'] = pd.cut(users['Age'], bins=bins, labels=labels,
users.to_csv('users_with_age_group.csv', index=False)
```

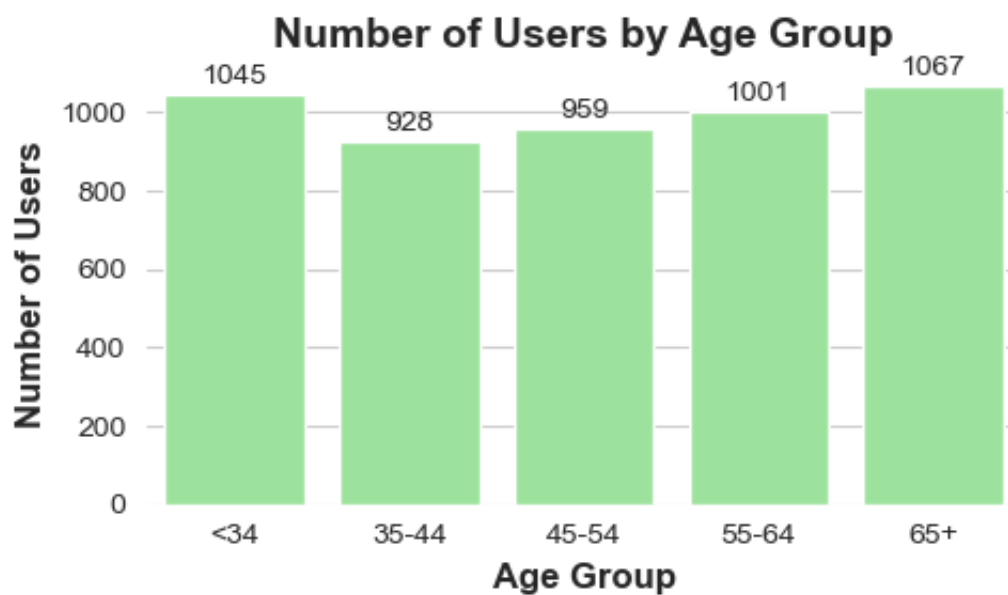
```
age_group_counts = users['Age_Group'].value_counts().sort_index()
print("Number of users per age group:")
display(age_group_counts)
```

Number of users per age group:

```
Age_Group
<34      1045
35-44     928
45-54     959
55-64    1001
65+     1067
```

Name: count, dtype: int64

```
In [26]: # Visualisation of Ages Groups
import seaborn as sns
plt.figure(figsize=(5, 3))
sns.set_style('whitegrid')
sns.barplot(x=age_group_counts.index, y=age_group_counts.values, color='green')
for i, value in enumerate(age_group_counts.values):
    plt.text(i, value + 10, str(value), ha='center', va='bottom', fontweight='bold')
plt.title('Number of Users by Age Group', fontsize=14, fontweight='bold')
plt.xlabel('Age Group', fontsize=12, fontweight='bold')
plt.ylabel('Number of Users', fontsize=12, fontweight='bold')
sns.despine(left=True, bottom=True)
plt.tight_layout()
plt.savefig('Ages groups.png', dpi = 300)
plt.show()
```



Observation

We divide customers into the following age groups:

- under 34 years old Young Adults
- 35–44 Adults 1
- 45–54 Adults 2
- 55–64 Older Adults

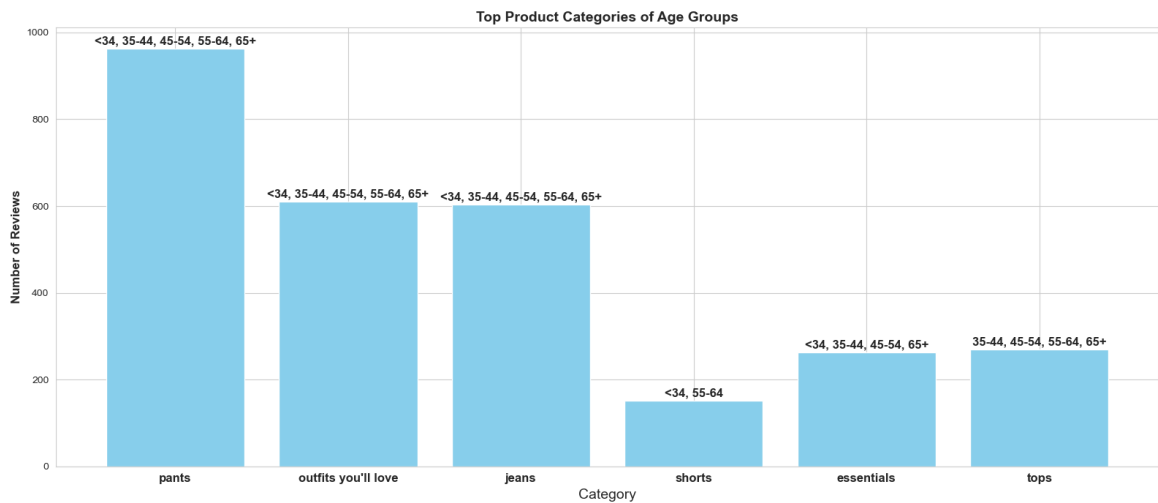
- Over 65 Elderly People

And the graph confirmed previous observation that the needs of each age group are being met by the company, as none of the groups appears to be prioritised.

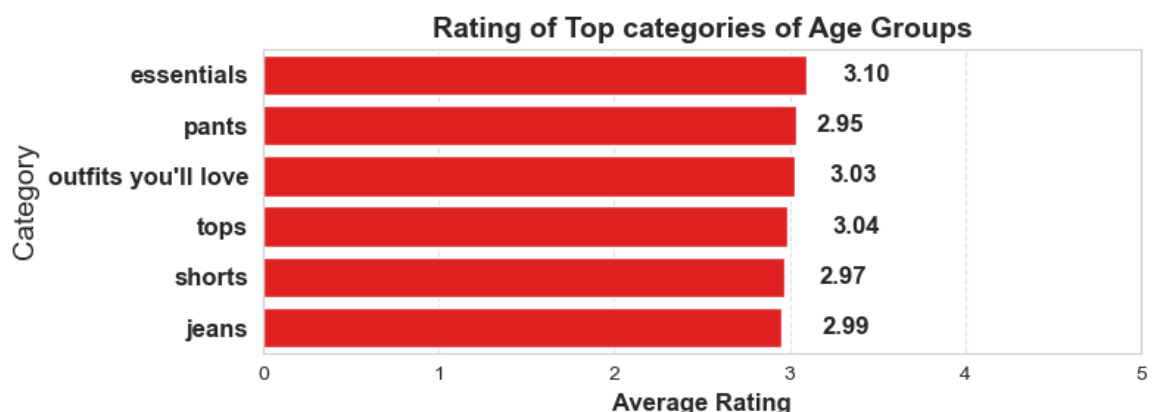
```
In [13]: # Top 5 categories for each Age group
jcpenny_products = pd.read_json('jcpenny_products.json', lines=True)
jcpenny_reviewers = pd.read_json('jcpenny_reviewers.json', lines=True)
reviews = pd.read_csv('reviews.csv', index_col=1)
products = pd.read_csv('products.csv', index_col=0)
users = pd.read_csv('users.csv', index_col=0)
from datetime import date
jcpenny_reviewers = jcpenny_reviewers.copy()
jcpenny_reviewers['DOB'] = pd.to_datetime(jcpenny_reviewers['DOB'])
def calculate_age(born):
    if pd.isnull(born):
        return None
    today = date.today()
    return today.year - born.year - ((today.month, today.day) < (born.month, born.day))
jcpenny_reviewers['Age'] = jcpenny_reviewers['DOB'].apply(calculate_age)
bins = [0, 34, 44, 54, 64, 100]
labels = ['<34', '35-44', '45-54', '55-64', '65+']
jcpenny_reviewers['Age_Group'] = pd.cut(jcpenny_reviewers['Age'], bins=bins, labels=labels)
# Products
if 'uniq_id' in jcpenny_products.columns:
    jcpenny_products.rename(columns={'uniq_id': 'Uniq_id'}, inplace=True)
# Join
reviews = reviews.reset_index().rename(columns={'index': 'Username'})
merged = (reviews.merge(jcpenny_reviewers[['Username', 'Age_Group']], on='Username'))
# Ignore useless categories
ignore_words = ['view all', 'view all brands', 'sale', 'clearance', 'discount']
merged = merged[merged['category'].notna() & (merged['category'].str.lower().isin(ignore_words) == False)]
# Top 5
category_counts = (merged.groupby(['Age_Group', 'category'], observed=True).count().reset_index())
top5_dict = {}
for age_group, group in category_counts.groupby('Age_Group', observed=True):
    top5 = group.sort_values('Review_Count', ascending=False).head(5)
    top5_list = [f"{row['category']} ({row['Review_Count']})" for _, row in top5.iterrows()]
    while len(top5_list) < 5:
        top5_list.append("")
    top5_dict[age_group] = top5_list
# Table
top5_df = pd.DataFrame(top5_dict)
display(top5_df)
top5_df.to_csv('Top5_Categories_by_Age.csv', index=False)
```

	<34	35-44	45-54	55-64	65+
0	pants (195)	pants (179)	pants (175)	pants (204)	pants (210)
1	outfits you'll love (139)	jeans (106)	outfits you'll love (115)	jeans (135)	jeans (130)
2	jeans (136)	outfits you'll love (104)	jeans (96)	outfits you'll love (125)	outfits you'll love (127)
3	shorts (76)	tops (67)	essentials (62)	shorts (75)	tops (80)
4	essentials (69)	essentials (62)	tops (58)	tops (64)	essentials (69)

```
In [32]: # Visualisation of Top categories of Age Groups
ignore_categories = ['view all', 'view all brands', 'sale', '', None]
data = category_counts[~category_counts['category'].isin(ignore_categories)]
top5_per_age_list = []
for age_group, group in data.groupby('Age_Group', observed=False):
    top5 = group.sort_values('Review_Count', ascending=False).head(5)
    top5_per_age_list.append(top5)
top5_per_age = pd.concat(top5_per_age_list, ignore_index=True)
categories = top5_per_age['category'].unique()
category_totals = []
age_labels = []
for cat in categories:
    subset = top5_per_age[top5_per_age['category'] == cat]
    total_reviews = subset['Review_Count'].sum()
    ages = ', '.join(subset['Age_Group'].astype(str).tolist())
    category_totals.append(total_reviews)
    age_labels.append(ages)
plt.figure(figsize=(16,7))
bars = plt.bar(categories, category_totals, color='skyblue')
for bar, label in zip(bars, age_labels):
    height = bar.get_height()
    plt.text(bar.get_x() + bar.get_width()/2, height + 5, label, ha='center')
plt.xlabel('Category', fontsize=14)
plt.ylabel('Number of Reviews', fontsize=12, fontweight='bold')
plt.title('Top Product Categories of Age Groups', fontsize=14, fontweight='bold')
plt.xticks(rotation=0, ha='center', fontsize=12, fontweight='bold')
plt.tight_layout()
plt.savefig('Top categories of Age Groups', dpi = 300)
plt.show()
```



```
In [35]: # Rating of the Top categories of Age groups
target_categories = ['pants','jeans','tops','essentials',"outfits you'll love"]
# Average ratings for these categories
if 'category' in jcpenny_products.columns and 'average_product_rating' in jcpenny_products.columns:
    avg_scores = (jcpenny_products[jcpenny_products['category'].isin(target_categories)]
                  .groupby('category', observed=False)['average_product_rating']
                  .mean()
                  .reset_index()
                  .rename(columns={'average_product_rating': 'Average Rating'})
                  .sort_values('Average Rating', ascending=False))
else:
    print("Missing columns: check if 'category' or 'average_product_rating' are in the data")
    avg_scores = pd.DataFrame()
plt.figure(figsize=(8, 3))
sns.barplot(data=avg_scores, x='Average Rating', y='category', color='red')
plt.title('Rating of Top categories of Age Groups', fontsize=14, fontweight='bold')
plt.xlabel('Average Rating', fontsize=12, fontweight='bold')
plt.ylabel('Category', fontsize=14)
plt.yticks(fontsize=12, fontweight='bold')
plt.xlim(0, 5)
plt.grid(axis='x', linestyle='--', alpha=0.5)
for index, row in avg_scores.iterrows():
    plt.text(row['Average Rating'] + 0.2, index, f"{row['Average Rating']:.2f}")
plt.tight_layout()
plt.savefig('Rating of Top categories of Age Groups', dpi = 300)
plt.show()
```



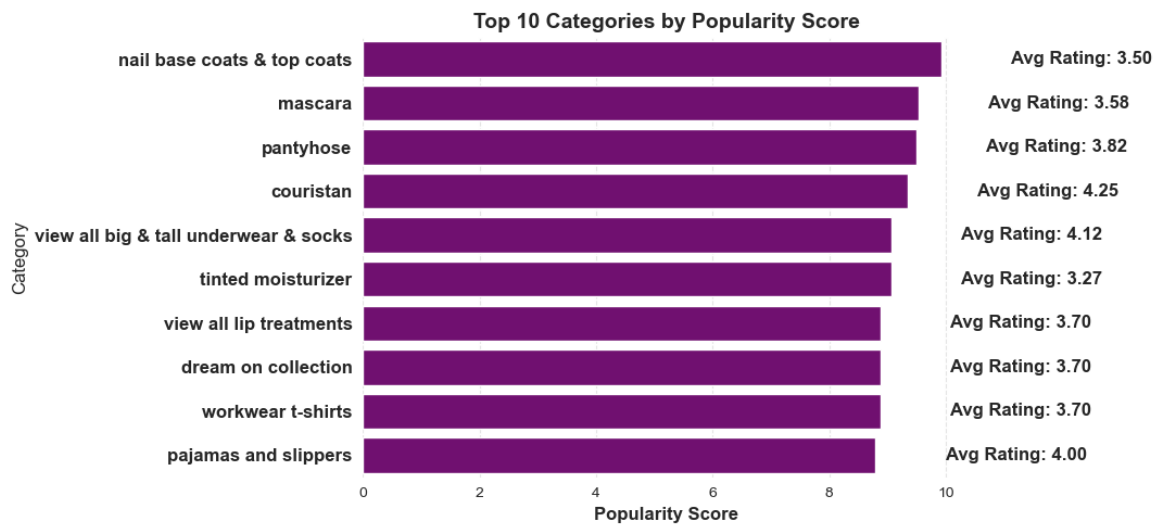
Observation

The most popular product categories, which indicate high demand, have very low ratings. This suggests a problem with product quality that does not meet customers' expectations.

```
In [16]: # Which categories are statistically high-performing, even with few
import numpy as np
df = jcpenny_products[jcpenny_products['total_number_reviews'] > 0]
df['popularity_score'] = df['average_product_rating'] * np.log1p(df['total_number_reviews'])
category_scores = (df.groupby('category', observed=False).agg({'popularity_score': 'sum',
                        'average_product_rating': 'mean'})).sort_values(by='popularity_score', ascending=False)
top_categories = category_scores.head(10).reset_index()
top_categories.rename(columns={'popularity_score': 'Popularity Score'})
display(top_categories)
```

	category	Popularity Score	Total Reviews	Average Rating
0	nail base coats & top coats	9.916247	16	3.500000
1	mascara	9.522681	27	3.575000
2	pantyhose	9.487825	11	3.818182
3	couristan	9.338204	8	4.250000
4	view all big & tall underwear & socks	9.063551	8	4.125000
5	tinted moisturizer	9.057123	15	3.266667
6	view all lip treatments	8.872213	10	3.700000
7	dream on collection	8.872213	10	3.700000
8	workwear t-shirts	8.872213	10	3.700000
9	pajamas and slippers	8.788898	8	4.000000

```
In [38]: # Visualisation of Top 10 Categories by Popularity Score
import seaborn as sns
plt.figure(figsize=(11, 5))
sns.barplot(data=top_categories, x='Popularity Score', y='category', color='black')
sns.despine(left=True, bottom=True)
plt.title('Top 10 Categories by Popularity Score', fontsize=14, fontweight='bold')
plt.xlabel('Popularity Score', fontsize=12, fontweight='bold')
plt.ylabel('Category', fontsize=12)
plt.yticks(fontweight='bold', fontsize=12)
for i, (score, rating) in enumerate(zip(top_categories['Popularity Score'], top_categories['Average Rating'])):
    plt.text(score + 1.2, i, f'Avg Rating: {rating:.2f}', va='center', color='red', fontweight='bold')
plt.grid(axis='x', linestyle='--', alpha=0.5)
plt.tight_layout()
plt.savefig('Top_10_Categories_by_Popularity_Score.png', dpi=300)
plt.show()
```

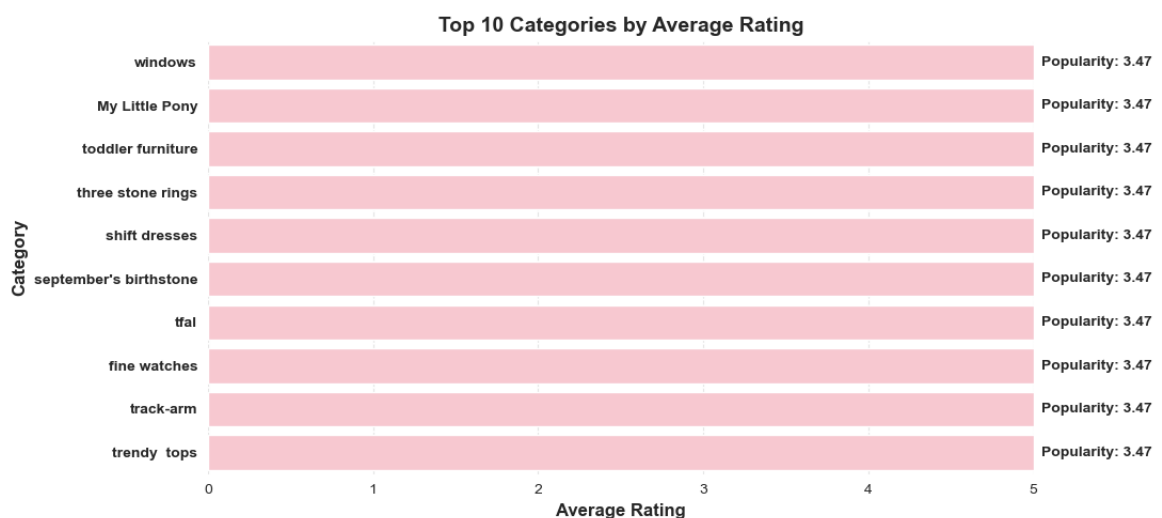
Observation

Although the resulting categories appear statistically high-performing, their relatively low ratings confirm earlier findings that popular products often fail to meet customer quality expectations. This suggests a broader issue with product quality across the range.

```
In [18]: # Top 10 Categories by Average Rating
df = jcpenny_products.copy()
df = df.dropna(subset=['average_product_rating', 'category'])
ignore_words = ['view all', 'view all brands', 'sale', 'clearance',
df = df[~df['category'].str.lower().isin(ignore_words)]
df['popularity_score'] = df['average_product_rating'] * np.log1p(df['total_number_reviews'])
top_categories_by_rating = (
    df.groupby('category', observed=False).agg({
        'average_product_rating': 'mean',
        'total_number_reviews': 'sum',
        'popularity_score': 'mean'}).reset_index())
top_categories_by_rating = top_categories_by_rating.sort_values('average_product_rating', ascending=False)
top_categories_table = top_categories_by_rating.copy()
top_categories_table.rename(columns={'category': 'Category', 'average_product_rating': 'Average Rating', 'total_number_reviews': 'Total Number Reviews', 'popularity_score': 'Popularity Score'})
top_categories_table['Average Rating'] = top_categories_table['Average Rating'].round(2)
top_categories_table['Popularity Score'] = top_categories_table['Popularity Score'].round(2)
top_categories_table = top_categories_table[['Category', 'Average Rating', 'Total Number Reviews', 'Popularity Score']]
top_categories_table = top_categories_table.reset_index(drop=True)
top_categories_by_rating = top_categories_by_rating.sort_values(by='average_product_rating', ascending=False)
display(top_categories_table)
top_categories_table.to_csv('Top 10 Categories by Average Rating.csv')
```

	Category	Average Rating	Popularity Score	Total Reviews
0	windows	5.0	3.47	1
1	My Little Pony	5.0	3.47	1
2	toddler furniture	5.0	3.47	1
3	three stone rings	5.0	3.47	1
4	shift dresses	5.0	3.47	1
5	september's birthstone	5.0	3.47	1
6	tfal	5.0	3.47	1
7	fine watches	5.0	3.47	1
8	track-arm	5.0	3.47	1
9	trendy tops	5.0	3.47	1

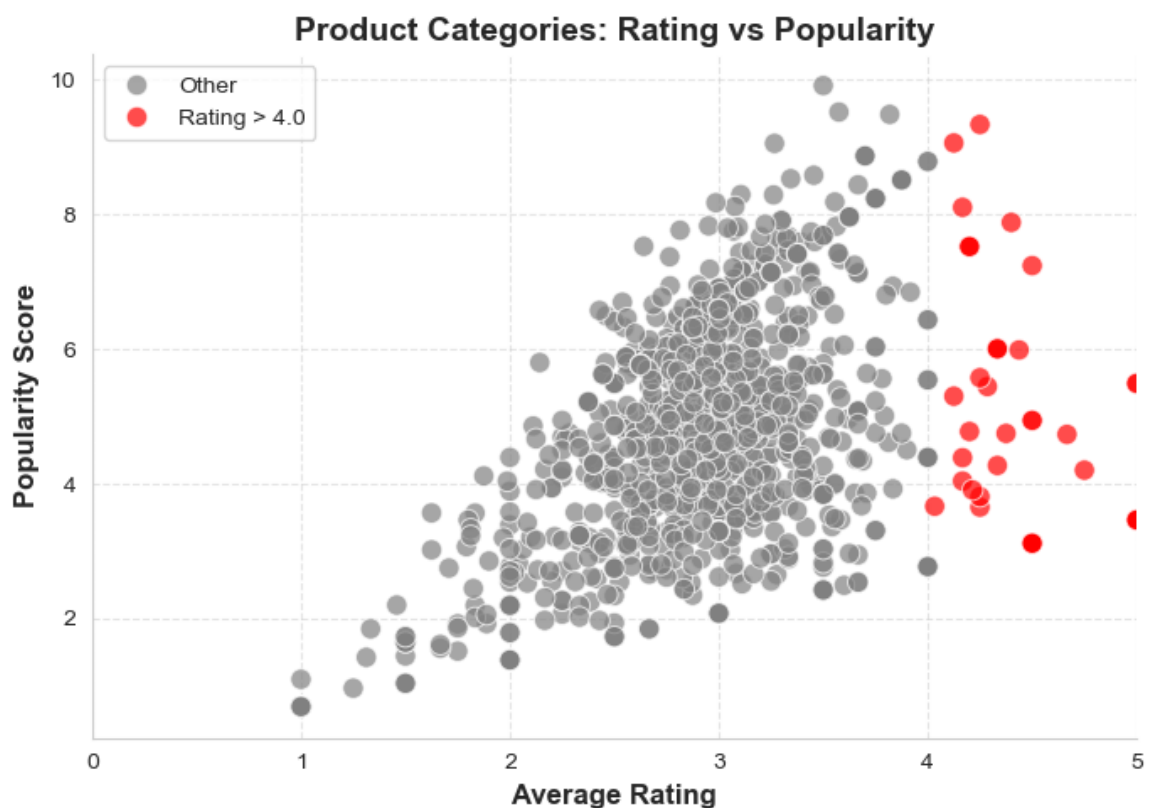
```
In [39]: # Visualisation of the Top 10 Categories by Average Rating
df_plot = top_categories_table.copy()
plt.figure(figsize=(11, 5))
sns.barplot(data=df_plot, x='Average Rating', y='Category', color='pink')
sns.despine(left=True, bottom=True)
plt.title('Top 10 Categories by Average Rating', fontsize=14, fontweight='bold')
plt.xlabel('Average Rating', fontsize=12, fontweight='bold')
plt.ylabel('Category', fontsize=12, fontweight='bold')
plt.yticks(fontweight='bold')
for i, (rating, score) in enumerate(zip(df_plot['Average Rating'], df_plot['Popularity Score'])):
    plt.text(rating + 0.05, i, f'Popularity: {score:.2f}', va='center')
plt.xlim(0, 5)
plt.grid(axis='x', linestyle='--', alpha=0.5)
plt.tight_layout()
plt.savefig('Top_10_Categories_by_Average_Rating.png', dpi=300)
plt.show()
```



Observation

Although several categories show high average ratings but low popularity scores, this likely reflects limited exposure or niche demand rather than product excellence alone. These items may represent untapped marketing opportunities, suggesting that greater visibility could increase overall sales performance.

```
In [41]: # Visualisation of High rating categories
df = jcpenney_products[jcpenney_products['total_number_reviews'] > 0]
df['popularity_score'] = df['average_product_rating'] * np.log1p(df['total_number_reviews'])
category_stats = (df.groupby('category', observed=False).agg({
    'average_product_rating': 'mean',
    'total_number_reviews': 'sum',
    'popularity_score': 'mean'})).reset_index()
category_stats['Category Type'] = np.where(category_stats['average_product_rating'] > 4.0, 'High Rating', 'Other')
plt.figure(figsize=(7, 5))
sns.scatterplot(data=category_stats, x='average_product_rating', y='popularity_score', hue='Category Type')
plt.title('Product Categories: Rating vs Popularity', fontsize=14)
plt.xlabel('Average Rating', fontsize=12, fontweight='bold')
plt.ylabel('Popularity Score', fontsize=12, fontweight='bold')
plt.xlim(0, 5)
plt.grid(True, linestyle='--', alpha=0.5)
plt.legend(title='', loc='upper left')
sns.despine()
plt.tight_layout()
plt.savefig('Rating_vs_Popularity_All_Categories.png', dpi=300)
plt.show()
```



```
In [23]: # Calculate percentage of categories and products with rating > 4.0
high_rating_count = (category_stats['average_product_rating'] > 4.0).sum()
total_categories = len(category_stats)
percentage_high_rating = (high_rating_count / total_categories) * 100
```

```
print(f"Categories with average rating > 4.0: {high_rating_count} out of {total_products}")
high_rating_products = (df['average_product_rating'] > 4.0).sum()
total_products = len(df)
percentage_high_products = (high_rating_products / total_products) * 100
print(f"Products with rating > 4.0: {high_rating_products} out of {total_products} ({percentage_high_products}%)")
```

Categories with average rating > 4.0: 70 out of 1169 (5.99%)
 Products with rating > 4.0: 628 out of 7964 (7.89%)

Observation

This scatter plot presents all product categories.

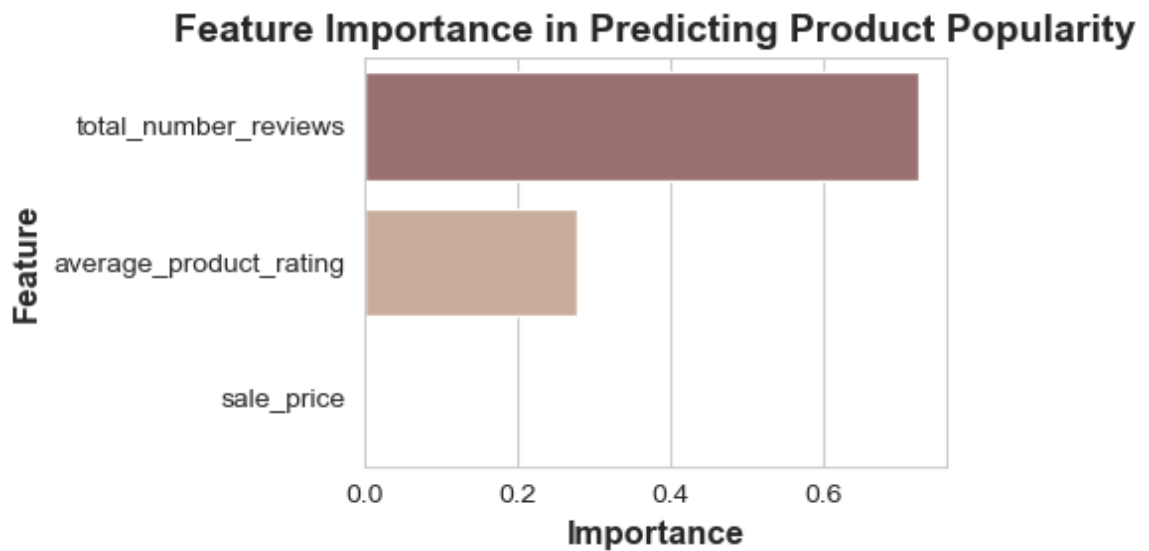
Categories highlighted in red represent those achieving exceptionally high customer satisfaction (average rating above 4.0).

The fact that only a few categories meet this standard underscores a broader challenge in maintaining consistent product quality.

These high-rated categories should be viewed as strategic benchmarks for improving the quality of other products and enhancing overall customer satisfaction.

```
In [42]: # Forecast of popularity factors
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split
df = jcpenny_products.copy()
df = df[df['total_number_reviews'] > 0]
df = df.dropna(subset=['average_product_rating', 'sale_price'])
import re
def clean_price(price):
    if pd.isna(price):
        return np.nan
    price = str(price).replace('$', '').strip()
    if '-' in price:
        parts = re.split(r'[-]', price)
        parts = [float(p) for p in parts if p.replace('.', '', 1).isdigit()]
        if len(parts) == 2:
            return np.mean(parts)
    try:
        return float(price)
    except:
        return np.nan
df['sale_price'] = df['sale_price'].apply(clean_price)
df = df.dropna(subset=['sale_price'])
df['popularity_score'] = df['average_product_rating'] * np.log1p(df['total_number_reviews'])
X = df[['sale_price', 'average_product_rating', 'total_number_reviews']]
y = df['popularity_score']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
model = RandomForestRegressor(random_state=42)
model.fit(X_train, y_train)
importances = pd.DataFrame({'Feature': X.columns, 'Importance': model.feature_importances_})
plt.figure(figsize=(5, 3))
sns.barplot(x='Importance', y='Feature', hue='Feature', data=importances)
plt.title('Feature Importance in Predicting Product Popularity', fontweight='bold')
plt.xlabel('Importance', fontsize=12, fontweight='bold')
```

```
plt.ylabel('Feature', fontsize=12, fontweight='bold')
plt.tight_layout()
plt.savefig('Forecast of popularity factors.png', dpi=300)
plt.show()
display(importances)
```



	Feature	Importance
2	total_number_reviews	0.724028
1	average_product_rating	0.275901
0	sale_price	0.000070

Observation

The feature importance analysis shows that the number of reviews is the strongest predictor of product popularity, followed by average customer rating, while price has almost no impact.

This indicates that encouraging customers to leave more reviews and maintaining high satisfaction levels will be far more effective for increasing product popularity than changing prices.