

Infosys Springboard Internship

Domain: Artificial Intelligence(AI)

Project Title:

AI STYLIST

Team Members:

- Akshay
- Dheeraj
- Lebin Bright
- Mohammed Ghouse
- Shameera Begum

Mentor:

Anil Shaw Sir

S.no	Content	Page.No
1	Dataset Overview	3
2	Introduction	4
3	Data Visualization	5
4	Data Cleaning	17
5	Stemming	19
6	Lemmitization	20
7	TF-IDF Vectorize with Cosine Similarity	21
8	Word2Vec Embeddings	22
9	TF-IDF-Weighted Word2Vec Embeddings	24
10	Brand Based Filtering using TF_IDF Vectorization	27
11	Collaborative Filtering	33
12	Overview of VGG16 Embedding-based Model	37
13	Recommndation System using Resnet 50 Embeddings	41

1.DatasetOverview

- **FileName:**fashion_products_data.ldjson
- **Format:**JSONLines(.ldjson)
- **Purpose:**Likelycontainsproduct-relateddata,typicallyusedfor fashionindustryanalysis.

2.DataframeSummary

- **Shape(Rows,Columns):**Tobedeterminedupondata inspection.
- **Columns:**Thedataframecontainsthefollowingcolumns(columnnamestobe extracted):
 - Example:product_id,name,category,price,stock_status.
- **ColumnData Types:**
Example:
 - product_id→Integer
 - name→String
 - price→Float
- **NullValues:**
Numberofmissingvalues ineachcolumn.

3.DescriptiveStatistics

- **NumericalColumns:**
 - Count,Mean,StandardDeviation,Min,Maxforcolumnslikeprice.
- **CategoricalColumns:**
 - Countofuniquecategories,mostfrequentvalues.

5.Insights

1. **DataCompleteness:**
 - Columnswithsignificantmissingdatashouldbeaddressed.
2. **PotentialAnalysis:**
 - Pricetrendsbycategory.
 - Stockavailabilityinsights.

DatasetInfo:

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 33 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   uniq_id                               30000 non-null  object
1   crawl_timestamp                       30000 non-null  object
2   asin                                   30000 non-null  object
3   product_url                           30000 non-null  object
4   product_name                           30000 non-null  object
5   image_urls_small                       29998 non-null  object
6   medium                                 29998 non-null  object
7   large                                  28841 non-null  object
8   browsenode                             29480 non-null  float64
9   brand                                  21857 non-null  object
10  sales_price                             27110 non-null  float64
11  weight                                  30000 non-null  object
12  rating                                  30000 non-null  float64
13  sales_rank_in_parent_category           25497 non-null  object
14  sales_rank_in_child_category            24851 non-null  object
15  delivery_type                           30000 non-null  object
16  meta_keywords                           30000 non-null  object
17  amazon_prime__y_or_n                    30000 non-null  object
18  parent__child_category__all             25497 non-null  object
19  best_seller_tag__y_or_n                 30000 non-null  object
...
31  formats__editions                        2 non-null      object
32  name_of_author_for_books                 1 non-null      object

```

Dtypes: float64(7), object(26)

Memory usage: 7.6+MB

Introduction:

In the realm of data analysis, particularly concerning fashion product datasets, several critical techniques and methodologies are employed to ensure that the data is both informative and actionable. This report delves into various topics covered in the Jupyter Notebook, including **data visualization**, **outliers**, **null values**, **stop words**, **short forms**, **punctuation**, **non-English words**, **stemming**, and **lemmatization**. Each section defines key concepts, discusses relevant libraries and functions, and evaluates their use cases, advantages, and disadvantages.

	uniq_id	crawl_timestamp	asin	product_url	product_name	image_urls_small	medium
0	26d41bdc1495de290bc8e6062d927729	2020-02-07 05:11:36 +0000	B07STS2W9T	https://www.amazon.in/Facon-Kalamkari-Handbloc...	LA' Facon Cotton Kalamkari Handblock Saree Blo...	https://images-na.ssl-images-amazon.com/images...	https://images-na.ssl-images-amazon.com/images...
1	410c62298852e68f34c35560f2311e5a	2020-02-07 08:45:56 +0000	B07N6TD2WL	https://www.amazon.in/Sf-Jeans-Pantaloons-T-Sh...	Sf Jeans By Pantaloons Men's Plain Slim fit T-...	https://images-na.ssl-images-amazon.com/images...	https://images-na.ssl-images-amazon.com/images...
2	52e31bb31680b0ec73de0d781a23cc0a	2020-02-06 11:09:38 +0000	B07WJ6WPN1	https://www.amazon.in/LOVISTA-Traditional-Prin...	LOVISTA Cotton Gota Patti Tassel Traditional P...	https://images-na.ssl-images-amazon.com/images...	https://images-na.ssl-images-amazon.com/images...
3	25798d6dc43239c118452d1bee0fb088	2020-02-07 08:32:45 +0000	B07PYSF4WZ	https://www.amazon.in/People-Printed-Regular-T...	People Men's Printed Regular fit T- Shirt	https://images-na.ssl-images-amazon.com/images...	https://images-na.ssl-images-amazon.com/images...
4	ad8a5a196d515ef09dfdaf082bdc37c4	2020-02-06 14:27:48 +0000	B082KXNM7X	https://www.amazon.in/Monte-Carlo-Cotton-Colla...	Monte Carlo Grey Solid Cotton Blend Polo Colla...	https://images-na.ssl-images-amazon.com/images...	https://images-na.ssl-images-amazon.com/images...

1. Data Visualization:

Definition:

Data visualization is the graphical representation of information and data. By using visual elements like charts, graphs, and maps, data visualization tools provide an accessible way to see and understand trends, outliers, and patterns in data.

Libraries and Functions

1. Matplotlib

- **Function:** plt.plot(), plt.bar(), plt.pie()
- **UseCase:** Creating static graphs for exploratory data analysis.
- **Advantages:** Highly customizable; supports a wide range of plots.
- **Disadvantages:** Steeper learning curve for complex visualizations.

2. Seaborn

- **Function:** sns.scatterplot(), sns.boxplot()
- **UseCase:** Statistical data visualization with attractive defaults.
- **Advantages:** Simplifies complex visualizations; integrates well with pandas.

- **Disadvantages:** Less flexible than Matplotlib for certain types of plots.

```
import pandas as pd
import matplotlib.pyplot as plt
import json

# Load the data from the .ldjson file
file_path = 'C:/Users/admin/fashion_products_data.ldjson' # Updated path
data = []

# Open and read data line by line
with open(file_path, 'r', encoding='utf-8') as file:
    for line in file:
        data.append(json.loads(line))

# Convert data into a DataFrame
df = pd.DataFrame(data)

# Data Cleaning and Preparation
# Convert sales_price and rating columns to numeric, handling errors
df['sales_price'] = pd.to_numeric(df['sales_price'], errors='coerce')
df['rating'] = pd.to_numeric(df['rating'], errors='coerce')

# Drop rows with NaN values in essential columns
df = df.dropna(subset=['product_url', 'sales_price', 'rating', 'brand'])

# Basic Analysis
print("Summary of the 'sales_price' column:")
print(df['sales_price'].describe())

print("\nTop 10 brands by number of products:")
print(df['brand'].value_counts().head(10))

# Visualization: Distribution of Sales Price
plt.figure(figsize=(10, 5))
plt.hist(df['sales_price'].dropna(), bins=30, color='skyblue', edgecolor='black')
plt.title("Distribution of Sales Price")
```

```

plt.xlabel("SalesPrice")
plt.ylabel("Frequency")
plt.grid(axis='y',linestyle='--',alpha=0.7)
plt.show()

# Visualization: Average Rating by Top 10 Brands
top_brands=df['brand'].value_counts().nlargest(10).index
df_top_brands = df[df['brand'].isin(top_brands)]
plt.figure(figsize=(12, 6))
box_data=[df_top_brands[df_top_brands['brand']==brand]['rating'].dropna()for
brand in top_brands]
plt.boxplot(box_data,labels=top_brands,patch_artist=True) plt.xticks(rotation=45)
plt.title("RatingDistributionbyTop10Brands")
plt.xlabel("Brand")
plt.ylabel("Rating")
plt.show()

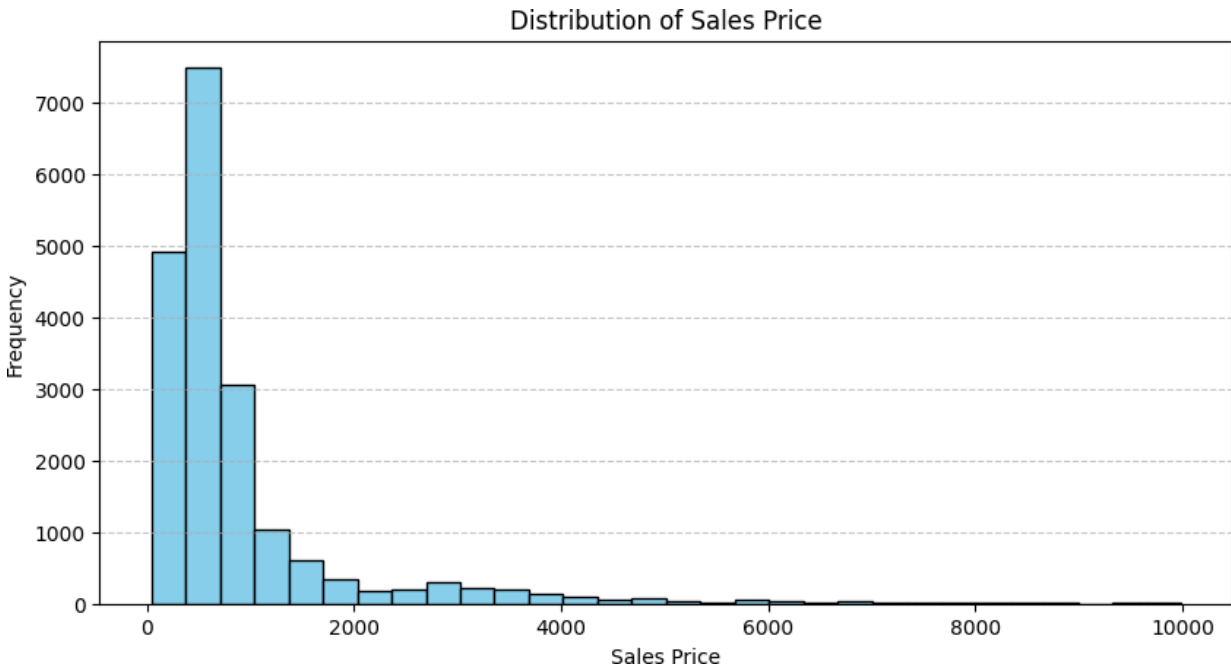
#Visualization:RelationshipbetweenSalesPriceandRating plt.figure(figsize=(8, 5))
forbrandintop_brands:
    brand_data = df[df['brand'] == brand]
    plt.scatter(brand_data['sales_price'],brand_data['rating'],label=brand,
alpha=0.6)

plt.title("SalesPricevsRating")
plt.xlabel("Sales Price")
plt.ylabel("Rating")
plt.legend(title="Brand",bbox_to_anchor=(1.05,1),loc='upperleft')

plt.grid(linestyle='--',alpha=0.7)
plt.show()

#CheckingforOutliersinSalesPrice
outliers=df[df['sales_price']>df['sales_price'].quantile(0.99)]
print("\nPotential Sales Price Outliers:")
print(outliers[['product_url', 'sales_price', 'brand', 'rating']])

```

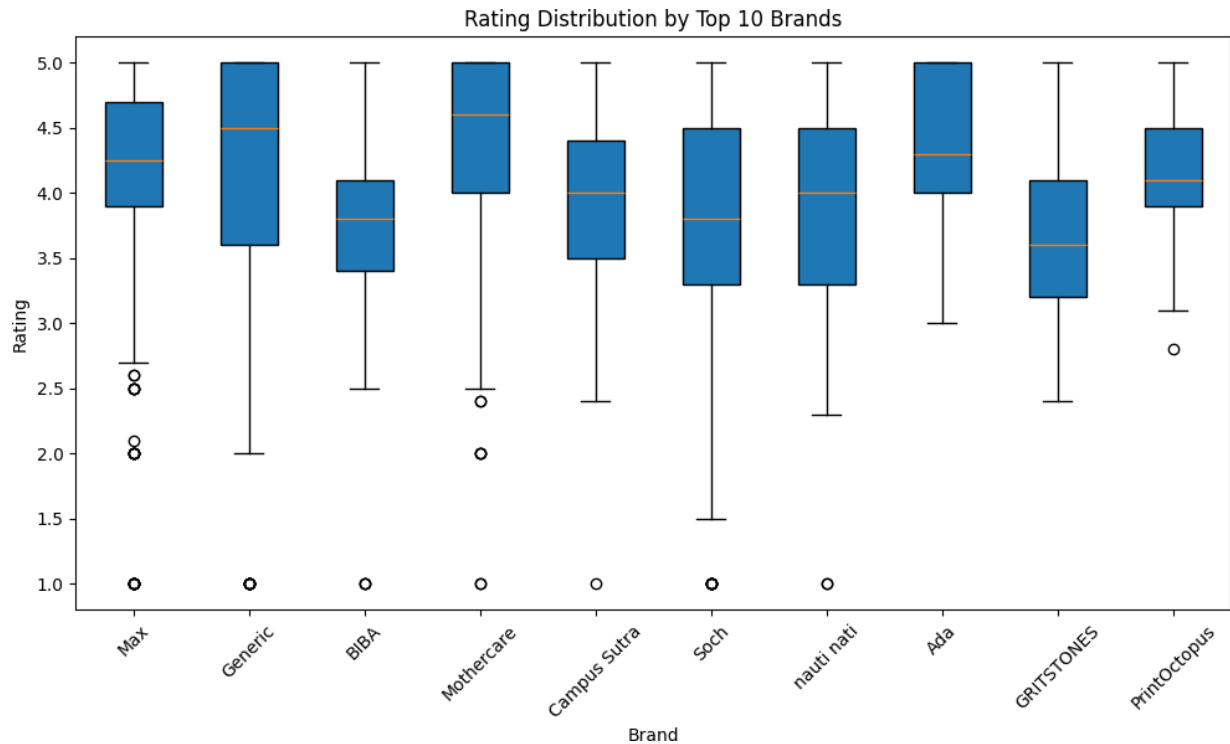


Histogram:

The histogram showcases the distribution of sales prices, revealing insights into the frequency of different price points.

Key Observation:

- From this above Histogram, we understood that it is a **Right-Skewed** distribution because the majority of sales occurs at lower prices, with the frequency tapering off as the price increases.
- The highest frequency of sales falls within the \$0 to \$2,000 range, **Peaking** around \$1,000.
- A significant number of sales are also present at higher price points, extending beyond \$8,000.
- This distribution suggests that the product or service being analyzed has a wider range of prices, appealing to both budget-conscious and higher-end consumers.



BoxPlot:

This Box Plot Represents the inter quartile range (IQR), which contains the middle 50% of the data. The bottom of the box is the first quartile (25th percentile), and the top is the third quartile (75th percentile).

- **The line inside the box:**

This indicates the median (50th percentile), which is the middle value of the data.

- **The dots:**

It represents outliers, which are data points that fall outside the whiskers.

Outliers:

Definition:

Outliers are data points that differ significantly from other observations. They can skew results and lead to misleading interpretations if not addressed properly.

Libraries and Functions

- **Pandas**
- **Function:** `df.describe()`, `df[df['column'] > threshold]`
- **Use Case:** Identifying outliers based on statistical measures like IQR or Z-score.
- **Advantages:** Easy to implement; integrates seamlessly with DataFrame operations.
- **Disadvantages:** May require manual threshold setting.

Example Code:

```
Q1=df['Sales'].quantile(0.25)
Q3=df['Sales'].quantile(0.75)
IQR =Q3- Q1
outliers=df[(df['Sales']<(Q1-1.5*IQR))|(df['Sales']> (Q3+1.5*IQR))]
```

1. We are Three Quartiles:

- **Q1 (First Quartile):** Represents the 25th percentile of the Sales data. This is the value below which 25% of the data lies.
- **Q2 (Second Quartile):** Represents the 50th percentile of the Sales data. This is the value below which 50% of the data lies.
- **Q3 (Third Quartile):** Represents the 75th percentile of the Sales data. This is the value below which 75% of the data lies.

2. Compute the Interquartile Range (IQR):

- **IQR:** Measures the spread of the middle 50% of the data. It is the difference between Q3 and Q1.
- **Formula:** $IQR = Q3 - Q1$

3. Define Outlier Boundaries:

- **Lower Bound:** Any data point less than $Q1 - 1.5 \times IQR$ is considered an outlier.

- **UpperBound:** Any datapoint greater than $Q3 + 1.5 \times IQR$ is considered an outlier.

The factor 1.5 is a standard multiplier used in statistics to define "**mild**" outliers.

4. Identify Outliers:

- **Condition:** The code uses boolean indexing to filter rows in the Sales column:
 - Rows where Sales is less than $Q1 - 1.5 \times IQR$ (extremely low values).
 - Rows where Sales is greater than $Q3 + 1.5 \times IQR$ (extremely high values).
- **Output:** The resulting outliers DataFrame contains all rows without outlier values in the Sales column.



ScatterPlot:

From the above diagram, The **Scatter Plot** shows the relationship between Sales Price and Rating for different brands.

- **Keypoints:**

- X-axis:**

- It Represents the Sales Price.

- Y-axis:**

- It Represents the Rating.

- Dots:**

- In this, Each dot represents a single product.

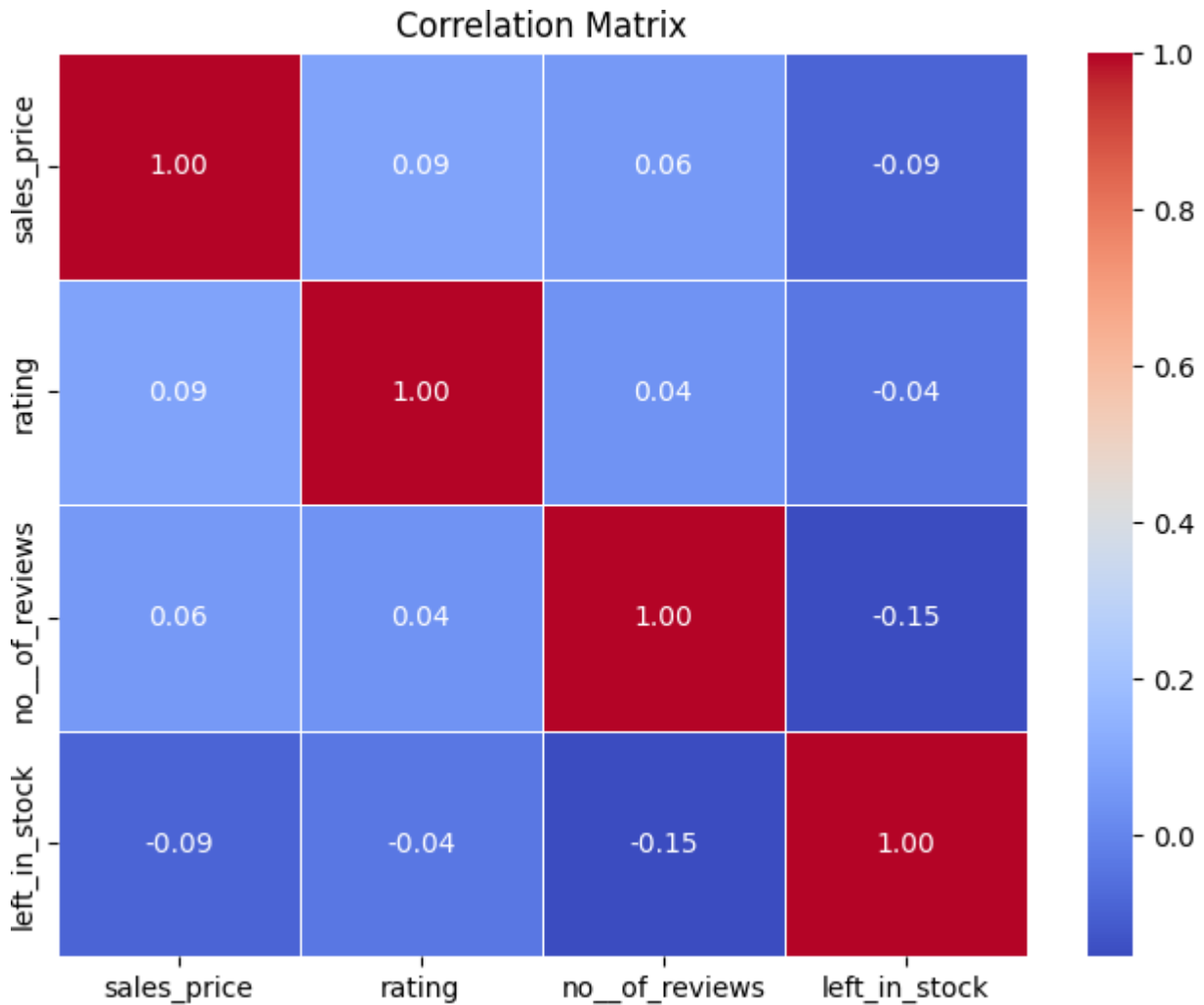
- Colors:**

- It Shows Different colors represent different brands.

- Interpretation:**

- The plot indicates that there isn't a strong correlation between Sales Price and Rating.

- Some products with high ratings have lower sales prices, while some with lower ratings have higher sales prices.
- Brand might be a more significant factor influencing sales price than rating.



Correlation Matrix:

The above image shows a correlation matrix, which is a table showing the correlation coefficients between multiple variables. In this case, the variables are: sales_price, rating, no_of_reviews, and left_in_stock.

- **Keypoints:**

- **Diagonal values:**

- The diagonal values are always 1, as a variable always perfectly correlates with itself.

➤ **Positive correlation:**

A positive value indicates that the two variables tend to move in the same direction. For instance, the 0.09 correlation between "sales_price" and "rating" suggests that higher-priced items tend to have slightly higher ratings.

➤ **Negative correlation:**

A negative value indicates that the two variables tend to move in opposite directions. For example, the -0.15 correlation between "no_of_reviews" and "left_in_stock" implies that items with more reviews tend to have less stock left.

➤ **Strength of correlation:**

The closer the value is to 1 or -1, the stronger the correlation. Values closer to 0 indicate a weaker relationship.

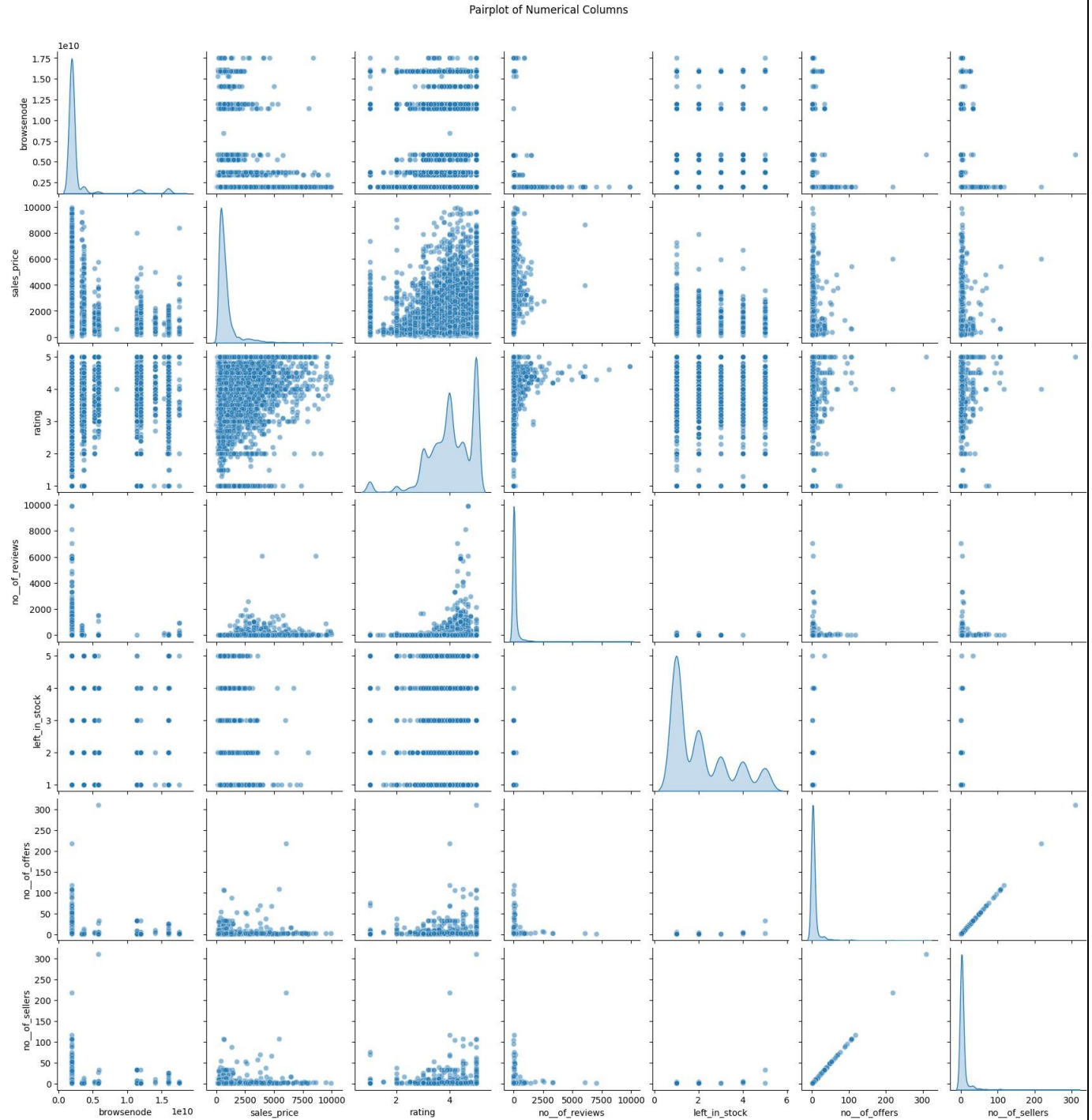
➤ **Interpreting the matrix:**

In this example, there are no strong correlations.

- The strongest positive correlation is between "sales_price" and "rating" (0.09).
- The strongest negative correlation is between "no_of_reviews" and "left_in_stock" (-0.15).
- Other relationships are weak, indicating little to no linear association between the variables.

```
#BIVARIATE ANALYSIS Assuming 'data' is your DataFrame

#1. Numerical vs. Numerical
# Pairplot to visualize pairwise relationships
sns.pairplot(data, diag_kind='kde', plot_kws={'alpha':0.5})
plt.suptitle("Pairplot of Numerical Columns", y=1.02)
plt.show()
```



PairPlot:

The image shows a pairplot, also known as a scatterplot matrix, which is a powerful tool for data visualization in statistics and machine learning. It helps to understand the relationships between multiple variables in a dataset at a glance.

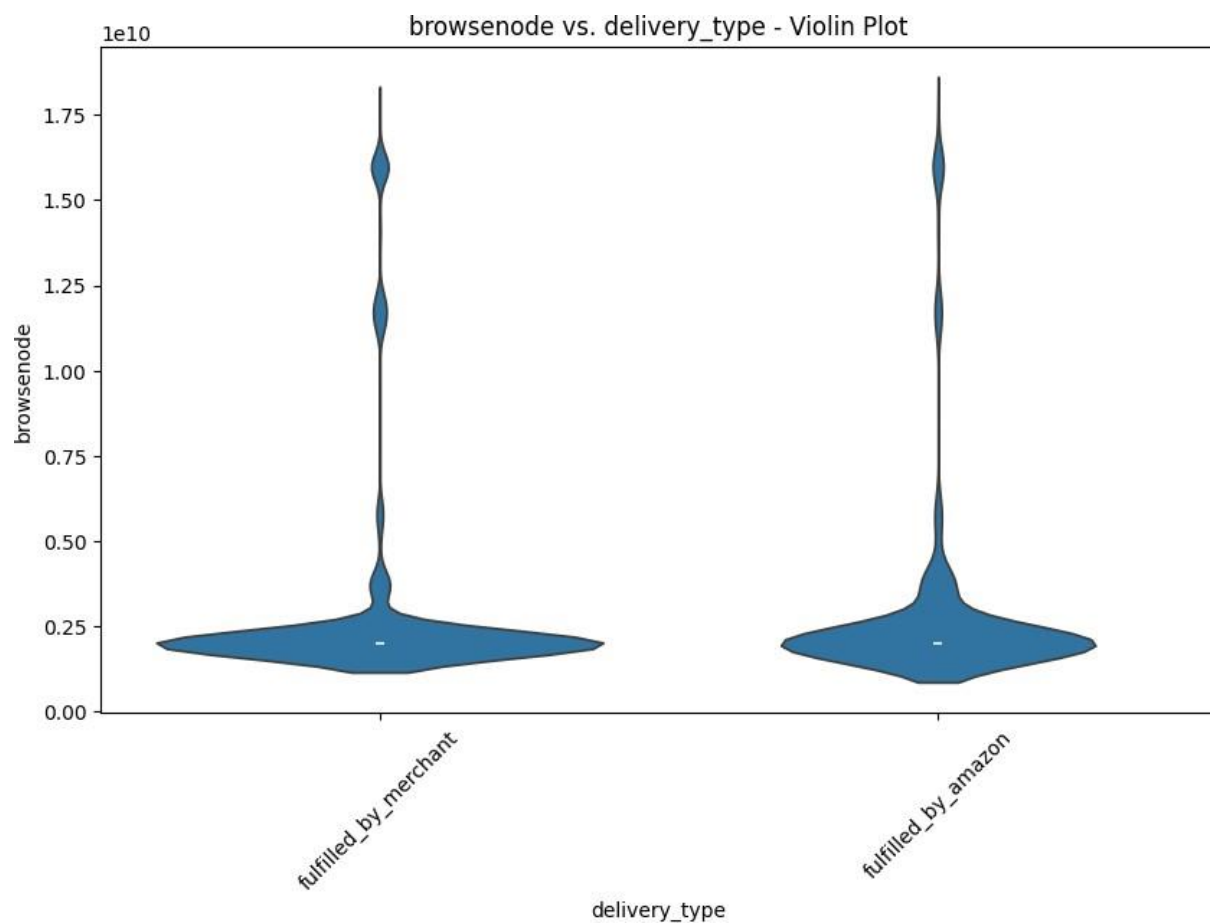
KeyPoints:

- Diagonal:

The diagonal cells display histograms for each variable, showing the distribution of their values. This helps to understand the frequency and spread of each variable individually.

- Off-Diagonal:

The off-diagonal cells display scatterplots for each pair of variables. Each dot in a scatterplot represents a single data point, with its x-coordinate corresponding to the value of one variable and its y-coordinate corresponding to the value of another variable.



Data Cleaning:

1. Null Handling

Null values are missing or undefined entries in a dataset. Effective handling is crucial for maintaining the integrity of the analysis.

Methods to Handle Nulls:

- **Removal of Nulls:** Drop rows or columns with null values if the missing data is minimal or non-essential.
 - `df.dropna()` to remove rows/columns with nulls.
- **Imputation:** Replace null values with meaningful substitutes like:
 - Mean, median, or mode for numerical data.
 - A placeholder value or the most frequent category for categorical data.
 - Example: `df['column'].fillna(df['column'].mean(), inplace=True)`.
- **Flagging:** Add a new column indicating rows with missing values for future reference.

Tools: pandas provides utilities like `isnull()`, `notnull()`, and `fillna()`.

2. Stop Words

Stop words are commonly used words (e.g., "the", "and", "is") that often do not add significant meaning to text analysis.

Why Remove Stop Words?

- They reduce noise in natural language processing (NLP) tasks.
- Improve computational efficiency and focus on meaningful terms.

LIBRARIES:

- Use predefined lists from libraries like `nlTK` or `spaCy`.
- Example with `nlTK`:

```
python
Copy code
from nltk.corpus import stopwords
stop_words = set(stopwords.words('english'))
filtered_text = [word for word in text.split() if word not in stop_words]
```

3. Non-English Words

Non-English words in a dataset can interfere with NLP tasks focused on English language data.

Detection and Removal:

- **Language Detection:** Use libraries like `langdetect` or `langid` to identify the language of text.
 - Example with `langdetect`:

```
python
Copycode
from langdetect import detect
if detect(text) == 'en':
```

4. Punctuation

Punctuation marks (e.g., !, ?, ,,) may not hold significant meaning in NLP tasks but could be valuable for sentiment or style analysis.

Handling Punctuation:

- **Removal:** Strip punctuation to focus on the core text.
 - Example with `string.punctuation`:
- **Retention:** Retain punctuation for tasks like sentiment analysis or where punctuation conveys meaning.
- **Custom Handling:** Preserve certain marks (e.g., @ for mentions) based on the use case.

```
python
Copycode
import string
text = text.translate(str.maketrans("", string.punctuation))
```

5. Short Form to Full Form

Expanding abbreviations or contractions is essential for standardizing text and improving clarity.

Examples:

- Short forms: can't, won't, idk
- Full forms: cannot, will not, I don't know

Approach:

- **Regex Replacement:** Define a dictionary of short forms and their expansions, then replace them using regular expressions.

```
python
Copycode
contractions = {"can't": "cannot", "won't": "will not", "idk": "I don't know"}
import re
```

```
pattern=re.compile(r'\b('+'|'.join(contractions.keys())+r')\b') text =  
pattern.sub(lambda x: contractions[x.group()], text)
```

Stemming:

Stemming is the process of reducing a word to its base or root form, often by removing suffixes. It focuses on achieving a normalized representation, regardless of whether the resulting form is a valid word.

How It Works:

- Reduces words to their stems by applying rule-based methods.
- Does not consider the context or meaning of the word.

Example:

"running" → "run"

"better" → "better" (unchanged, depending on the stemmer used)

Common Algorithms:

- **Porter Stemmer:** A widely used rule-based stemmer.

Python:

```
from nltk.stem import PorterStemmer
```

```
stemmer = PorterStemmer()
```

```
words = ["running", "runner", "runs"]
```

```
stemmed_words = [stemmer.stem(word) for word in words] print(stemmed_words) #
```

Output: ['run', 'runner', 'run']

- **Snowball Stemmer:** A more advanced version of the Porter Stemmer, supporting multiple languages.

Python:

```
from nltk.stem import SnowballStemmer

stemmer = SnowballStemmer("english")

words=["happily","happiness","happiest"]

stemmed_words = [stemmer.stem(word) for word in words]

print(stemmed_words)#Output:['happili','happi', 'happiest']
```

Lemmatization:

Lemmatization is the process of reducing a word to its base or root form, known as a lemma, while considering the word's context and its part of speech (POS). Unlike stemming, lemmatization ensures that the resulting word is a valid dictionary word.

How it works:

- Based on the Dictionary form
- Normalize words for consistent analysis (e.g., "running" → "run")
- It chooses the **Frequency** word

Library:

NLTK (Natural Language Toolkit)

Example Output:

```
.. Product Names (with stopwords removed, lemmatized, and Porter stemming) and Ratings:
   product_name rating
0    la facon cotton kalamkari handblock saree blou...  5.0
1      sf jean pantaloons man plain slim fit t shirt  3.6
2    lovista cotton gota patti tassel traditional p...  3.5
3      people man print regular fit t shirt  3.0
4    monte carlo grey solid cotton blend polo colla...  5.0
...
29995    indian virasat woman rayon anarkali kurta  5.0
29996    urban ranger pantaloons man slim fit t shirt  3.0
29997      peter england man regular fit t shirt  4.0
29998    pinky pari woman embroider short denim straig...  4.0
29999    gutsy man full sleeve all over print navy t s...  4.0

[30000 rows x 2 columns]
```

The system leverages various natural language processing (NLP) and machine learning techniques to recommend similar products based on their names, embeddings, and brand.

The report covers the following approaches:

1. **TF-IDF Vectorizer with Cosine Similarity**
2. **Word2Vec Embeddings**
3. **TF-IDF-Weighted Word2Vec Embeddings**
4. **Brand-Based Filtering with TF-IDF-Weighted Word2Vec**

1. TF-IDF Vectorizer with Cosine Similarity

Theoretical Foundation

TF-IDF (Term Frequency-Inverse Document Frequency) is a statistical method that transforms textual data into numerical vectors by:

- Measuring term importance within individual documents
- Reducing the weight of commonly occurring words
- Creating a unique vector representation for each product description

Mathematical Formulation:

Key Formulas:

1. Term Frequency (TF): $TF(t,d) = f_{t,d}$
2. Inverse Document Frequency (IDF): $IDF(t) = \log\left(\frac{N}{df_t}\right)$
3. TF-IDF Score: $TF-IDF(t,d) = TF(t,d) \times IDF(t)$
4. Cosine Similarity: $\cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$

Implementation Approach

- Convert product names into numerical vectors
- Calculate vector similarities using cosine similarity
- Recommend products with highest similarity scores

Advantages

- Simple and interpretable

- Computationally efficient
- Works well with short text descriptions

Limitations

- Does not capture semantic relationships
- Sensitive to exact word matches
- Struggles with synonyms and contextual variations

Together, TF-IDF identifies terms that are important for distinguishing between products. Cosine similarity, which measures the angle between two vectors, is used to compare the product names by their TF-IDF vectors. The closer the cosine value is to 1, the more similar the products are.

Key Code Snippets

```
# Initialize TfidfVectorizer
vectorizer = TfidfVectorizer()

# Fit and transform the product names
X_tfidf = vectorizer.fit_transform(data['product_name'])

# Function to recommend products
def recommend_products(product_id, num_recommendations, vectors):
    product_index = data[data['uniq_id'] == product_id].index[0]
    distances = pairwise_distances(vectors[product_index], vectors, metric='cosine')
    recommended_indices = distances.argsort()[0][1:num_recommendations + 1]
    return data.iloc[recommended_indices][['uniq_id', 'product_name', 'sales_price', 'rating', 'medium']]
```

Results:



2. Word2Vec Embeddings

Theoretical Foundation

Word2Vec generates dense vector representations that capture semantic relationships between words by analyzing their contextual usage.

Key Techniques

- **Continuous Bag of Words (CBOW):** Predicts target word from context
- **Skip-Gram:** Predicts context words from target word

Mathematical Representation

Product Embedding Calculation: $\text{ProductEmbedding} = \frac{1}{n} \sum_{i=1}^n \text{WordEmbedding}_i$

Implementation Strategy

- Train embeddings on large product description corpus
- Average word embeddings to create product-level representations
- Compute similarities between product embeddings

Advantages

- Captures semantic meaning
- Handles contextual variations
- Discovers latent relationships between words

Limitations

- Requires large training corpus
- Computationally intensive
- Treats all words equally

Overview: Word2Vec is a powerful model that generates vector representations of words based on the context in which they appear. By averaging the embeddings of words in a product name, we create a product-level embedding that captures the semantic meaning of the product name. This method is especially useful for understanding the context of words, allowing for more context-aware recommendations.

Explanation: Word2Vec works by analyzing the context of words in large corpora to generate word embeddings that capture semantic relationships between words. In the case of product names, each word's vector is averaged to generate a product-level vector that represents the entire product description.

Key Code Snippets

```

# Tokenize product names
data['product_name_tokens'] = data['product_name'].str.split()

# Train Word2Vec model
model = Word2Vec(
    sentences=data['product_name_tokens'],
    vector_size=100,
    window=5,
    min_count=1,
    workers=4,
    sg=1 # Skip-gram model
)

# Compute product embeddings
def get_product_embeddings(tokens, model):
    vectors = [model.wv[word] for word in tokens if word in model.wv]
    return sum(vectors) / len(vectors) if vectors else None

data['embedding'] = data['product_name_tokens'].apply(lambda tokens: get_product_embeddings(tokens, model))

```

Results:



3. TF-IDF-

WeightedWord2VecEmbeddingsHybridApp

roach

Combines strengths of TF-IDF and Word2Vec by weighting word embeddings based on their importance.

Mathematical Formula

$$\text{WeightedEmbedding} = \sum_{i=1}^n (\text{TF-IDF}_i \times \text{WordEmbedding}_i) \quad \text{WeightedEmbedding} = \frac{\sum_{i=1}^n (\text{TF-IDF}_i \times \text{WordEmbedding}_i)}{n}$$

Implementation Methodology

- Calculate TF-IDF scores for each word
- Weight Word2Vec embeddings using TF-IDF scores
- Create product embeddings with weighted averaging

Advantages

- Balances semantic meaning with term importance

- More nuanced recommendations
- Reduces impact of less significant words

Limitations

- More complex implementation
- Requires careful parameter tuning

Overview: This approach combines the strengths of TF-IDF and Word2Vec by weighting the embeddings of words based on their TF-IDF scores. Words that are more important to the product description (as determined by TF-IDF) will have a greater influence on the resulting product embedding.

Explanation: By multiplying each word's Word2Vec embedding by its corresponding TF-IDF score, this method ensures that the more important terms have a higher impact on the product's final embedding. This approach combines the semantic power of Word2Vec with the importance weighting provided by TF-IDF, creating more accurate and meaningful embeddings for product recommendations.

Key Code Snippets

```
# Train a TF-IDF vectorizer
vectorizer = TfidfVectorizer(tokenizer=lambda x: x, lowercase=False)
tfidf_matrix = vectorizer.fit_transform(data['product_name_tokens'])

# Compute TF-IDF-weighted Word2Vec embeddings
def get_tfidf_weighted_embeddings(tokens, model, tfidf_matrix, idx):
    vectors = []
    tfidf_scores = tfidf_matrix[idx].toarray()[0]
    for word in tokens:
        if word in model.wv and word in vectorizer.vocabulary_:
            weight = tfidf_scores[vectorizer.vocabulary_[word]]
            vectors.append(model.wv[word] * weight)
    return sum(vectors) / len(vectors) if vectors else None

data['embedding'] = [
    get_tfidf_weighted_embeddings(tokens, model, tfidf_matrix, idx)
    for idx, tokens in enumerate(data['product_name_tokens'])
]
```

Results:

Query Product



Rec 1



Rec 2



Rec 3



Rec 4



Rec 5



4. Brand-Based Filtering with TF-IDF-Weighted Word2Vec

Core Concept

Filters and recommends products within the same brand ecosystem using advanced embedding techniques.

Implementation Strategy

- Segment products by brand
- Apply TF-IDF-Weighted Word2Vec within brand subset
- Recommends similar products from same brand

Advantages

- Maintains brand consistency
- Provides more targeted recommendations
- Useful for brand-specific marketing strategies

Limitations

- Reduces recommendation diversity
- May miss cross-brand similarities

a Overview

This method refines recommendations by considering only products from the same brand. TF-IDF-Weighted Word2Vec embeddings are used to compute similarities within the brand.

Key Code Snippets:

```
# Function to recommend products from the same brand
def recommend_products_by_brand(product_id, num_recommendations):
    selected_product = data[data['uniq_id'] == product_id]
    selected_brand = selected_product['brand'].values[0]
    same_brand_products = data[data['brand'] == selected_brand]

    if len(same_brand_products) <= num_recommendations:
        return same_brand_products

    selected_embedding = selected_product['embedding'].values[0].reshape(1, -1)
    embeddings = same_brand_products['embedding'].tolist()
    similarities = cosine_similarity(selected_embedding, embeddings).flatten()
    indices = similarities.argsort()[-num_recommendations:][::-1]
    return same_brand_products.iloc[indices]
```

Results:

Query Product



Rec 1



Rec 2



Rec 3



Rec 4



Rec 5



ComparativeAnalysis:

TF-IDF:

Query Image



Recommendation 1



Recommendation 2



Recommendation 3



Recommendation 4



Recommendation 5



Word2Vec:

Query Product



Rec 1



Rec 2



Rec 3



Rec 4



Rec 5



TF-IDF-WeightedWord2Vec:

Query Product



Rec 1



Rec 2



Rec 3



Rec 4



Rec 5



Brand-BasedFiltering alongwithtfidfand word2vec:

Query Product



Rec 1



Rec 2



Rec 3



Rec 4



Rec 5



Technique	Semantic Understanding	Computational Complexity	Recommendation Precision
TF-IDF	Low	Low	Moderate
Word2Vec	High	High	Good
TF-IDF-Weighted Word2Vec	VeryHigh	VeryHigh	Excellent
Brand-Based Filtering	VeryHigh	moderate	Brand-Specific

Detailed Comparative Analysis of Recommendation Models:

Model	Visual Observations (Image-Specific)	Strengths	Weaknesses
TF-IDF Vectorizer	<p>Highly inaccurate recommendations.</p> <ul style="list-style-type: none"> - Includes products with vastly different patterns, styles, and even genders. - No alignment to the plaid check pattern or color scheme of the query product. - Example: Dresses or shirts with solid colors and styles are recommended, missing the plaid pattern entirely. 	<p>Easy to implement and computationally efficient.</p> <p>Uses term frequency to match basic product metadata (e.g., tags).</p>	<p>Fails to understand visual/contextual features like patterns, colors, and fabric style.</p> <p>Recommends items outside the intended product category or gender.</p>
Word2Vec	<p>Recommendations are closer in style to the query product.</p>	<p>Excels at capturing semantic similarity</p>	<p>Fails to consistently capture specific</p>

	<ul style="list-style-type: none"> - Focuses on solid or darker colors. - Misses the plaid pattern or the red/white/blue color scheme of the query. - Example: Dark solid shirts and light-colored plain shirts are included, which lack plaid alignment. 	<p>between products. Avoids unrelated recommendations.</p>	<p>features like color combinations or intricate plaid patterns. May prioritize overall similarity over exact details.</p>
TF-IDF-Weighted Word2Vec	<p>Shows much better contextual understanding.</p> <ul style="list-style-type: none"> - Matches plaid pattern but not exact red/white/blue color combination. - Example: A green plaid shirt appears, matching the pattern style but deviating from the query's color scheme. 	<p>Combines global contextual relevance with local keyword importance, leading to more accurate recommendations. Balances style and relevance.</p>	<p>Computationally more intensive. Still lacks perfect color matching, important for aesthetic-focused users.</p>
Brand-Based Filtering	<p>Visually aligned with the brand's design language.</p> <ul style="list-style-type: none"> - Focuses on design and stylistic consistency. - May ignore plaid patterns or specific color preferences if absent in the brand catalog. - Example: 	<p>Ensures consistency with the brand identity, appealing to brand-loyal customers.</p>	<p>Limited by the diversity of the brand catalog. May ignore stylistic elements outside the brand's focus.</p>

	<p>Recommends shirts with solid colors, stripes, or partial plaid, focusing on brand style over query-product features.</p>		
--	---	--	--

Collaborative Filtering

Collaborative Filtering (CF) is a popular recommendation system technique that makes predictions about a user's preferences based on their previous interactions and the preferences of other users. The fundamental assumption of collaborative filtering is that users who have agreed on past preferences are likely to agree on future preferences. This approach is widely used in applications like e-commerce, streaming services, and social media.

Types of Collaborative Filtering

1. User-Based Collaborative Filtering

- This method finds similarities between users based on their behavior or preferences.
- For example, if User A and User B have rated several items similarly, User A's preferences can predict User B's preferences for items they haven't rated yet.
- **Steps:**
 1. Calculate the similarity between users using metrics like cosine similarity, Pearson correlation, or Jaccard index.
 2. Use these similarities to generate recommendations.

2. Item-Based Collaborative Filtering

- This method focuses on similarities between items instead of users.
- For example, if two items are frequently rated or bought together, one item can recommend the other.
- **Steps:**
 1. Calculate similarity between items based on user interaction patterns.
 2. Recommend items similar to those the user has already interacted with.

Key Techniques

1. Memory-Based Collaborative Filtering

- Uses statistical methods directly on the data.
- Simpler and easier to implement.
- Works well for small datasets but suffers from scalability issues as the dataset grows.

2. Model-Based Collaborative Filtering

- Employs machine learning models like matrix factorization (e.g., Singular Value Decomposition, Alternating Least Squares) to discover latent factors that influence user-item interactions.
- Scalable and effective for large datasets.

Advantages of Collaborative Filtering

1. **Personalized Recommendations:** Generate tailored suggestions based on user preferences and behaviors.
2. **No Domain Knowledge Required:** CF relies solely on user-item interaction data without needing domain-specific information about items.
3. **Dynamic Adaptation:** Updates recommendations as new data becomes available.

Challenges

1. **Cold-Start Problem:** Difficult to recommend for new users or new items due to a lack of interaction data.
2. **Data Sparsity:** In systems with large datasets, most users interact with a small subset of items, leading to sparse user-item matrices.
3. **Scalability:** Computational cost increases with the number of users and items in the system.

Applications

- **E-commerce:** Suggesting products based on purchase history.
- **Streaming Services:** Recommending movies or songs based on viewing or listening history.
- **Social Media:** Proposing friends, groups, or posts based on interactions.

```

from tensorflow.keras.applications import VGG16
from tensorflow.keras.preprocessing import image
from tensorflow.keras.applications.vgg16 import preprocess_input
from tensorflow.keras.layers import GlobalMaxPooling2D
from tensorflow import keras
import numpy as np
import pandas as pd
import os

# Image processing and model setup
img_width, img_height, chnl = 224, 224, 3 # Set image dimensions
vgg = VGG16(include_top=False, weights='imagenet', input_shape=(img_width, img_height, chnl))
vgg.trainable = False
model1 = keras.Sequential([vgg, GlobalMaxPooling2D()]) # Use keras from tensorflow
model1.summary()

# Function to create image path
def img_path(img_name):
    return os.path.join("/content/ai_stylist_data/images", img_name)

# Function to get embeddings from the model
def model_predict(model, img_name):
    img = image.load_img(img_path(img_name), target_size=(img_width, img_height)) # Load image
    x = image.img_to_array(img) # Convert image to array
    x = np.expand_dims(x, axis=0) # Expand dimensions to match model input
    x = preprocess_input(x) # Preprocess image input
    return model.predict(x).reshape(-1) # Predict and reshape output

# Path to the embeddings file
embeddings_file_path = '/content/ai_stylist_data/embeddings_subset.csv'
# Check if the embeddings file already exists
if not os.path.exists(embeddings_file_path):
    print("Embeddings file not found, generating embeddings...")

```

```

# Generate embeddings for the subset of images
df_embedding_subset = df_copy_subset['imagePath'].apply(lambda x: model_predict(model1, x)) # Apply model to each image
# Convert the embeddings into a DataFrame
df_embedding_subset = df_embedding_subset.apply(pd.Series)
# Add embeddings to the copied subset DataFrame
df_copy_subset = pd.concat([df_copy_subset, df_embedding_subset], axis=1)
# Save the resulting DataFrame with embeddings for the first 2,000 images
df_copy_subset.to_csv(embeddings_file_path, index=False)
# Display the first few rows of the resulting DataFrame
print(df_copy_subset.head())
else:
    print("Embeddings file already exists, skipping the embedding generation.")

```

Output:

Downloading data from https://storage.googleapis.com/tensorflow/keras-applications/vgg16/vgg16_weights_tf_dim_ordering_58889256/58889256 0s 0us/step
Model: "sequential"

Layer (type)	Output Shape	Param #
vgg16 (Functional)	(None, 7, 7, 512)	14,714,688
global_max_pooling2d (GlobalMaxPooling2D)	(None, 512)	0

Total params: 14,714,688 (56.13 MB)
Trainable params: 0 (0.00 B)
Non-trainable params: 14,714,688 (56.13 MB)
Embeddings file not found, generating embeddings...
1/1 ----- 3s 1s/step
1/1 ----- 0s 17ms/step
1/1 ----- 0s 17ms/step
1/1 ----- 0s 16ms/step
1/1 ----- 0s 17ms/step
1/1 ----- 0s 16ms/step
1/1 ----- 0s 17ms/step

```
# all columns
merged_df.head()
```

	id	gender	masterCategory	subCategory	articleType	baseColour	season	year	usage	productDisplayItem	...	102	103	104	105	106	107	108	109	110	111
0	1570	Men	Apparel	Topwear	T-shirt	Navy Blue	Fall	2011	Casual	Turtle Check Men Navy Blue Shirt	-	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
1	1638	Men	Apparel	Bottomwear	Jeans	Blue	Summer	2012	Casual	Men's Slim Fit Blue Denim Jeans	-	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
2	1632	Women	Accessories	Watches	Watches	Silver	Winter	2013	Casual	Thin Women Silver Watch	-	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
3	1737	Men	Apparel	Bottomwear	Tank Tops	Black	Fall	2011	Casual	Men's V-neck Black Tank Top	-	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
4	1770	Men	Apparel	Topwear	T-shirt	Grey	Summer	2012	Casual	Men's V-neck Grey T-shirt	-	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000

Rows = 50 columns

Overview of VGG16 Embedding-Based Model

The VGG16 embedding-based model leverages the VGG16 architecture, a popular convolutional neural network (CNN), pre-trained on the ImageNet dataset, to extract rich feature representations from images. These embeddings can serve as input for various tasks, such as content-based recommendation systems, image classification, clustering, or retrieval.

VGG16 Architecture

1. **Developed by:** Visual Geometry Group (VGG) at Oxford University.
 2. **Key Features:**
 - Consists of 16 layers (13 convolutional layers and 3 fully connected layers).
 - Uses small convolution filters (3x3) with stride 1.
 - Incorporates max pooling for down-sampling.
 - Employs a fixed input size of 224x224 pixels for images.
 3. **Pre-Trained Weights:** Trained on ImageNet, a large-scaled dataset containing millions of labeled images.
-

Embedding Extraction with VGG16

1. **Feature Extraction:**
 - Remove the fully connected layers (classification head).
 - Use the output of the last convolutional layer or the penultimate layer as the embedding.
 - These embeddings capture high-level visual features such as shapes, textures, and patterns.
 2. **Process:**
 - Load the VGG16 model pre-trained on ImageNet.
 - Preprocess images (resize to 224x224, normalize pixel values, etc.).
 - Pass each image through the model to obtain embeddings.
 3. **Storage:**
 - Save embeddings in a structured format (e.g., NumPy arrays, database, or file storage) for future use.
-

Applications

1. **Content-Based Recommendation Systems:**
 - Find visually similar items, such as clothing, artwork, or furniture.
 - Use similarity measures (e.g., cosine similarity, Euclidean distance) to compare embeddings.
2. **Image Retrieval:**

- Search for images in a database that are similar to a given query image.

- Useful in applications like stock photo libraries or visual search engines.

3. **Image Clustering:**

- Group images with similar visual features for categorization or pattern discovery.
- Applied in domains like e-commerce and medical imaging.

4. **Hybrid Recommendation Systems:**

- Combine embeddings with metadata (e.g., product descriptions, ratings) for enhanced recommendations.

Advantages of VGG16 Embeddings

1. **Pre-Trained Model:**

- Leverages the knowledge from ImageNet, reducing the need for extensive training on specific datasets.

2. **High-Level Features:**

- Captures complex visual patterns and representations from images.

3. **Versatility:**

- Suitable for various downstream tasks beyond classification.

CODE

```
# Function to recommend products based on embeddings
def recommend_products(input_img_name, top_n=5):
    # Get input image embedding
    input_embedding = model_predict(model1, input_img_name).reshape(1, -1)

    # Compute cosine similarity
    similarities = cosine_similarity(input_embedding, embeddings).flatten()

    # Get top N recommendations (excluding the input image itself if it's in the dataset)
    recommended_indices = np.argsort(similarities)[::-1][:top_n]
    recommended_scores = similarities[recommended_indices]

    # Display recommendations
    input_img_path = img_path(input_img_name)
    display_recommendations(input_img_path, recommended_indices, recommended_scores)

# Example usage
input_image_name = '30039.jpg' # Replace with your image filename
recommend_products(input_image_name, top_n=5)
```

```
# Extract embeddings as a numpy array
embeddings = df.iloc[:, -512:].values # Assuming embeddings are the last 512 columns
# Function to display images with product names and similarity scores
def display_recommendations(input_img_path, recommended_indices, scores):
    fig, axes = plt.subplots(1, len(recommended_indices) + 1, figsize=(15, 5))

    # Display input image
    img = image.load_img(input_img_path, target_size=(224, 224))
    axes[0].imshow(img)
    axes[0].set_title("Input Image")
    axes[0].axis('off')
    # Display recommended images
    for i, idx in enumerate(recommended_indices):
        rec_img_path = img_path(df.iloc[idx]['imagePath']) # Get image path
        rec_img = image.load_img(rec_img_path, target_size=(224, 224))
        axes[i + 1].imshow(rec_img)
        axes[i + 1].set_title(f"{df.iloc[idx]['productDisplayName']}\nScore: {scores[i]:.2f}")
        axes[i + 1].axis('off')

    plt.tight_layout()
    plt.show()
```

```
# Example usage
input_image_name = '30039.jpg' # Replace with your image filename
recommend_products(input_image_name, top_n=5)
```

Output:



RecommendationSystemUsingResNet50Embeddings Overview

A recommendation system leveraging ResNet50 embeddings is a hybrid approach that combines deep learning and traditional recommendation techniques. ResNet50, a popular convolutional neural network (CNN), is pre-trained on the ImageNet dataset and is widely used for feature extraction from images. By generating embeddings (numerical representations) for images, ResNet50 captures high-level visual features, which can be utilized to recommend similar items based on image content.

Workflow for the System

1. Dataset Preparation

- Collect a dataset of images associated with the items to be recommended (e.g., product images for an e-commerce platform).
- Include metadata like item IDs, names, prices, or descriptions.

2. Feature Extraction Using ResNet50

- Load the ResNet50 model pre-trained on ImageNet.
- Remove the classification head to obtain the feature extractor model.
- Pass each image through the model to extract embeddings from the penultimate layer.
- Store these embeddings in a database for further processing.

3. Similarity Computation

- Calculate the similarity between embeddings to find visually similar items.
- Use similarity metrics like:
 - **Cosine Similarity:** Measures the cosine of the angle between two embedding vectors.
 - **Euclidean Distance:** Measures the straight-line distance between two points in the embedding space.
- Create a similarity matrix where each item is compared with every other item.

4. Recommendation Generation

- For a given item, find the top N similar items based on similarity scores.
- Display these recommendations along with relevant metadata (e.g., product name, price, or rating).

5. Integration with Metadata

- Enhance recommendations by combining image-based features with other metadata, such as textual descriptions or user behavior.
- Use hybrid approaches to improve the system's performance.

Advantages of Using ResNet50

1. **Pre-Trained Expertise:** ResNet50's pre-trained weights on ImageNet provide robust feature extraction capabilities, reducing the need for extensive training.
2. **Visual Recommendations:** Useful for domains where visual appearance significantly influences user preferences (e.g., fashion, furniture, art).
3. **Scalability:** Embeddings allow for efficient similarity computations and scalable recommendation generation.

Challenges

1. **Cold Start for Non-Visual Items:** The system may not perform well for items without meaningful visual content.
2. **Computational Cost:** Feature extraction and similarity computation can be resource-intensive for large datasets.
3. **Limited Context Understanding:** Purely visual embeddings might miss contextual relevance provided by metadata or user behavior.

Applications

- **E-commerce:** Recommending visually similar products, such as clothing or accessories.
- **Art Platforms:** Suggesting artwork or designs based on visual style.
- **Furniture Stores:** Helping users find furniture pieces with similar aesthetics.

```
# Example recommendation function
def recommend_similar_images(image_id, top_n=5):
    # Find the embedding of the query image
    query_embedding = df[df['id'] == image_id]['embeddings'].values[0]

    # Compute cosine similarity
    all_embeddings = np.stack(df['embeddings'].values)
    similarity_scores = cosine_similarity([query_embedding], all_embeddings).flatten()

    # Get top N similar images
    top_indices = similarity_scores.argsort()[::-1][1:] # Exclude the query image itself
    similar_images = df.iloc[top_indices]

    # Display the images
    plt.figure(figsize=(15, 5))
    for i, (_, row) in enumerate(similar_images.iterrows()):
        img = image.load_img(row['fullImagePath'], target_size=(224, 224))
        plt.subplot(1, top_n, i + 1)
        plt.imshow(img)
        plt.axis('off')
        plt.title(f"Similarity: {similarity_scores[top_indices[i]]:.2f}")
    plt.show()
```

```
# Test the recommendation function
test_image_id = 57958 # Replace with an actual ID from the styles.csv
recommend_similar_images(test_image_id, top_n=5)
```

Output:

Similarity: 0.82



Similarity: 0.79



Similarity: 0.79



Similarity: 0.78



Similarity: 0.77

