

## Importing the required libraries

```
1 # Importing libraries for handling matrices, dataframes, and visualizations
2 import numpy as np
3 import pandas as pd
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 %matplotlib inline
7
8 #miscellaneous libraries used
9 import re
10 import itertools
11 from collections import Counter
12 from sys import maxsize
13 import math
14 import warnings
15 warnings.filterwarnings("ignore")
16
17 #importing NLP and Data Visualisation related libraries
18 import string
19 from string import ascii_letters, punctuation, digits
20 import nltk
21 nltk.download('all')
22 from nltk.corpus import stopwords
23 from nltk.stem import WordNetLemmatizer
24 from nltk.stem.porter import PorterStemmer
25 from nltk.tokenize import word_tokenize
26 from wordcloud import WordCloud, STOPWORDS
27 from textblob import TextBlob
28 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer, TfidfTransformer
29 from sklearn.model_selection import train_test_split
30 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
31 from sklearn.metrics import roc_curve, auc, roc_auc_score
```



```
[nltk_data] | Unzipping corpora/wordnet_ic.zip.
[nltk_data] | Downloading package words to /root/nltk_data...
[nltk_data] | Unzipping corpora/words.zip.
[nltk_data] | Downloading package ycoe to /root/nltk_data...
[nltk_data] | Unzipping corpora/ycoe.zip.
[nltk_data] |
[nltk_data] Done downloading collection all
```

## Loading the Data

```
1 products = pd.read_csv("/content/drive/MyDrive/UpGrad_Hackathon/train_product_data.csv")
2 products.head()
```

		uniq_id	crawl_timestamp	product_url	product_name	product_category_tree	
0	c2d766ca982eca8304150849735ffef9	2016-03-25 22:59:23 +0000	http://www.flipkart.com/alisha-solid-women-s-c...	Alisha Solid Women's Cycling Shorts	Clothing	SRTEH2FF9KEI	
1	f449ec65dcb041b6ae5e6a32717d01b	2016-03-25 22:59:23 +0000	http://www.flipkart.com/aw-bellies/p/itmeh4grg...	AW Bellies	Footwear	SHOEH4GRSUB.	
2	0973b37acd0c664e3de26e97e5571454	2016-03-25 22:59:23 +0000	http://www.flipkart.com/alisha-solid-women-s-c...	Alisha Solid Women's Cycling Shorts	Clothing	SRTEH2F6HUZI	
3	ce5a6818f7707e2cb61fcdcbba61f5ad	2016-03-25 22:59:23 +0000	http://www.flipkart.com/alisha-solid-women-s-c...	Alisha Solid Women's Cycling Shorts	Clothing	SRTEH2FVVKRE	
4	29c8d290caa451f97b1c32df64477a2c	2016-03-25 22:59:23 +0000	http://www.flipkart.com/dilli-bazaaar-bellies-...	dilli bazaaar Bellies, Corporate Casuals, Casuals	Footwear	SHOEH3DZBFRi	

Next steps:

[View recommended plots](#)[New interactive sheet](#)

```
1 products_test = pd.read_csv("/content/drive/MyDrive/UpGrad_Hackathon/test_data.csv")
2 products_test.head()
```

	uniq_id	crawl_timestamp	product_url	product_name	pid	retail_pri
0	4fb99d98225f415e7ece96938e95628f	2015-12-20 08:26:17 +0000	http://www.flipkart.com/v-v-art-brass-bracelet...	V&V ART Brass Bracelet	BBAE6NYHCDTEZJTB	470
1	4ea284c8d38b2ea97a1c2a26f34e057c	2015-12-20 08:26:17 +0000	http://www.flipkart.com/kalpaveda-copper-cuff/...	Kalpaveda Copper Copper Cuff	BBAEDFFKZJTY7SZZ	1200
2	ee6ce2c7045c54257e2a0b590e09c296	2015-12-20 08:26:17 +0000	http://www.flipkart.com/thelostpuppy-book-cove...	Thelostpuppy Book Cover for Apple iPad Air	ACCEA4DZH6M5SFVH	2199
3	e797ba3b5f2e2d1fdc520e48486ab60e	2015-12-20 08:26:17 +0000	http://www.flipkart.com/riana-copper-bangle/p/...	Riana Copper Copper Bangle	BBAEAXFQHMHF3EYZ	2499
4	f4d8d43858c8858c68d75ce07ac641c0	2015-12-20 08:26:17 +0000	http://www.flipkart.com/inox-jewelry-stainless...	Inox Jewelry Stainless Steel Cuff	BBAECH63WYDG6TE2	1629

Next steps:

[View recommended plots](#)[New interactive sheet](#)

## Data Exploration and Preparation

```
1 products.shape
```

```
➦ (14999, 15)
```

```
1 products_test.shape
```

```
➦ (2534, 14)
```

```
1 products.info()
```

```
➦ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 14999 entries, 0 to 14998
Data columns (total 15 columns):
 #   Column                                Non-Null Count  Dtype  
---  -
 0   uniq_id                              14999 non-null  object 
 1   crawl_timestamp                      14999 non-null  object 
 2   product_url                          14999 non-null  object 
 3   product_name                         14999 non-null  object 
 4   product_category_tree                14999 non-null  object 
 5   pid                                  14999 non-null  object 
 6   retail_price                         14942 non-null  float64
 7   discounted_price                     14942 non-null  float64
 8   image                                14996 non-null  object 
 9   is_FK_Advantage_product              14999 non-null  bool    
10   description                          14998 non-null  object 
11   product_rating                       14999 non-null  object 
12   overall_rating                       14999 non-null  object 
13   brand                                10289 non-null  object 
14   product_specifications                14993 non-null  object 
dtypes: bool(1), float64(2), object(12)
memory usage: 1.6+ MB
```

```
1 products_test.shape
```

```
➦ (2534, 14)
```

```
1 print(products.isna().sum())
```

```
➦ uniq_id                0
  crawl_timestamp        0
  product_url            0
  product_name           0
  product_category_tree  0
  pid                   0
  retail_price          57
  discounted_price       57
  image                 3
  is_FK_Advantage_product 0
  description            1
  product_rating         0
  overall_rating         0
  brand                 4710
  product_specifications 6
dtype: int64
```

```
1 print(products_test.isna().sum())
```

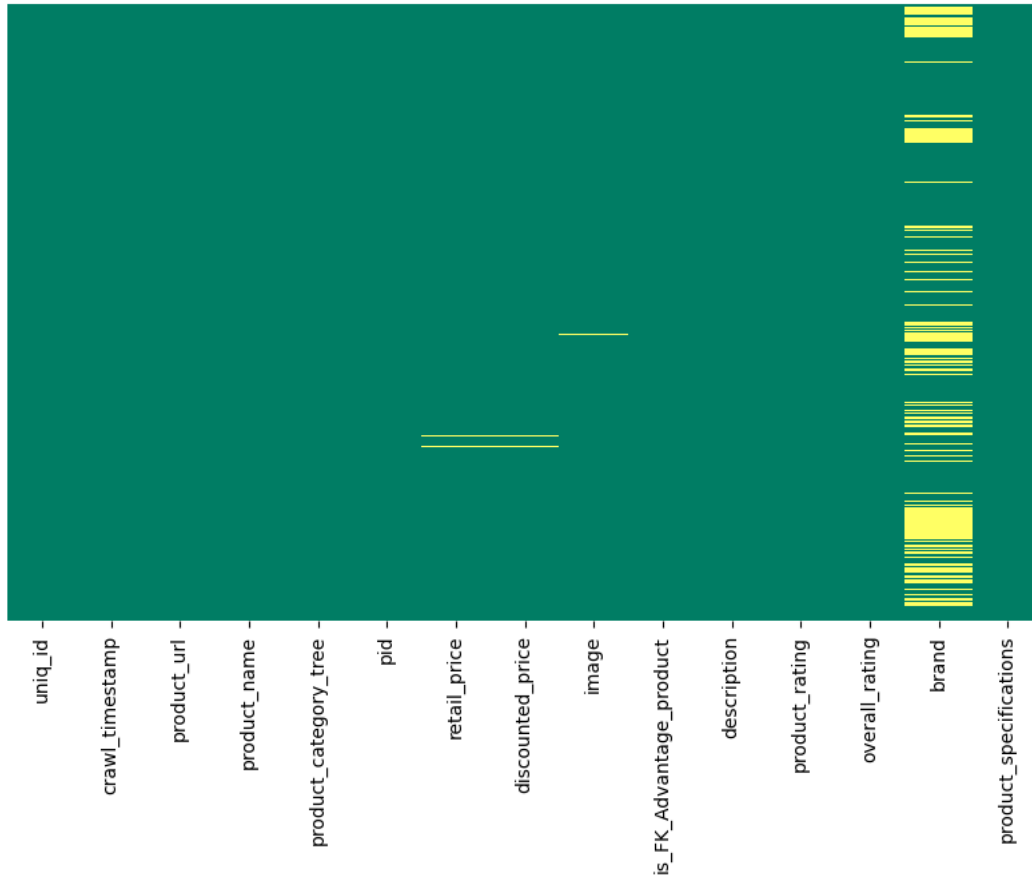
```
➦ uniq_id                0
  crawl_timestamp        0
  product_url            0
  product_name           0
  pid                   0
  retail_price           4
  discounted_price        4
  image                 0
  is_FK_Advantage_product 0
  description            0
  product_rating         0
  overall_rating         0
  brand                 522
  product_specifications 5
dtype: int64
```

```
1 #heatmap showing the distribution of all Nan's of the data
```

```
2 plt.figure(figsize=(10,6))
```

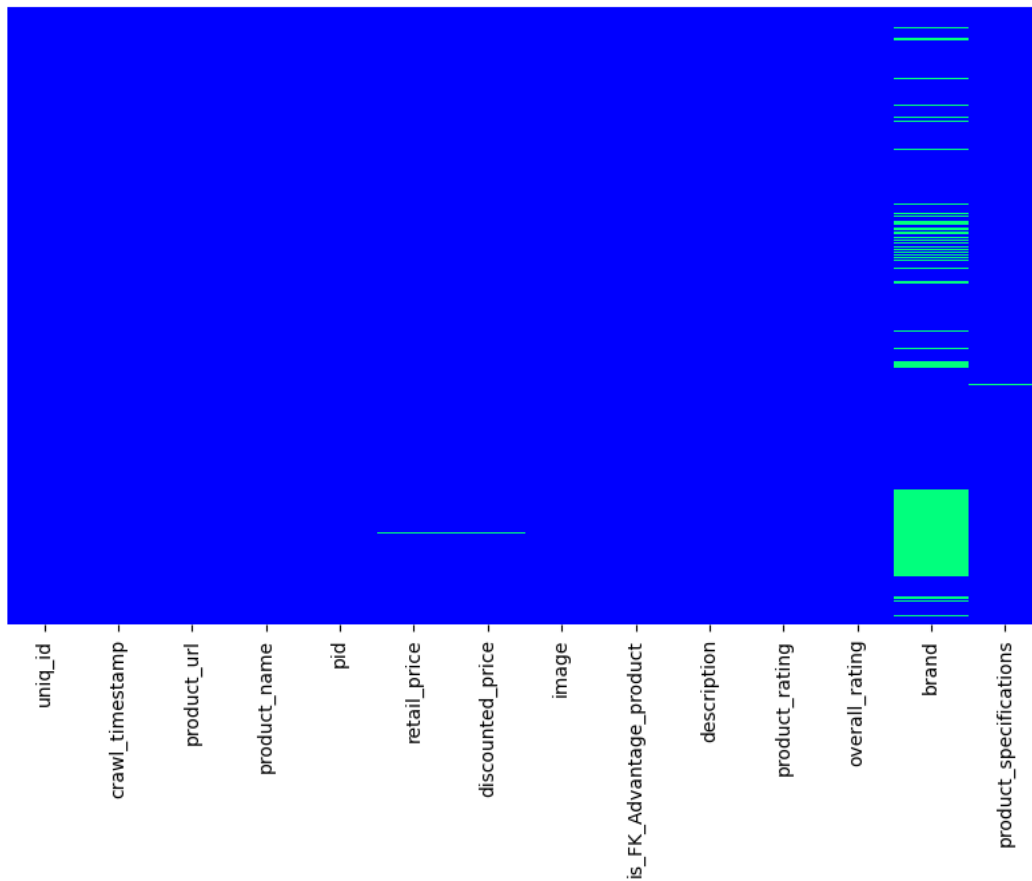
```
3 sns.heatmap(products.isnull(), cbar=False, yticklabels=False, cmap = 'summer')
```

<Axes: >



```
1 #heatmap showing the distribution of all Nan's of the test data
2 plt.figure(figsize=(10,6))
3 sns.heatmap(products_test.isnull(), cbar=False, yticklabels=False, cmap = 'winter')
```

<Axes: >



```
1 print(products.duplicated().sum())
```

↻ 0

```
1 print(products_test.duplicated().sum())
```

↻ 0

```
1 #keeping only those datapoints whose description is not NaN
2 products = products[products['description'].notna()]
```

```
1 #keeping only those datapoints in test data whose description is not NaN
2 products_test = products_test[products_test['description'].notna()]
```

```
1 #keeping only those datapoints product_category_tree is not NaN
2 products = products[products['product_category_tree'].notna()]
```

```
1 #dropping the datapoints with duplicate descriptions
2 products = products.drop_duplicates("description",keep='first', inplace=False, ignore_index=True)
3 products = products.reset_index(drop=True)
4 products
```

↻

		uniq_id	crawl_timestamp	product_url	product_name	product_category_tree	
0	c2d766ca982eca8304150849735ffef9	2016-03-25 22:59:23 +0000	http://www.flipkart.com/alisha-solid-women-s-c...	Alisha Solid Women's Cycling Shorts	Clothing	S	
1	f449ec65dcbc041b6ae5e6a32717d01b	2016-03-25 22:59:23 +0000	http://www.flipkart.com/aw-bellies/p/itmeh4grg...	AW Bellies	Footwear	S	
2	0973b37acd0c664e3de26e97e5571454	2016-03-25 22:59:23 +0000	http://www.flipkart.com/alisha-solid-women-s-c...	Alisha Solid Women's Cycling Shorts	Clothing	S	
3	ce5a6818f7707e2cb61fdcdbba61f5ad	2016-03-25 22:59:23 +0000	http://www.flipkart.com/alisha-solid-women-s-c...	Alisha Solid Women's Cycling Shorts	Clothing	S	
4	29c8d290caa451f97b1c32df64477a2c	2016-03-25 22:59:23 +0000	http://www.flipkart.com/dilli-bazaaar-bellies-...	dilli bazaaar Bellies, Corporate Casuals, Casuals	Footwear	S	
...	...	...	...	...	...	...	
13158	16de377e88660863bc028949aedb8557	2015-12-20 08:26:17 +0000	http://www.flipkart.com/thelostpuppy-book-cove...	Thelostpuppy Book Cover for Apple iPad Air 2	Mobiles & Accessories	A	
13159	9402d23592adc0795e1ca71a661c9a5f	2015-12-20 08:26:17 +0000	http://www.flipkart.com/babes-brass-cuff/p/itm...	Babes Brass Cuff	Jewellery	BE	
13160	87bcd46bb48bfc1045d7ee84aef7b7a	2015-12-20 08:26:17 +0000	http://www.flipkart.com/kenway-retail-brass-co...	Kenway Retail Brass Copper Cuff	Jewellery	BE	
13161	1336909e5468b63c9b1281350eba647d	2015-12-20 08:26:17 +0000	http://www.flipkart.com/kenway-retail-brass-co...	Kenway Retail Brass Copper Cuff	Jewellery	BE	
13162	d6eff0e0c938cc39c4451083994a2227	2015-12-20 08:26:17 +0000	http://www.flipkart.com/kenway-retail-brass-co...	Kenway Retail Brass Copper Cuff	Jewellery	BE	

13163 rows × 15 columns

Next steps:

View recommended plots

New interactive sheet

```
1 #dropping the datapoints from test data with duplicate descriptions
2 products_test = products_test.drop_duplicates("description",keep='first', inplace=False, ignore_index=True)
3 products_test = products_test.reset_index(drop=True)
4 products_test
```

		uniq_id	crawl_timestamp	product_url	product_name	pid	reta
0		4fb99d98225f415e7ece96938e95628f	2015-12-20 08:26:17 +0000	http://www.flipkart.com/v-v-art-brass-bracelet...	V&V ART Brass Bracelet	BBAE6NYHCDTEZJTB	
1		4ea284c8d38b2ea97a1c2a26f34e057c	2015-12-20 08:26:17 +0000	http://www.flipkart.com/kalpaveda-copper-cuff/...	Kalpaveda Copper Copper Cuff	BBAEDFFKZJTY7SZZ	
2		ee6ce2c7045c54257e2a0b590e09c296	2015-12-20 08:26:17 +0000	http://www.flipkart.com/thelostpuppy-book-cove...	Thelostpuppy Book Cover for Apple iPad Air	ACCEA4DZH6M5SFVH	
3		e797ba3b5f2e2d1fdc520e48486ab60e	2015-12-20 08:26:17 +0000	http://www.flipkart.com/riana-copper-bangle/p/...	Riana Copper Copper Bangle	BBAEAXFQHMF3EYZ	
4		f4d8d43858c8858c68d75ce07ac641c0	2015-12-20 08:26:17 +0000	http://www.flipkart.com/inox-jewelry-stainless...	Inox Jewelry Stainless Steel Cuff	BBAECH63WYDG6TE2	
...		...	...	...	...	...	...
2097		3ab6fae88a53a66dd7c3cbf6c9fbd3c	2015-12-01 10:15:43 +0000	http://www.flipkart.com/wallmantra-large-vinyl...	Wallmantra Large Vinyl Stickers Sticker	STIEBU65TYHDZPGX	
2098		d5a16fb788c38554feb734c15d66be6b	2015-12-01 10:15:43 +0000	http://www.flipkart.com/wallmantra-extra-large...	Wallmantra Extra Large Vinyl Stickers Sticker	STIE9F5UWNR43SZ4	
2099		43c9e22c8e9d67c0ef63f6b2d11671d7	2015-12-01 10:15:43 +0000	http://www.flipkart.com/wallmantra-extra-large...	Wallmantra Extra Large Vinyl Stickers Sticker	STIEBU65VMEQGTZY	
2100		b90031c6daba26d176aeda12eb3960d3	2015-12-01 10:15:43 +0000	http://www.flipkart.com/wallmantra-extra-large...	Wallmantra Extra Large Vinyl Stickers Sticker	STIEBU65HUNJ9GZB	
2101		d8b681d31a99ae133659764b3fc2e06a	2015-12-01 10:15:43 +0000	http://www.flipkart.com/uberlyfe-extra-large-v...	Uberlyfe Extra Large Vinyl Sticker	STIE4NXGSXG5GFR2	
2102 rows × 14 columns							

Next steps:



View recommended plots

New interactive sheet

```
1 #listing all the columns of dataset
2 print(products.columns.tolist())
```

```
['uniq_id', 'crawl_timestamp', 'product_url', 'product_name', 'product_category_tree', 'pid', 'retail_price', 'discounted_price', 'i
```

```
1 #listing all the columns of test dataset
2 print(products_test.columns.tolist())
```

```
['uniq_id', 'crawl_timestamp', 'product_url', 'product_name', 'pid', 'retail_price', 'discounted_price', 'image', 'is_FK_Advantage_f
```

```
1 #dropping the unnecessary columns
2 products = products.drop(['uniq_id',
```

```

3         'crawl_timestamp',
4         'product_url',
5         'product_name',
6         'pid',
7         'retail_price',
8         'discounted_price',
9         'image',
10        'product_rating',
11        'overall_rating',
12        'brand',
13        'is_FK_Advantage_product',
14        'product_specifications'], axis = 1)

```

```

1 #dropping the unnecessary columns of test data
2 products_test = products_test.drop(['uniq_id',
3         'crawl_timestamp',
4         'product_url',
5         'product_name',
6         'pid',
7         'retail_price',
8         'discounted_price',
9         'image',
10        'product_rating',
11        'overall_rating',
12        'brand',
13        'is_FK_Advantage_product',
14        'product_specifications'], axis = 1)

```

```
1 products.head()
```

	product_category_tree	description
0	Clothing	Key Features of Alisha Solid Women's Cycling S...
1	Footwear	Key Features of AW Bellies Sandals Wedges Heel...
2	Clothing	Key Features of Alisha Solid Women's Cycling S...
3	Clothing	Key Features of Alisha Solid Women's Cycling S...
4	Footwear	Key Features of dilli bazaaar Bellies, Corpora...

Next steps:

☒ View recommended plots

☐ New interactive sheet

```
1 products_test.head()
```

	description
0	V&V ART Brass Bracelet - Buy V&V ART Brass Bra...
1	Kalpaveda Copper Copper Cuffln ...
2	Thelostpuppy Book Cover for Apple iPad Air (Mu...
3	Riana Copper Copper Bangle - Buy Riana Copper ...
4	Inox Jewelry Stainless Steel Cuffln ...

Next steps:

☒ View recommended plots

☐ New interactive sheet

## Descriptive Analysis

A bar graph of the 30 and 20 most frequent words occurring in the train and test dataset is made. This has helped us in adding some words to our stopwords list like shipping, delivery, flipkart, etc (which are then removed) as they do not have much meaning/contribution in the prediction of product category.

```

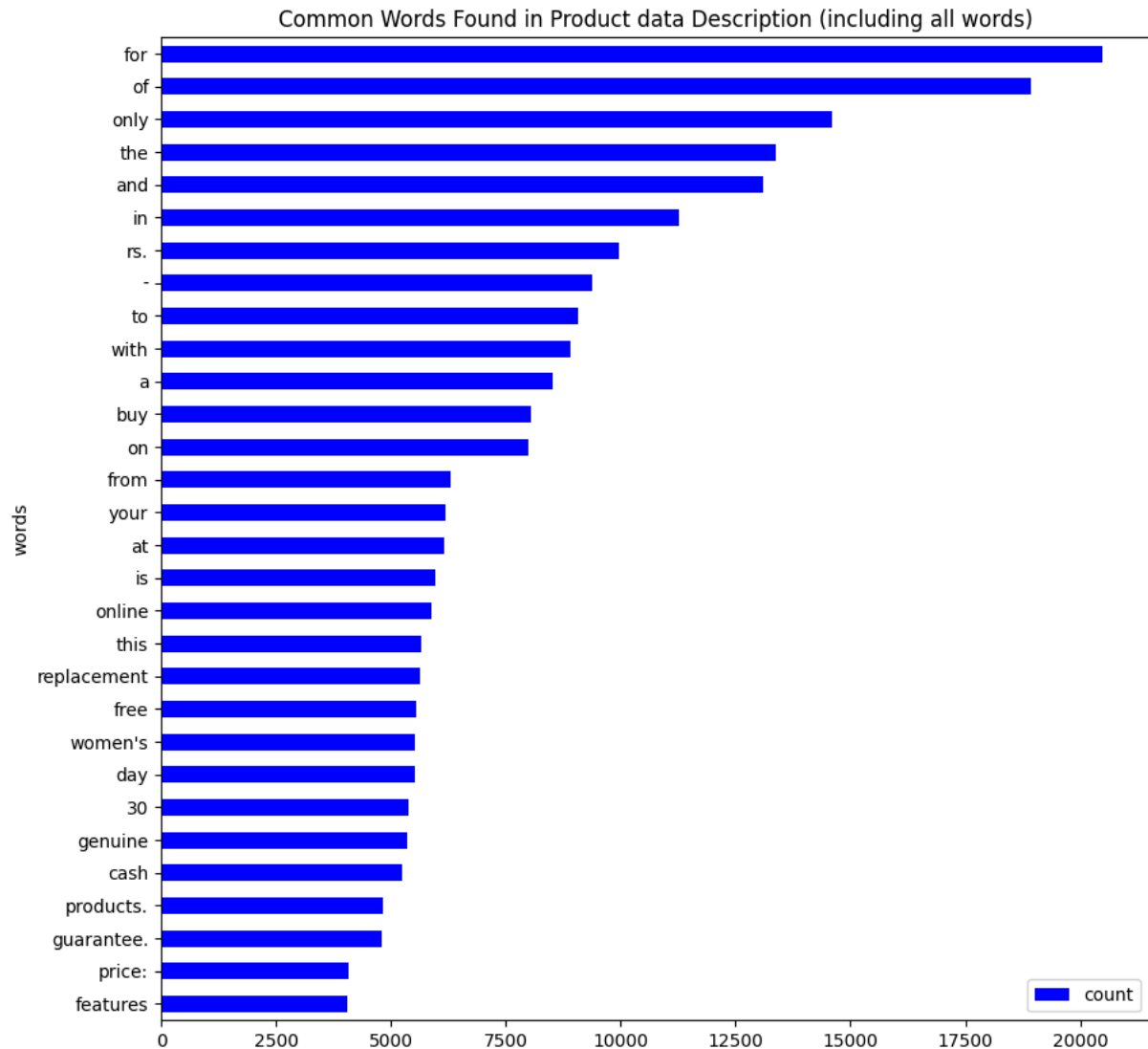
1 # for Train Data
2 def most_frequent_words(description):
3
4     for i in range(len(description)):
5         description[i] = description[i].lower().split()
6
7     all_words = list(itertools.chain(*description))
8     word_counts = Counter(all_words)
9
10    clean_description = pd.DataFrame(word_counts.most_common(30),
11                                     columns=['words', 'count'])
12

```

```

13 return clean_description
14
15 description = products["description"].tolist()
16 most_common_description = most_frequent_words(description)
17
18 fig, ax = plt.subplots(figsize=(10, 10))
19
20 # plotting the bargraph in decreasing sorted order
21 most_common_description.sort_values(by='count').plot.barh(x='words',
22                                                         y='count',
23                                                         ax=ax,
24                                                         color="blue")
25
26 ax.set_title("Common Words Found in Product data Description (including all words)")
27 plt.show()

```



```

1 # for Test Data
2 def most_frequent_words(description):
3
4     for i in range(len(description)):
5         description[i] = description[i].lower().split()
6
7     all_words = list(itertools.chain(*description))
8     word_counts = Counter(all_words)
9
10    clean_description = pd.DataFrame(word_counts.most_common(20),
11                                    columns=['words', 'count'])
12
13    return clean_description
14
15 description = products_test["description"].tolist()
16 most_common_description = most_frequent_words(description)
17
18 fig, ax = plt.subplots(figsize=(10, 10))
19
20 # plotting the bargraph in decreasing sorted order
21 most_common_description.sort_values(by='count').plot.barh(x='words',

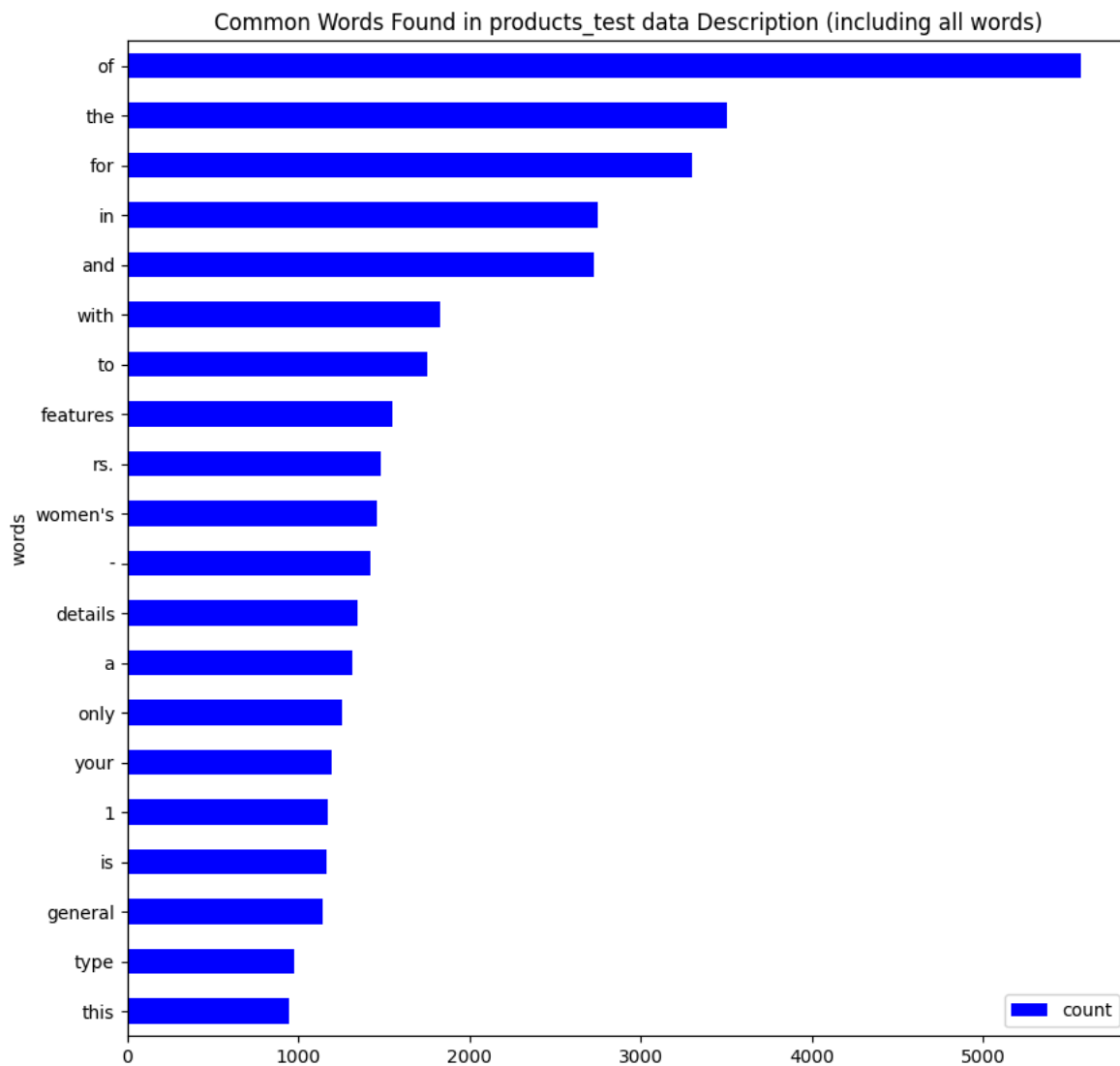
```



```

22         y='count',
23         ax=ax,
24         color="blue")
25
26 ax.set_title("Common Words Found in products_test data Description (including all words)")
27 plt.show()

```



### Exploratory Data Analysis

```

1 # Split the product_category_tree column by ">>" and extract the first level category
2 products['product_category'] = products['product_category_tree'].apply(lambda x: x.split('>>')[0].strip())
3
4 # Get the frequency count of each category
5 category_counts = products['product_category'].value_counts()
6
7 # Display the top categories
8 categories = category_counts
9 print(categories)

```



```

product_category
Clothing          4663
Jewellery         2658
Footwear          988
Automotive        935
Home Decor & Festive Needs  606
Kitchen & Dining    606
Computers         529
Watches           526
Mobiles & Accessories  501
Tools & Hardware    321
Toys & School Supplies  260
Pens & Stationery   223
Baby Care         195
Bags, Wallets & Belts  152
Name: count, dtype: int64

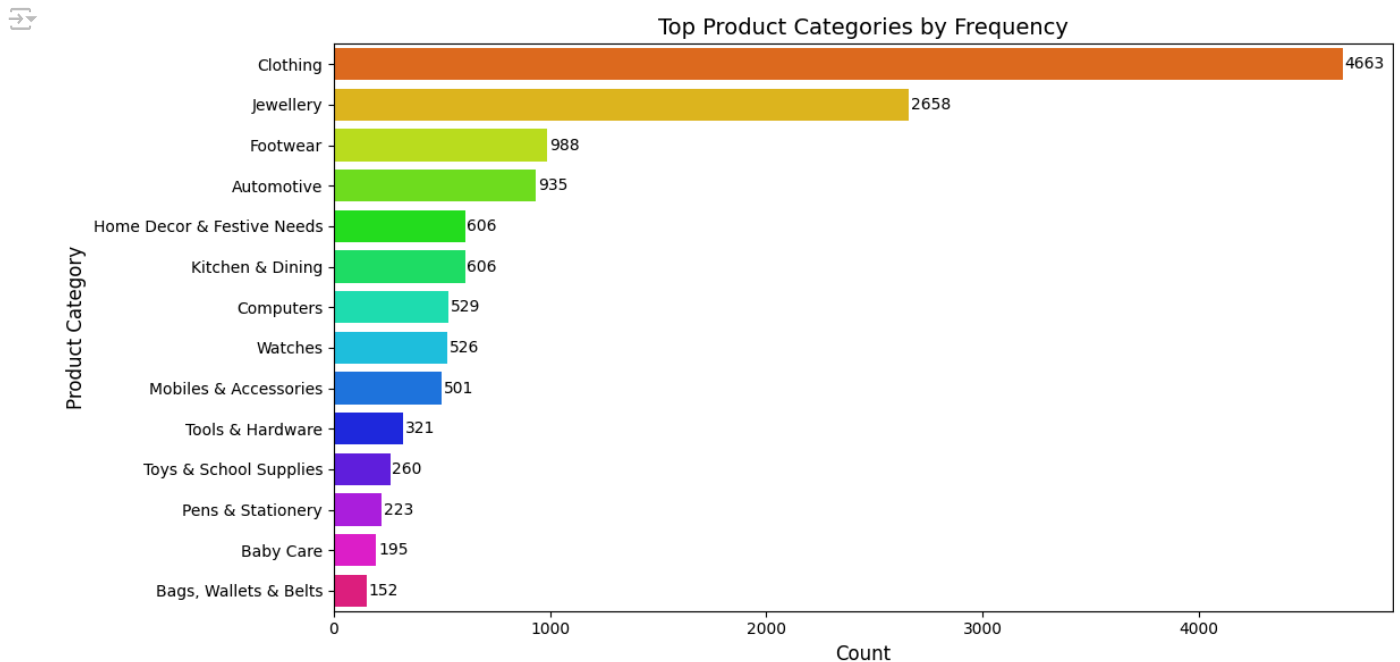
```

```

1 # Data for the bar plot
2 categories = {
3     "Clothing": 4663,
4     "Jewellery": 2658,
5     "Footwear": 988,
6     "Automotive": 935,
7     "Home Decor & Festive Needs": 606,
8     "Kitchen & Dining": 606,
9     "Computers": 529,
10    "Watches": 526,
11    "Mobiles & Accessories": 501,
12    "Tools & Hardware": 321,
13    "Toys & School Supplies": 260,
14    "Pens & Stationery": 223,
15    "Baby Care": 195,
16    "Bags, Wallets & Belts": 152,
17 }

1 # Convert to sorted lists
2 category_names = list(categories.keys())
3 category_counts = list(categories.values())
4
5 # Create a colorful bar plot
6 plt.figure(figsize=(12, 6))
7 colors = sns.color_palette("hsv", len(categories))
8 sns.barplot(x=category_counts, y=category_names, palette=colors)
9
10 # Add text annotations for counts
11 for index, value in enumerate(category_counts):
12     plt.text(value + 10, index, str(value), va='center', fontsize=10, color='black')
13
14 # Add labels and title
15 plt.xlabel("Count", fontsize=12)
16 plt.ylabel("Product Category", fontsize=12)
17 plt.title("Top Product Categories by Frequency", fontsize=14)
18 plt.tight_layout()
19
20 # Display the plot
21 plt.show()

```



### Text Length Analysis

Analysis of the length of the Product Description is done to help us get an idea about the minimum, maximum and average length of the same. This is done in order to decide whether we have to discard some datapoints having text length less than or greater to a threshold.

```

1 #finding the length of the description
2 max_desc_len = -1
3 desc_len_sum = 0

```

```

4 min_desc_len = maxsize
5
6 product_description = products["description"].tolist()
7
8 for i in range(len(product_description)):
9     try:
10         max_desc_len = max(max_desc_len, len(product_description[i]))
11         min_desc_len = min(min_desc_len, len(product_description[i]))
12         desc_len_sum += len(product_description[i])
13     except:
14         pass
15
16 print("Max description length is {}".format(max_desc_len))
17 print("Min description length is {}".format(min_desc_len))
18 print("Average description length is {}".format(desc_len_sum / len(product_description)))

```

```

↗ Max description length is 5309.
  Min description length is 74.
  Average description length is 410.81508774595454.

```

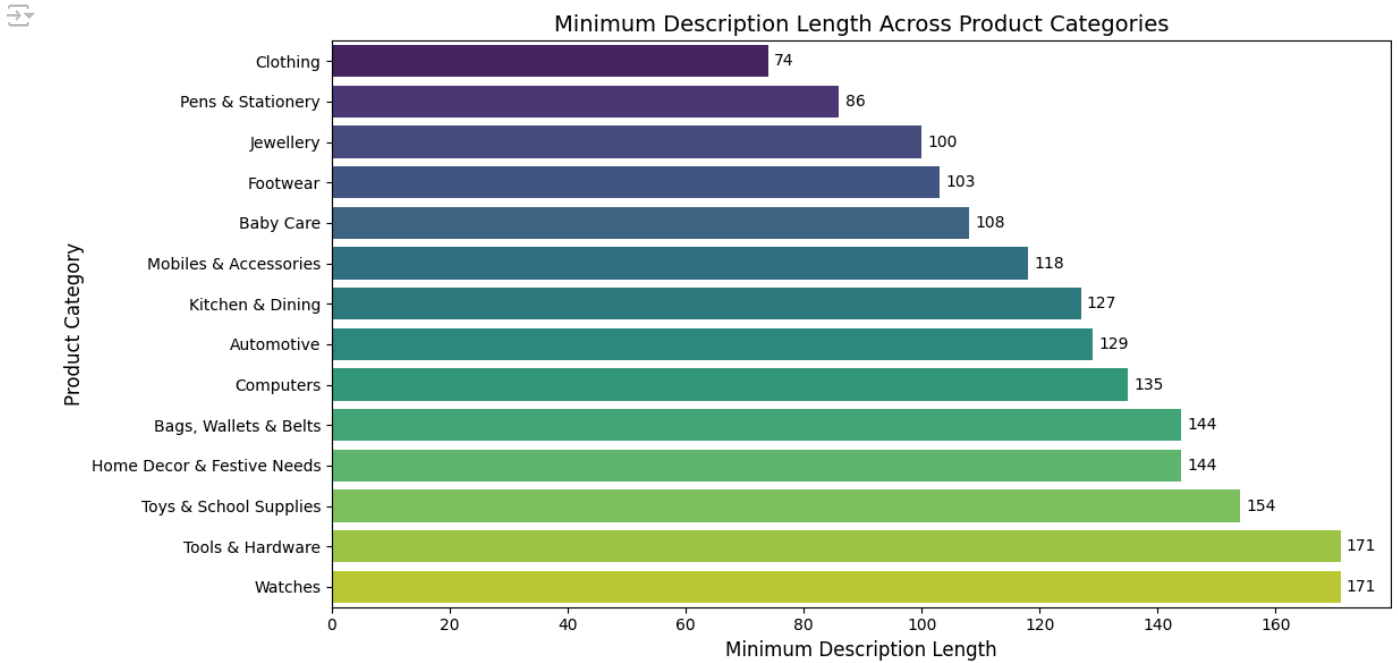
### Visualisation of the Minimum description length across all the categories

From the graph below we can see that there are discrepancies in the minimum length across all the categories. **Watches** have a minimum length almost greater than 170 while **clothing** has the minimum length.

```

1 # Calculate the length of each description
2 products['description_length'] = products['description'].apply(len)
3
4 # Get the minimum description length for each category
5 min_desc_length_per_category = products.groupby('product_category')['description_length'].min().sort_values()
6
7 # Prepare data for visualization
8 categories = min_desc_length_per_category.index
9 min_lengths = min_desc_length_per_category.values
10
11 # Create a bar plot for minimum description lengths across categories
12 plt.figure(figsize=(12, 6))
13 sns.barplot(x=min_lengths, y=categories, palette="viridis")
14
15 # Add text annotations for minimum lengths
16 for index, value in enumerate(min_lengths):
17     plt.text(value + 1, index, str(value), va='center', fontsize=10, color='black')
18
19 # Add labels and title
20 plt.xlabel("Minimum Description Length", fontsize=12)
21 plt.ylabel("Product Category", fontsize=12)
22 plt.title("Minimum Description Length Across Product Categories", fontsize=14)
23 plt.tight_layout()
24
25 # Display the plot
26 plt.show()
27

```



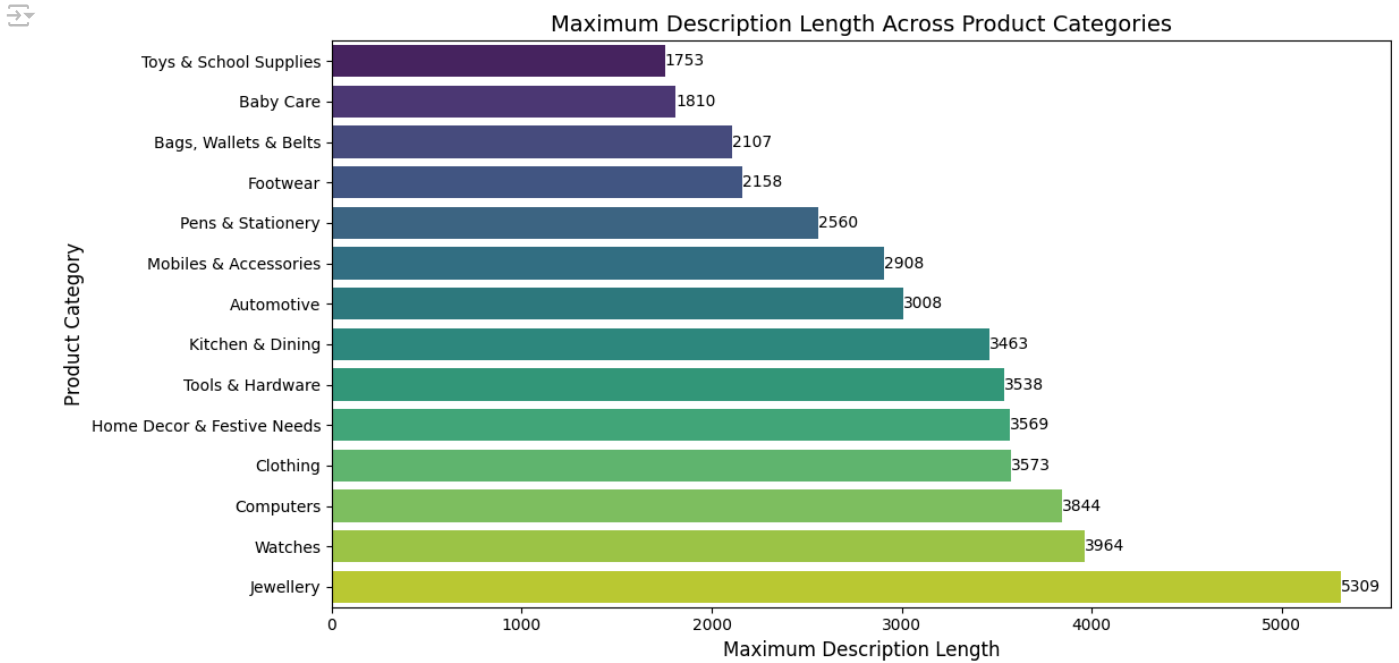
#### Visualisation of the Maximum description length across all the categories

From the graph below we can see that almost all the product description lengths are greater than 1500 with **jewellery** category having the greatest length for product description (greater than 5000)

```

1 # Calculate the length of each description
2 products['description_length'] = products['description'].apply(len)
3
4 # Get the minimum description length for each category
5 max_desc_length_per_category = products.groupby('product_category')['description_length'].max().sort_values()
6
7 # Prepare data for visualization
8 categories = max_desc_length_per_category.index
9 max_lengths = max_desc_length_per_category.values
10
11 # Create a bar plot for minimum description lengths across categories
12 plt.figure(figsize=(12, 6))
13 sns.barplot(x=max_lengths, y=categories, palette="viridis")
14
15 # Add text annotations for minimum lengths
16 for index, value in enumerate(max_lengths):
17     plt.text(value + 1, index, str(value), va='center', fontsize=10, color='black')
18
19 # Add labels and title
20 plt.xlabel("Maximum Description Length", fontsize=12)
21 plt.ylabel("Product Category", fontsize=12)
22 plt.title("Maximum Description Length Across Product Categories", fontsize=14)
23 plt.tight_layout()
24
25 # Display the plot
26 plt.show()
27

```



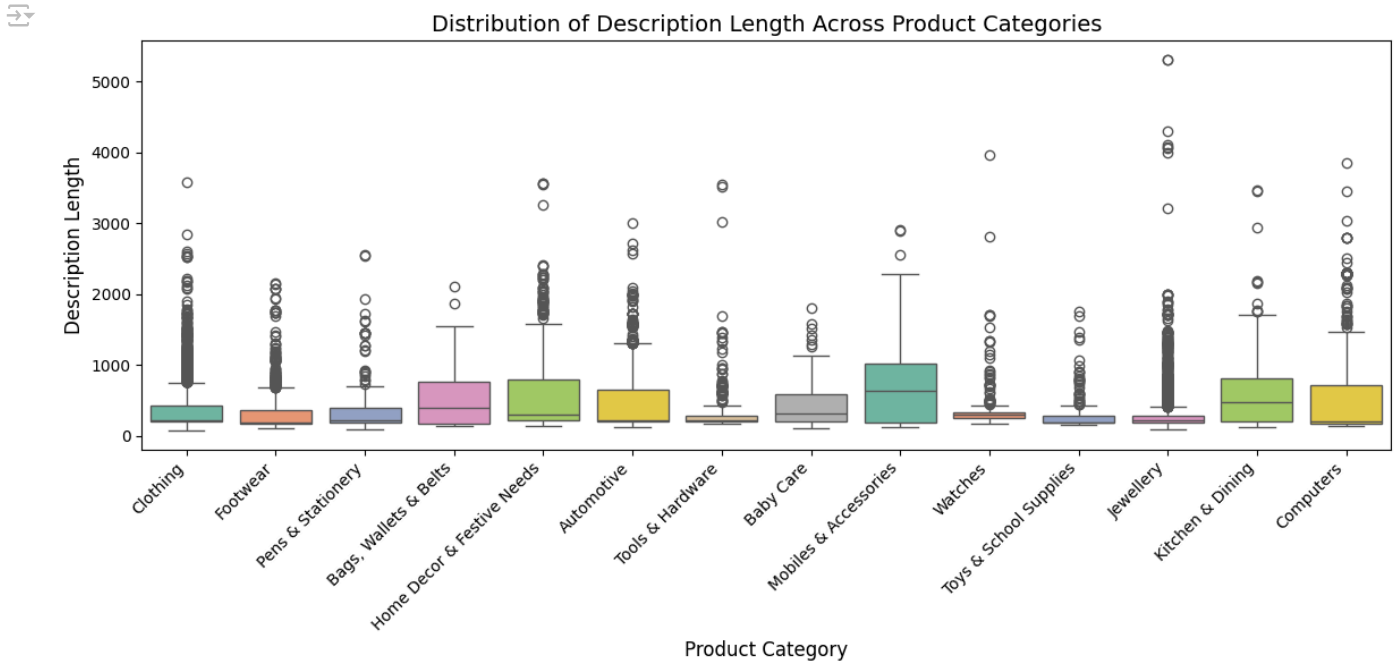
#### Boxplot of the Average description length across all the categories

From the boxplot, we can see that the average length for all the categories lies more or less around 1000 words. I decided to not go with a certain minimum/maximum words threshold to prevent loss of information.

```

1 # Calculate the description length
2 products['description_length'] = products['description'].apply(len)
3
4 # Create a boxplot
5 plt.figure(figsize=(12, 6))
6 sns.boxplot(x='product_category', y='description_length', data=products, palette="Set2")
7
8 # Rotate category labels for better readability
9 plt.xticks(rotation=45, ha='right')
10
11 # Add labels and title
12 plt.xlabel("Product Category", fontsize=12)
13 plt.ylabel("Description Length", fontsize=12)
14 plt.title("Distribution of Description Length Across Product Categories", fontsize=14)
15 plt.tight_layout()
16
17 # Display the plot
18 plt.show()

```



## Word Clouds

### Word cloud consisting of the most frequent words in the Product Description

This wordcloud shows the 200 most common words in the raw dataset that was provided. This wordcloud helped a lot to get an idea about removing words such as **[flpkart, replacement, genuine, product, shipping, cash etc]** as these words are common to the context of all the categories and will not contribute much to predicting the category of a particular product.

```

1 product_content = ""
2
3 for i in products["description"]:
4     i = str(i)
5     separate = i.split()
6     for j in range(len(separate)):
7         separate[j] = separate[j].lower()
8
9     product_content += " ".join(separate)+" "
10
11 stop_words = set(STOPWORDS)
12 final_wordcloud = WordCloud(width = 3000, height = 1600,
13                             max_words=100,
14                             background_color = 'black',
15                             stopwords = stop_words,
16                             min_font_size = 10).generate(product_content)
17
18 plt.figure(figsize = (10, 10), facecolor = None)
19 plt.title("100 frequent words in the Product Description Corpus", fontsize=20)
20 plt.imshow(final_wordcloud)
21 plt.axis("off")
22 plt.tight_layout(pad = 0)
23 plt.show()

```



```

1 # Get unique categories
2 categories = products['product_category'].unique()
3
4 # Stop words for filtering
5 stop_words = set(STOPWORDS)
6
7 # Create word clouds for each category
8 for category in categories:
9     # Filter descriptions for the current category
10    category_descriptions = products[products['product_category'] == category]['description']
11
12    # Concatenate all descriptions for the category
13    category_content = " ".join([str(desc).lower() for desc in category_descriptions])
14
15    # Generate word cloud
16    wordcloud = WordCloud(
17        width=1500,
18        height=800,
19        max_words=50,
20        background_color='black',
21        stopwords=stop_words,
22        min_font_size=10
23    ).generate(category_content)
24
25    # Display the word cloud
26    plt.figure(figsize=(5, 5), facecolor=None)
27    plt.title(f"{category} Most Frequent Words in {category} Descriptions", fontsize=16)
28    plt.imshow(wordcloud)
29    plt.axis("off")
30    plt.tight_layout(pad=0)
31
32 plt.show()

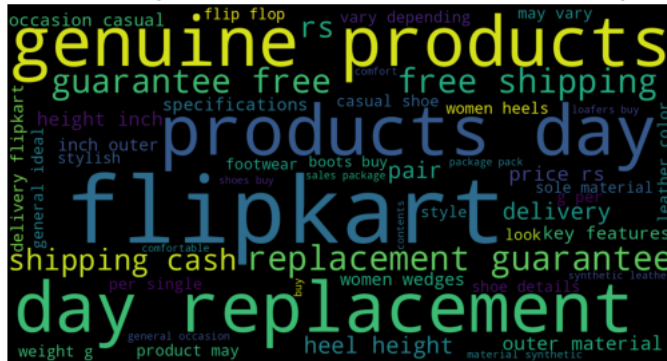
```



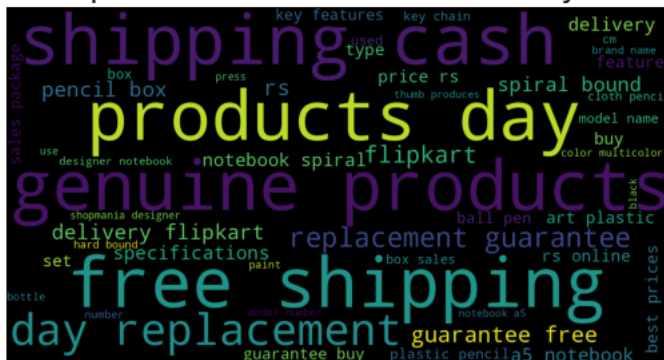
## 50 Most Frequent Words in Clothing Descriptions



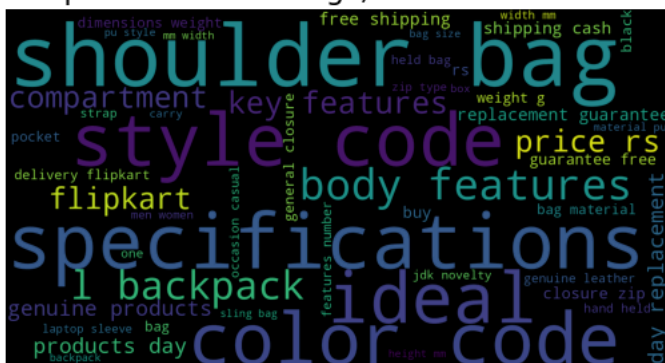
## 50 Most Frequent Words in Footwear Descriptions



## 50 Most Frequent Words in Pens & Stationery Descriptions



## 50 Most Frequent Words in Bags, Wallets & Belts Descriptions

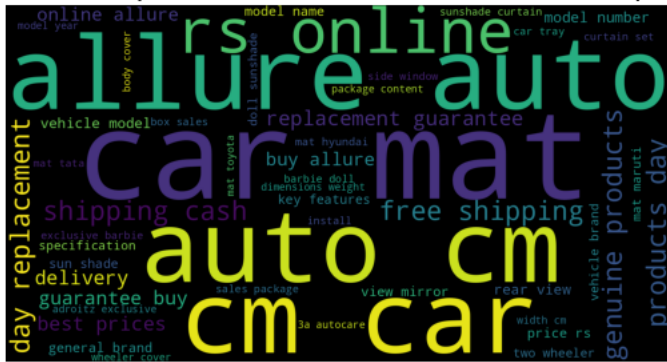


## 50 Most Frequent Words in Home Decor & Festive Needs Descriptions

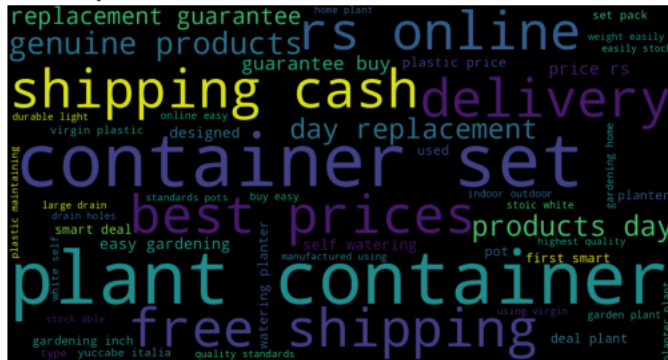




### 50 Most Frequent Words in Automotive Descriptions



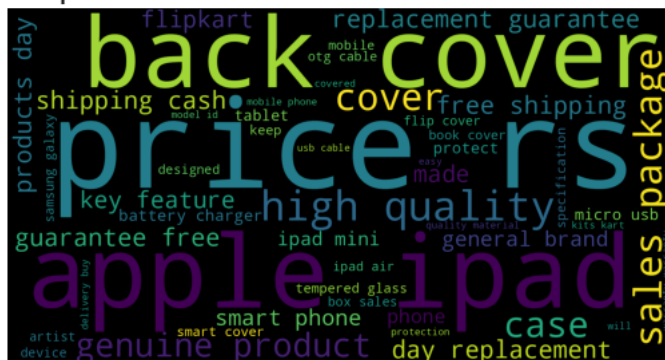
## 50 Most Frequent Words in Tools & Hardware Descriptions



## 50 Most Frequent Words in Baby Care Descriptions



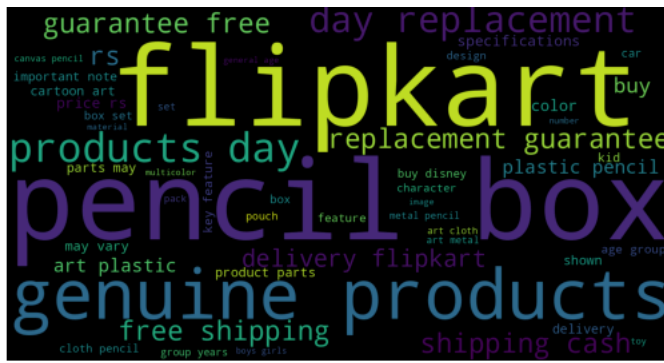
## 50 Most Frequent Words in Mobiles & Accessories Descriptions



## 50 Most Frequent Words in Watches Descriptions



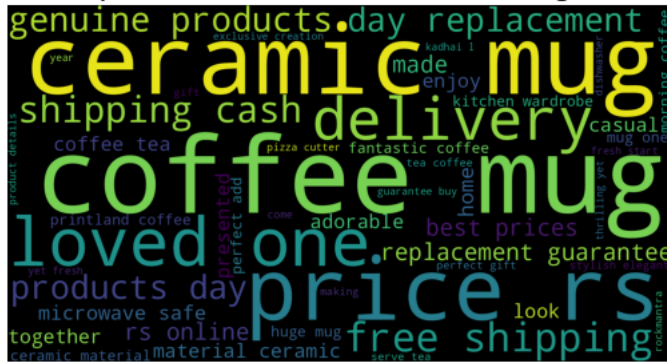
## 50 Most Frequent Words in Toys & School Supplies Descriptions



## 50 Most Frequent Words in Jewellery Descriptions



## 50 Most Frequent Words in Kitchen & Dining Descriptions



## 50 Most Frequent Words in Computers Descriptions



## ✓ Data Cleaning and Pre Processing in Train & Test Data

### Character Contraction

Character contraction is done to look at what percentage of the dataset is in English Characters. Punctuations, numbers, hyperlinks, etc will all be removed during further cleaning of the dataset.

During this analysis of the characters, we can see that there are several emoticons, letters from different languages (Chinese), etc that can be seen. These are then removed from the corpus.

### ✓ For train data

```
1 corpus_train = ' '.join(products['description']).lower()
2 characters = Counter(corpus_train)
3 sorted_characters = sorted(characters.items(), key=lambda i: i[1], reverse=True)
4
5 total=0
6 for i in ascii_letters+punctuation+digits:
7     total+=characters[i]
8
9 print("The % of data consisting of only English Characters is {}".format(100*total/len(' '.join(products['description']))))
10 print("\n")
11 print(characters)
```

↻ The % of data consisting of only English Characters is 82.31818239676973.

Counter({' ': 922665, 'e': 458959, 'a': 343552, 'o': 320427, 'r': 297788, 't': 296895, 'i': 295022, 'n': 277981, 's': 262981, 'l': 2

```
1 # getting all the words ending an apostrophe and single letter
2 contractions = Counter(re.findall("[a-z]+'[a-z]+", corpus_train))
3 apostrophe_end = sorted(contractions.items(), key=lambda i: i[1], reverse=True)
4 print("\n")
5 print(apostrophe_end)
6
7 # getting all the words starting with a single letter and an apostrophe
8 contractions=Counter(re.findall("[a-z]'[a-z]+", corpus_train))
9 apostrophe_start = sorted(contractions.items(), key=lambda i: i[1], reverse=True)
10 print("\n")
11 print(apostrophe_start)
12
13 #getting all the URLs
14 urls = re.findall('(*http[s]?://(?:[a-zA-Z]|[0-9]|[$-_@.&#]|[*\(\),]|(?:%[0-9a-fA-F][0-9a-fA-F]))*)', corpus_train)
15 print("\n")
16 print(urls)
```

↻

```
[("women's", 5572), ("men's", 2705), ("girl's", 841), ("boy's", 524), ("don't", 131), ("you're", 114), ("it's", 69), ("doesn't", 40
[("n's", 8314), ("l's", 859), ("y's", 564), ("n't", 216), ("r's", 160), ("u're", 114), ("t's", 105), ("a's", 91), ("e's", 57), ("d'r
['https://www.dropbox.com/s/xkth19lhya1jvvm/mile%201430%20black%204.jpg', 'https://www.dropbox.com/s/xkth19lhya1jvvm/mile%201430%20
```

```
1 custom_contracts = {
2     "women's" : "women",
3     "men's" : "men",
4     "girl's" : "girl",
5     "boy's" : "boy",
6     "don't" : "do not",
7     "product's" : "product",
8     "it's" : "its",
9     "bra's" : "bras",
10    "won't": "will not",
11    "doesn't" : "does not",
12    "l's" : " ",
13    "n's" : " ",
14    "y's" : " ",
15    "n't" : "not",
16    "r's" : "rs",
17    "u're" : "your",
18    "a's" : " ",
19    "e's": " "
```

In the following code snippet, the following things have been taken care of:

- Lowercasing
- Custom Contraction Mapping
- Keeping only the ascii characters in the corpus
- Removal of URLs/ Hyperlinks
- Removal of numbers and punctuations
- Custom Stopword Removal
- Lemmatization
- Removal of extra whitespaces

## Visualisation of the cleaned Dataset

40 most common words are plotted in the form of a bargraph after removal of the unnecessary data. From the bargraph, we can clearly see that these are the words which actually can help us in identifying the particular category of a product.

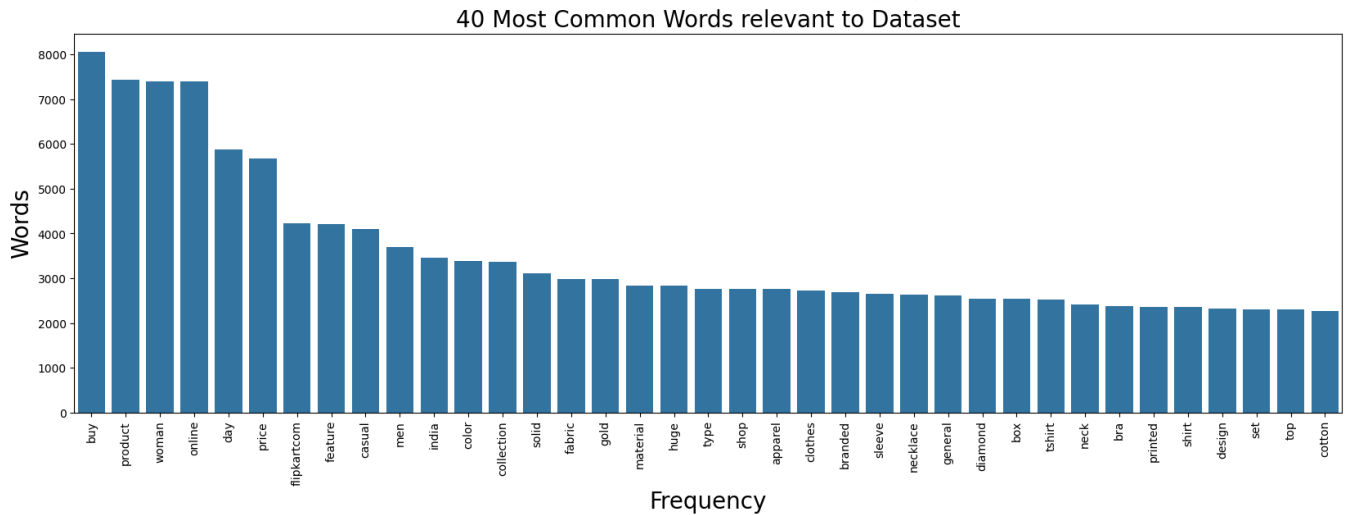
20/31

```

17 plt.xlabel("Frequency", fontsize=20)
18 plt.xticks(rotation=90)
19 sns.barplot(x=x,y=y)

```

<Axes: title={'center': '40 Most Common Words relevant to Dataset'}, xlabel='Frequency', ylabel='Words'>



```

1 temp_cleaned_description = [" ".join(sentence) for sentence in cleaned_description]
2 products["cleaned_desc"] = temp_cleaned_description
3 products.head(10)

```

	product_category_tree	description	product_category	description_length	cleaned_desc
0	Clothing	Key Features of Alisha Solid Women's Cycling S...	Clothing	410	key feature alisha solid woman cycling short c...
1	Footwear	Key Features of AW Bellies Sandals Wedges Heel...	Footwear	650	key feature belly sandal wedge heel casualsaw ...
2	Clothing	Key Features of Alisha Solid Women's Cycling S...	Clothing	403	key feature alisha solid woman cycling short c...
3	Clothing	Key Features of Alisha Solid Women's Cycling S...	Clothing	416	key feature alisha solid woman cycling short c...
4	Footwear	Key Features of dilli bazaar Bellies, Corpora...	Footwear	428	key feature dilli bazaar belly corporate casu...
5	Clothing	Key Features of Alisha Solid Women's Cycling S...	Clothing	419	key feature alisha solid woman cycling short c...
6	Footwear	Key Features of Ladela Bellies	Footwear	358	key feature ladela belly brand

Next steps:

[View recommended plots](#)

[New interactive sheet](#)

## For Test data

```

1 #for test data
2 corpus_test = ' '.join(products_test['description']).lower()
3 characters = Counter(corpus_test)
4 sorted(characters.items(), key=lambda i: i[1], reverse=True)
5
6 total=0
7 for i in ascii_letters+punctuation+digits:
8     total+=characters[i]
9
10 print("The % of data consisting of only English Characters is {}".format(100*total/len(' '.join(products_test['description']))))
11 print("\n")
12 print(characters)

```

The % of data consisting of only English Characters is 82.85368314127203.

Counter({' ': 190236, 'e': 101870, 'a': 74963, 't': 66102, 'i': 65585, 'o': 65292, 'r': 62096, 's': 60424, 'n': 57436, 'l': 49120,

```

1 # getting all the words ending an apostrophe and single letter
2 contractions = Counter(re.findall("[a-z]+'[a-z]+", corpus_test))
3 apostrophe_end = sorted(contractions.items(), key=lambda i: i[1], reverse=True)
4 print("\n")
5 print(apostrophe_end)
6
7 # getting all the words starting with a single letter and an apostrophe
8 contractions=Counter(re.findall("[a-z]'[a-z]+", corpus_test))
9 apostrophe_start = sorted(contractions.items(), key=lambda i: i[1], reverse=True)
10 print("\n")
11 print(apostrophe_start)
12
13 #getting all the URLs
14 urls = re.findall('(http[s]?://(?:[a-zA-Z]|[0-9]|[$-_@.&#]|[*\(\),]|(?:%[0-9a-fA-F][0-9a-fA-F])))*', corpus_test)
15 print("\n")
16 print(urls)

```



```

[('women's', 1488), ('men's', 419), ('boy's', 248), ('girl's', 173), ('don't', 24), ('l'appel', 16), ('it's', 12), ('doesn't', 11),

[('n's', 1907), ('y's', 254), ('l's', 173), ('n't', 41), ('a's', 23), ('l'appel', 16), ('r's', 15), ('t's', 15), ('d's', 11), ('e's

['http://www.ninecolours.com/suits/net-digital-print-anarkali-suit-in-blue-colour-sm0561055#sthash.k8ea4wgy.dpuf,specifications']

```

```

1 custom_contracts = {
2     "women's" : "women",
3     "men's" : "men",
4     "girl's" : "girl",
5     "boy's" : "boy",
6     "don't" : "do not",
7     "product's" : "product",
8     "it's" : "its",
9     "bra's" : "bras",
10    "won't": "will not",
11    "doesn't" : "does not",
12    "l's" : " ",
13    "n's" : " ",
14    "y's" : " ",
15    "n't" : "not",
16    "r's" : "rs",
17    "u're" : "your",
18    "a's" : " ",
19    "e's": " "
20 }
21
22 custom_stopwords = [w for w in set(stopwords.words("english"))]
23 custom_stopwords += list(punctuation)
24 stopwords_dataset = ["replacement","shipping","delivery","cash", "rs", "flipkart", "genuine", "details", "guarantee","free", "genuine
25 custom_stopwords.extend(stopwords_dataset)
26 wordnet_lemmatizer = WordNetLemmatizer()

```

```

1 def clean(text):
2
3     for i in range(len(text)):
4         text[i] = text[i].lower()
5         text[i] = text[i].replace("\n", " ")
6         for keys,values in custom_contracts.items():
7             text[i] = text[i].replace(keys,values)
8         text[i] = re.sub("[a-z]'[a-z]+", " ", text[i])
9         #removing the extra whitespaces
10        text[i] = re.sub(' +', ' ', text[i])
11        #keeping only the ascii characters -> handles emoticons, letters from other languages, etc
12        text[i] = re.sub(r'[^\x00-\x7F]+', ' ', text[i])
13        #removing the urls
14        text[i] = re.sub('(http[s]?://(?:[a-zA-Z]|[0-9]|[$-_@.&#]|[*\(\),]|(?:%[0-9a-fA-F][0-9a-fA-F])))*', ' ', text[i])
15        text[i] = ''.join([j for j in text[i] if not j.isdigit()])
16        text[i] = text[i].split()
17        text[i] = ' '.join([word for word in text[i] if word not in custom_stopwords])
18        #removing the punctuations
19        text[i] = re.sub(r'[^\w\s]', '', text[i])
20        #lemmatization
21        text[i] = [wordnet_lemmatizer.lemmatize(w) for w in word_tokenize(text[i])]
22        #removing the words which have a length less than 3
23        text[i] = [word for word in text[i] if len(word)>=3]
24

```

```

25 return text
26
27 raw_description = products_test["description"].tolist()
28 cleaned_description = clean(raw_description)
29 print(cleaned_description[:5])

```

```

[[ 'art', 'brass', 'bracelet', 'buy', 'art', 'brass', 'bracelet', 'flipkartcom', 'product', 'day', 'guarantee', 'shipping', 'delivery'

```

```

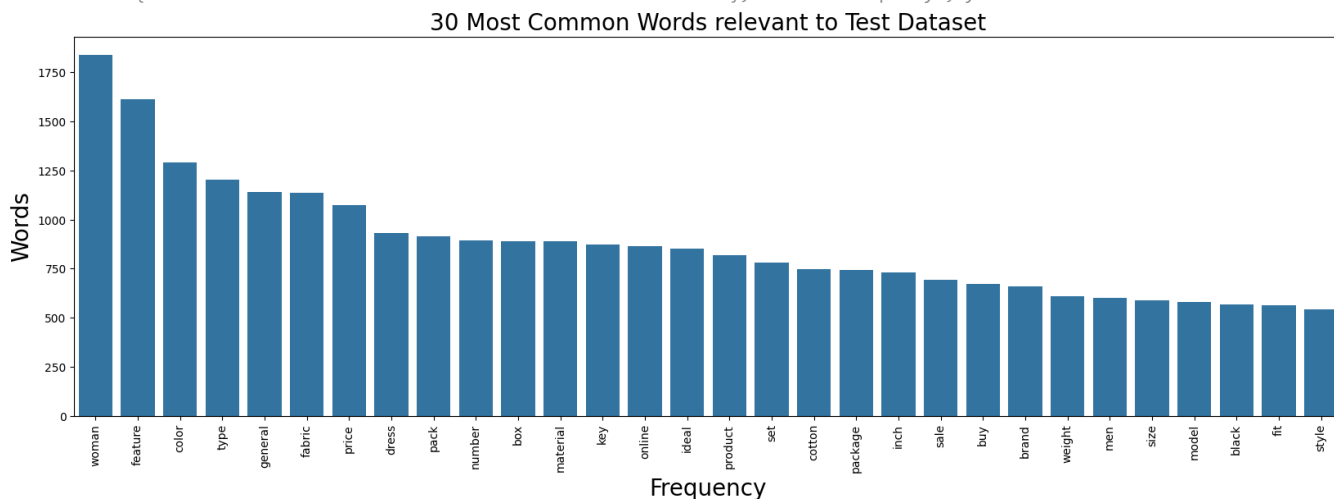
1 corpus_test = []
2 for text in cleaned_description:
3     for word in text:
4         corpus_test.append(word)
5
6 counter = Counter(corpus_test)
7 most=counter.most_common()
8 x, y = [], []
9 for word,count in most[:30]:
10     if (word not in custom_stopwords):
11         x.append(word)
12         y.append(count)
13
14 plt.figure(figsize=(20,6))
15 plt.title("30 Most Common Words relevant to Test Dataset", fontsize=20)
16 plt.ylabel("Words", fontsize=20)
17 plt.xlabel("Frequency", fontsize=20)
18 plt.xticks(rotation=90)
19 sns.barplot(x=x,y=y)

```

```

<Axes: title={'center': '30 Most Common Words relevant to Test Dataset'}, xlabel='Frequency', ylabel='Words'>

```



```

1 temp_cleaned_description = [" ".join(sentence) for sentence in cleaned_description]
2 products_test["cleaned_desc"] = temp_cleaned_description
3 products_test.head(10)

```

	description	cleaned_desc
0	V&V ART Brass Bracelet - Buy V&V ART Brass Bra...	art brass bracelet buy art brass bracelet flip...
1	Kalpaveda Copper Copper Cuffn ...	kalpaveda copper copper cuff price attract emi...
2	Thelostpuppy Book Cover for Apple iPad Air (Mu...	thelostpuppy book cover apple ipad air multico...
3	Riana Copper Copper Bangle - Buy Riana Copper ...	riana copper copper bangle buy riana copper co...
4	Inox Jewelry Stainless Steel Cuffn ...	inox jewelry stainless steel cuff price inox j...
5	Thelostpuppy Book Cover for Apple iPad Air 2 (...)	thelostpuppy book cover apple ipad air multico...
6	Ridhi Sidhi Collection Brass Bangle Set (Pack ...	ridhi sidhi collection brass bangle set pack p...
7	Theskinmantra Sleeve for All versions of Apple...	theskinmantra sleeve version apple ipad multic...
8	TheLostPuppy Back Cover for Apple iPad Air 2 (...)	thelostpuppy back cover apple ipad air multico...
9	Intex Happy Animal Chair Assortment Inflatable...	intex happy animal chair assortment inflatable...

Next steps:

View recommended plots

New interactive sheet

1 products.head()

	product_category_tree	description	product_category	description_length	cleaned_desc	
0	Clothing	Key Features of Alisha Solid Women's Cycling S...	Clothing	410	key feature alisha solid woman cycling short c...	
1	Footwear	Key Features of AW Bellies Sandals Wedges Heel...	Footwear	650	key feature belly sandal wedge heel casualsaw ...	
2	Clothing	Key Features of Alisha Solid Women's Cycling S...	Clothing	403	key feature alisha solid woman cycling short c...	
3	Clothing	Key Features of Alisha Solid Women's Cycling S...	Clothing	403	key feature alisha solid woman cycling short c...	

-----Key Features of Alisha Solid-----key feature alisha solid woman-----

Next steps:

View recommended plots

New interactive sheet

1 products\_test.head()

	description	cleaned_desc	
0	V&V ART Brass Bracelet - Buy V&V ART Brass Bra...	art brass bracelet buy art brass bracelet flip...	
1	Kalpaveda Copper Copper Cuffln ...	kalpaveda copper copper cuff price attract emi...	
2	Thelostpuppy Book Cover for Apple iPad Air (Mu...	thelostpuppy book cover apple ipad air multico...	
3	Riana Copper Copper Bangle - Buy Riana Copper ...	riana copper copper bangle buy riana copper co...	
4	Inox Jewelry Stainless Steel Cuffln ...	inox jewelry stainless steel cuff price inox j...	

-----

**Encoding of the Product Classes** In order to plot the ROC Curves and find the AUC score, there was a need to have a proper encoding for the 14 primary categories (in both directions). Hence, two of the following dictionaries are created to create a mapping.

```
1 #helper dictionaries created which are later used to manipulate the testing output into suitable form before plotting the ROC Curves
2
3 category_mapping = {
4     0 : "Clothing",
5     1 : "jewellery",
6     2 : "Footwear",
7     3 : "Automotive",
8     4 : "Home Decor & Festive Needs",
9     5 : "Kitchen & Dining",
10    6 : "Computers",
11    7 : "Watches",
12    8 : "Mobiles & Accessories",
13    9 : "Tools & Hardware",
14   10 : "Toys & School Supplies",
15   11 : "Pens & Stationery",
16   12 : "Baby Care",
17   13 : "Bags, Wallets & Belts"
18 }
19 reverse_category_mapping = {
20     "Clothing": 0,
21     "Jewellery": 1,
22     "Footwear": 2,
23     "Automotive": 3,
24     "Home Decor & Festive Needs": 4,
25     "Kitchen & Dining": 5,
26     "Computers": 6,
27     "Watches": 7,
28     "Mobiles & Accessories": 8,
29     "Tools & Hardware": 9,
30     "Toys & School Supplies": 10,
31     "Pens & Stationery": 11,
32     "Baby Care": 12,
33     "Bags, Wallets & Belts": 13,
34 }
```

1) Logistic Regression (Binary Classification Method)

```
1 from sklearn.linear_model import LogisticRegression
2
3 def logistic_regression(train_data, test_data):
4     # Define the feature (description) and target (product_category) columns
```



```
5 x_train = train_data['cleaned_desc'] # cleaned description in train data
6 y_train = train_data['product_category'] # target labels in train data
7
8 x_test = test_data['cleaned_desc'] # cleaned description in test data
9
10 # Splitting the dataset into training and test parts (for internal validation)
11 x_train, x_val, y_train, y_val = train_test_split(x_train, y_train, test_size=0.2, random_state=42)
12
13 # Bag of words implementation
14 cv = CountVectorizer()
15 x_train = cv.fit_transform(x_train).toarray()
16 x_val = cv.transform(x_val).toarray()
17
18 # TF-IDF implementation
19 vector = TfidfTransformer()
20 x_train = vector.fit_transform(x_train).toarray()
21 x_val = vector.transform(x_val).toarray()
22
23 # Initialize the logistic regression model
24 lr_model = LogisticRegression(max_iter=1000)
25
26 # Fitting the model with training data
27 lr_model.fit(x_train, y_train)
28
29 # Predict on the validation set
30 lr_predict = lr_model.predict(x_val)
31
32 # Evaluation metrics for the validation data
33 print("Validation Accuracy: ", accuracy_score(y_val, lr_predict))
34 print("\n***** CONFUSION MATRIX *****")
35 print(confusion_matrix(y_val, lr_predict))
36 print("\n***** CLASSIFICATION REPORT *****")
37 print(classification_report(y_val, lr_predict))
38
39 # Now apply the model to the test data (products_test)
40 x_test = cv.transform(x_test) # Apply CountVectorizer on test data
41 x_test = vector.transform(x_test) # Apply TF-IDF on test data
42
43 # Predict the product category for test data
44 lr_test_predict = lr_model.predict(x_test)
45 lr_test_pred_prob = lr_model.predict_proba(x_test)
46
47 # Returning the predictions for the test data
48 return lr_test_predict, lr_test_pred_prob
49
50 # Call the function with the train and test data
51 lr_predictions, lr_pred_probabilities = logistic_regression(products, products_test)
52
53 # You can now add these predictions to the test dataframe (products_test)
54 products_test['predicted_category'] = lr_predictions
55
56 # Show the test data with the predicted categories
57 products_test.head(10)
```

Validation Accuracy: 0.9699962020508925

```
***** CONFUSION MATRIX *****
[[196  0  0  1  0  0  0  0  1  1  0  0  0  0]
 [ 0 32  1 12  0  1  4  2  0  0  0  0  0  0]
 [ 0  0 26  4  0  2  0  0  0  0  0  0  0  0]
 [ 0  1  0 90 8  0  0  0  2  0  0  0  0  0]
 [ 0  0  1  0 89 0  0  0  0  3  0  0  0  0]
 [ 0  0  0  3  0 17 6  0  0  0  0  0  0  0]
 [ 0  0  0  0  0  0 13 2  0  2  0  0  0  0]
 [ 0  0  0  1  0  0  0 55 7  0  0  0  0  0]
 [ 0  0  0  1  0  2  0  0 11 5  0  0  0  0]
 [ 1  0  0  1  2  0  0  0  0 10 5  0  0  0]
 [ 0  0  0  1  0  0  3  2  0  0 26  0  7  0]
 [ 0  0  0  4  0  0  1  0  0  0  0 66  0  0]
 [ 3  0  0  2  1  0  0  0  0  0  2  0 42  0]
 [ 0  0  0  4  0  0  0  0  0  0  0  0  0 84]]
```

```
***** CLASSIFICATION REPORT *****
              precision    recall  f1-score   support

   Automotive              0.98       0.98       0.98        199
    Baby Care              0.97       0.62       0.75         52
  Bags, Wallets & Belts      0.93       0.81       0.87         32
      Clothing              0.96       1.00       0.98        911
      Computers              0.97       0.96       0.96         93
      Footwear              0.97       0.98       0.98        179
Home Decor & Festive Needs    0.94       0.99       0.96        134
      Jewellery              0.99       1.00       0.99        558
    Kitchen & Dining          0.97       0.97       0.97        118
  Mobiles & Accessories      0.96       0.96       0.96        109
    Pens & Stationery         0.93       0.67       0.78         39
    Tools & Hardware          1.00       0.93       0.96         71
  Toys & School Supplies      0.86       0.84       0.85         50
        Watches              1.00       0.95       0.98         88

   accuracy                   0.97       2633
  macro avg              0.96       0.90       0.93       2633
 weighted avg              0.97       0.97       0.97       2633
```

	description	cleaned_desc	predicted_category
0	V&V ART Brass Bracelet - Buy V&V ART Brass Bra...	art brass bracelet buy art brass bracelet flip...	Jewellery
1	Kalpaveda Copper Copper Cuffln ...	kalpaveda copper copper cuff price attract emi...	Jewellery
2	Thelostpuppy Book Cover for Apple iPad Air (Mu...	thelostpuppy book cover apple ipad air multico...	Mobiles & Accessories
3	Riana Copper Copper Bangle - Buy Riana Copper ...	riana copper copper bangle buy riana copper co...	Jewellery
4	Inox Jewelry Stainless Steel Cuffln ...	inox jewelry stainless steel cuff price inox j...	Jewellery
5	Thelostpuppy Book Cover for Apple iPad Air 2 (...)	thelostpuppy book cover apple ipad air multico...	Mobiles & Accessories
6	Ridhi Sidhi Collection Brass Bangle Set (Pack ...	ridhi sidhi collection brass bangle set pack p...	Jewellery
7	Theskinmantra Sleeve for All versions of Apple...	theskinmantra sleeve version apple ipad multico...	Mobiles & Accessories
8	TheLostPuppy Back Cover for Apple iPad Air 2 (...)	thelostpuppy back cover apple ipad air multico...	Mobiles & Accessories
9	Intex Happy Animal Chair Assortment Inflatable...	intex happy animal chair assortment inflatable...	Clothing

Next steps:

 View recommended plots

 New interactive sheet

2) Logistic Regression (Multiclass Classification Method)

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.model_selection import train_test_split
3 from sklearn.feature_extraction.text import CountVectorizer, TfidfTransformer
4 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
5 from sklearn.preprocessing import LabelBinarizer
6 from sklearn.metrics import roc_curve, auc
7 import matplotlib.pyplot as plt
8
9 def logistic_regression_multiclass(train_data, test_data):
10     # Define the feature (description) and target (product_category) columns
11     x_train = train_data['cleaned_desc'] # cleaned description in train data
12     y_train = train_data['product_category'] # target labels in train data
13
14     x_test = test_data['cleaned_desc'] # cleaned description in test data
15
16     # Splitting the dataset into training and validation parts (for internal validation)
17     x_train, x_val, y_train, y_val = train_test_split(x_train, y_train, test_size=0.2, random_state=42)
18
19     # Apply CountVectorizer on the training data (fit_transform) and validation/test data (transform)
20     cv = CountVectorizer()
```

```
21 x_train = cv.fit_transform(x_train)
22 x_val = cv.transform(x_val)
23 x_test_transformed = cv.transform(x_test)
24
25 # Apply TF-IDF on the transformed data (training, validation, and test data)
26 vector = TfidfTransformer()
27 x_train = vector.fit_transform(x_train)
28 x_val = vector.transform(x_val)
29 x_test_transformed = vector.transform(x_test_transformed)
30
31 # Initialize the logistic regression model for multiclass classification
32 lr_model = LogisticRegression(solver='lbfgs', multi_class='multinomial', max_iter=1000)
33
34 # Fitting the model with training data
35 lr_model.fit(x_train, y_train)
36
37 # Predict on the validation set
38 lr_predict = lr_model.predict(x_val)
39
40 # Evaluation metrics for the validation data
41 print("Validation Accuracy: ", accuracy_score(y_val, lr_predict))
42 print("\n***** CONFUSION MATRIX *****")
43 print(confusion_matrix(y_val, lr_predict))
44 print("\n***** CLASSIFICATION REPORT *****")
45 print(classification_report(y_val, lr_predict))
46
47 # Predict the product category for test data (products_test)
48 lr_test_predict = lr_model.predict(x_test_transformed)
49 lr_test_pred_prob = lr_model.predict_proba(x_test_transformed)
50
51 # Returning the predictions for the test data
52 return lr_test_predict, lr_test_pred_prob
53
54 # Call the function with the train and test data
55 lr_predictions, lr_pred_probabilities = logistic_regression_multiclass(products, products_test)
56
57 # You can now add these predictions to the test dataframe (products_test)
58 products_test['predicted_category'] = lr_predictions
59
60 # Show the test data with the predicted categories
61 products_test.head(10)
62
```

Validation Accuracy: 0.9699962020508925

```
***** CONFUSION MATRIX *****
[[196  0  0  1  0  0  0  0  1  1  0  0  0  0]
 [ 0 32  1 12  0  1  4  2  0  0  0  0  0  0]
 [ 0  0 26  4  0  2  0  0  0  0  0  0  0  0]
 [ 0  1  0 90 8  0  0  0  2  0  0  0  0  0]
 [ 0  0  1  0 89 0  0  0  0  3  0  0  0  0]
 [ 0  0  0  3  0 17 6  0  0  0  0  0  0  0]
 [ 0  0  0  0  0  0 13 2  0  2  0  0  0  0]
 [ 0  0  0  1  0  0  0 55 7  0  0  0  0  0]
 [ 0  0  0  1  0  2  0  0 11 5  0  0  0  0]
 [ 1  0  0  1  2  0  0  0  0 10 5  0  0  0]
 [ 0  0  0  1  0  0  3  2  0  0 26  0  7  0]
 [ 0  0  0  4  0  0  1  0  0  0  0 66  0  0]
 [ 3  0  0  2  1  0  0  0  0  0  2  0 42  0]
 [ 0  0  0  4  0  0  0  0  0  0  0  0  0 84]]
```

```
***** CLASSIFICATION REPORT *****
              precision    recall  f1-score   support

Automotive           0.98       0.98       0.98       199
Baby Care            0.97       0.62       0.75        52
Bags, Wallets & Belts 0.93       0.81       0.87        32
Clothing              0.96       1.00       0.98       911
Computers             0.97       0.96       0.96        93
Footwear              0.97       0.98       0.98       179
Home Decor & Festive Needs 0.94       0.99       0.96       134
Jewellery             0.99       1.00       0.99       558
Kitchen & Dining       0.97       0.97       0.97       118
Mobiles & Accessories 0.96       0.96       0.96       109
Pens & Stationery       0.93       0.67       0.78        39
Tools & Hardware       1.00       0.93       0.96        71
Toys & School Supplies 0.86       0.84       0.85        50
Watches               1.00       0.95       0.98        88

accuracy              0.97       0.97       0.97      2633
macro avg             0.96       0.90       0.93      2633
weighted avg          0.97       0.97       0.97      2633
```

	description	cleaned_desc	predicted_category
0	V&V ART Brass Bracelet - Buy V&V ART Brass Bra...	art brass bracelet buy art brass bracelet flip...	Jewellery
1	Kalpaveda Copper Copper Cuffln ...	kalpaveda copper copper cuff price attract emi...	Jewellery
2	Thelostpuppy Book Cover for Apple iPad Air (Mu...	thelostpuppy book cover apple ipad air multico...	Mobiles & Accessories
3	Riana Copper Copper Bangle - Buy Riana Copper ...	riana copper copper bangle buy riana copper co...	Jewellery
4	Inox Jewelry Stainless Steel Cuffln ...	inox jewelry stainless steel cuff price inox j...	Jewellery
5	Thelostpuppy Book Cover for Apple iPad Air 2 (...)	thelostpuppy book cover apple ipad air multico...	Mobiles & Accessories
6	Ridhi Sidhi Collection Brass Bangle Set (Pack ...	ridhi sidhi collection brass bangle set pack p...	Jewellery
7	Theskinmantra Sleeve for All versions of Apple...	theskinmantra sleeve version apple ipad multico...	Mobiles & Accessories
8	TheLostPuppy Back Cover for Apple iPad Air 2 (...)	thelostpuppy back cover apple ipad air multico...	Mobiles & Accessories
9	Intex Happy Animal Chair Assortment Inflatable...	intex happy animal chair assortment inflatable...	Clothing

Next steps:

[View recommended plots](#)

[New interactive sheet](#)

### 3) Multinomial Naive Bayes Classifier

```
1 from sklearn.naive_bayes import MultinomialNB
2 from sklearn.model_selection import train_test_split
3 from sklearn.feature_extraction.text import CountVectorizer, TfidfTransformer
4 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
5 from sklearn.preprocessing import LabelBinarizer
6
7 def naive_bayes(train_data, test_data):
8     # Define the feature (description) and target (product_category) columns
9     x_train = train_data['cleaned_desc'] # cleaned description in train data
10    y_train = train_data['product_category'] # target labels in train data
11
12    x_test = test_data['cleaned_desc'] # cleaned description in test data
13
14    # Splitting the dataset into training and validation parts (for internal validation)
15    x_train, x_val, y_train, y_val = train_test_split(x_train, y_train, test_size=0.2, random_state=42)
16
17    # Apply CountVectorizer on the training data (fit_transform) and validation/test data (transform)
18    cv = CountVectorizer()
19    x_train = cv.fit_transform(x_train)
20    x_val = cv.transform(x_val)
```

```
21 x_test_transformed = cv.transform(x_test)
22
23 # Apply TF-IDF on the transformed data (training, validation, and test data)
24 vector = TfidfTransformer()
25 x_train = vector.fit_transform(x_train)
26 x_val = vector.transform(x_val)
27 x_test_transformed = vector.transform(x_test_transformed)
28
29 # Initialize the Multinomial Naive Bayes model for multiclass classification
30 nb_model = MultinomialNB()
31
32 # Fitting the model with training data
33 nb_model.fit(x_train, y_train)
34
35 # Predict on the validation set
36 nb_predict = nb_model.predict(x_val)
37
38 # Evaluation metrics for the validation data
39 print("Validation Accuracy: ", accuracy_score(y_val, nb_predict))
40 print("\n***** CONFUSION MATRIX *****")
41 print(confusion_matrix(y_val, nb_predict))
42 print("\n***** CLASSIFICATION REPORT *****")
43 print(classification_report(y_val, nb_predict))
44
45 # Predict the product category for test data (products_test)
46 nb_test_predict = nb_model.predict(x_test_transformed)
47 nb_test_pred_prob = nb_model.predict_proba(x_test_transformed)
48
49 # Returning the predictions for the test data
50 return nb_test_predict, nb_test_pred_prob
51
52 # Call the function with the train and test data
53 nb_predictions, nb_pred_probabilities = naive_bayes(products, products_test)
54
55 # You can now add these predictions to the test dataframe (products_test)
56 products_test['predicted_category'] = nb_predictions
57
58 # Show the test data with the predicted categories
59 products_test.head(10)
```

Validation Accuracy: 0.8997341435624763

```
***** CONFUSION MATRIX *****
[[194  0  0  3  0  0  0  1  0  1  0  0  0  0]
 [ 1  0  0 39  0  0  2 10  0  0  0  0  0  0]
 [ 0  0  0 21  0  0  0 11  0  0  0  0  0  0]
 [ 0  0  0 90 9  0  0  2  0  0  0  0  0  0]
 [ 0  0  0  0 69  0  0 21  0  3  0  0  0  0]
 [ 0  0  0  7  0 167  0  5  0  0  0  0  0  0]
 [ 1  0  0  5  0  1 103 24  0  0  0  0  0  0]
 [ 0  0  0  1  0  0  0 557  0  0  0  0  0  0]
 [ 2  0  0  5  0  0  0  9 102  0  0  0  0  0]
 [ 2  0  0  2  0  0  0  6  0 99  0  0  0  0]
 [ 1  0  0  4  0  0  0 25  0  1  8  0  0  0]
 [ 0  0  0  1  1  0  1  5  0  0  0 63  0  0]
 [ 3  0  0  5  0  0  0 28  0  0  0  0 14  0]
 [ 0  0  0  4  0  0  0  0  0  0  0  0  0 84]]
```

```
***** CLASSIFICATION REPORT *****
              precision    recall  f1-score   support

   Automotive              0.95       0.97       0.96        199
   Baby Care              0.00       0.00       0.00         52
  Bags, Wallets & Belts    0.00       0.00       0.00         32
      Clothing            0.90       1.00       0.95       911
      Computers            0.99       0.74       0.85         93
      Footwear            0.99       0.93       0.96       179
Home Decor & Festive Needs  0.97       0.77       0.86       134
      Jewellery           0.79       1.00       0.88       558
   Kitchen & Dining       1.00       0.86       0.93       118
  Mobiles & Accessories    0.95       0.91       0.93       109
   Pens & Stationery       1.00       0.21       0.34         39
   Tools & Hardware        1.00       0.89       0.94         71
  Toys & School Supplies   1.00       0.28       0.44         50
      Watches            1.00       0.95       0.98         88

 accuracy                   0.90        2633
 macro avg              0.82        0.68        0.72        2633
 weighted avg           0.88        0.90        0.88        2633
```

	description	cleaned_desc	predicted_category
0	V&V ART Brass Bracelet - Buy V&V ART Brass Bra...	art brass bracelet buy art brass bracelet flip...	Jewellery
1	Kalpaveda Copper Copper Cuffln ...	kalpaveda copper copper cuff price attract emi...	Jewellery
2	Thelostpuppy Book Cover for Apple iPad Air (Mu...	thelostpuppy book cover apple ipad air multico...	Mobiles & Accessories
3	Riana Copper Copper Bangle - Buy Riana Copper ...	riana copper copper bangle buy riana copper co...	Jewellery
4	Inox Jewelry Stainless Steel Cuffln ...	inox jewelry stainless steel cuff price inox j...	Jewellery
5	Thelostpuppy Book Cover for Apple iPad Air 2 (...)	thelostpuppy book cover apple ipad air multico...	Mobiles & Accessories
6	Ridhi Sidhi Collection Brass Bangle Set (Pack ...	ridhi sidhi collection brass bangle set pack p...	Jewellery
7	Theskinmantra Sleeve for All versions of Apple...	theskinmantra sleeve version apple ipad multic...	Mobiles & Accessories
8	TheLostPuppy Back Cover for Apple iPad Air 2 (...)	thelostpuppy back cover apple ipad air multico...	Mobiles & Accessories
9	Intex Happy Animal Chair Assortment Inflatable...	intex happy animal chair assortment inflatable...	Clothing

Next steps:

[View recommended plots](#)

[New interactive sheet](#)

#### 4) Linear Support Vector Machine

```
1 from sklearn.svm import LinearSVC
2 from sklearn.model_selection import train_test_split
3 from sklearn.feature_extraction.text import CountVectorizer, TfidfTransformer
4 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
5
6 def linear_svm(train_data, test_data):
7     # Define the feature (description) and target (product_category) columns
8     x_train = train_data['cleaned_desc'] # cleaned description in train data
9     y_train = train_data['product_category'] # target labels in train data
10
11     x_test = test_data['cleaned_desc'] # cleaned description in test data
12
13     # Splitting the dataset into training and validation parts (for internal validation)
14     x_train, x_val, y_train, y_val = train_test_split(x_train, y_train, test_size=0.2, random_state=42)
15
16     # Apply CountVectorizer on the training data (fit_transform) and validation/test data (transform)
17     cv = CountVectorizer()
18     x_train = cv.fit_transform(x_train)
19     x_val = cv.transform(x_val)
20     x_test_transformed = cv.transform(x_test)
```

```
21
22 # Apply TF-IDF on the transformed data (training, validation, and test data)
23 vector = TfidfTransformer()
24 x_train = vector.fit_transform(x_train)
25 x_val = vector.transform(x_val)
26 x_test_transformed = vector.transform(x_test_transformed)
27
28 # Initialize the Linear SVC model
29 svc_model = LinearSVC(random_state=42, max_iter=2000)
30
31 # Fitting the model with training data
32 svc_model.fit(x_train, y_train)
33
34 # Predict on the validation set
35 svc_predict = svc_model.predict(x_val)
36
37 # Evaluation metrics for the validation data
38 print("Validation Accuracy: ", accuracy_score(y_val, svc_predict))
39 print("\n***** CONFUSION MATRIX *****")
40 print(confusion_matrix(y_val, svc_predict))
41 print("\n***** CLASSIFICATION REPORT *****")
42 print(classification_report(y_val, svc_predict))
43
44 # Predict the product category for test data (products_test)
45 svc_test_predict = svc_model.predict(x_test_transformed)
46
47 # Returning the predictions for the test data
48 return svc_test_predict
49
50 # Call the function with the train and test data
51 svc_predictions = linear_svm(products, products_test)
52
53 # You can now add these predictions to the test dataframe (products_test)
54 products_test['predicted_category'] = svc_predictions
55
56 # Show the test data with the predicted categories
57 products_test.head(10)
```