n [1]: n [2]: ut[2]:	import numpy as np import pandas as pd import ast import plotly.express as px from plotly import graph_objects as go df = pd.read_csv("flipkart_com-ecommerce_sample.csv") df uniq_id crawl_timestamp
	1 7f7036a6d550aaa89d34c77bd39a5e48 2016-03-25 22:59:23 +0000 http://www.flipkart.com/fabhomedecor-fabric-do Fabric Double Sofa Bed Room Furniture >> Sofa Bed SBEEH3QGU7MFYJFY 32157.0 22646.0 ["http://img6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/image/sofa-bed/j/mg6a.flixcart.com/sofa-bed/j/mg6a.flixcart.com/sofa-bed/
	1997 93e9d343837400ce0d7980874ece471c 2015-12-01 http://www.flipkart.com/elite-collection mediu
[3]: rt[3]:	20002 rows × 15 columns df . head() relail_price discounted_price image in the product_under solid women's Coloning > Lingeriee, Cycling Shorts 1 7f7036a6d550aaa89d34c77bd39a5e48 2016-03-25 http://www.flipkart.com/fabhomedeccor Fabric Double Room Furniture >> Sola SBEEH3QGU7MFYJFY 32157.0 22646.0 ["http://img6a.flixcart.com/image/sofa-bed/lyft
	2 f449ec65dcbc041b6ae5e6a32717d01b
[4]: :[4]: [5]:	df.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 20002 entries, 0 to 20001 Data columns (total 15 columns): # Column Non-Null Count Dtype</class>
	3
[6]:	<pre>df.isnull().sum() uniq_id</pre>
[7]: [7]:	brand product_specifications 16 dtype: int64 df["retail_price"].fillna(df["retail_price"].median(),inplace=True) df["discounted_price"].fillna(df["discounted_price"].median(),inplace=True) df.head() retail_price "].fillna(df["discounted_price"].median(),inplace=True) df.head() retail_price discounted_price
	1 7f7036a6d550aaa89d34c77bd39a5e48 2016-03-25 22:59:23 +0000 http://www.flipkart.com/abhomedecorfabric-do Fabric Double Fabric Double Sofa Bed ["Furniture >> Living Fabric Double Sofa Bed SBEEH3QGU7MFYJFY 32157.0 22646.0 ["http://img6a.flixcart.com/image/sofa-bed/j/f 2 1449ec65dcbc041b6ae5e6a32717d01b 2016-03-25 22:59:23 +0000 http://www.flipkart.com/aw-bellies/p/itmeh4grg AW Bellies Footwear >> Women's Footwear >> Ballerinas Solid-Women's Clothing >> Livingerie, Cycling Shorts SHOEH4GRSUBJGZXE 999.0 499.0 ["http://img5a.flixcart.com/image/short/6/2/h 3 0973b37acd0c664e3de26e97e5571454 2016-03-25 22:59:23 +0000 http://www.flipkart.com/sicons-all-purpose-arm Clothing >> Livingerie, Cycling Shorts SRTEH2F6HUZMQ6SJ 699.0 267.0 ["http://img5a.flixcart.com/image/short/6/2/h/ 4 bc940ea42ee6bef5ac7cea3fb5cfbee7 2016-03-25 22:59:23 +0000 http://www.flipkart.com/sicons-all-purpose-arm Process Armica Dog Shampoo PSOEH3ZYDMSYARJ5 220.0 210.0 ["http://img5a.flixcart.com/image/pet-shampoo/ shampoo/
[8]: [9]:	x=df['retail_price']-df['discounted_price'] y=(x/df['retail_price'])*100 df['discount_percentage']=y df['timestamp'] = pd.to_datetime(df['crawl_timestamp'], errors='coerce') # Extract the time part of the timestamp (handling NaT gracefully) df['time'] = df['timestamp'].apply(lambda x: x.time() if pd.notnull(x) else None) # Extract the date part of the timestamp (handling NaT gracefully) df['date'] = df['timestamp'].apply(lambda x: x.date() if pd.notnull(x) else None)
[9]:	df . head()
[10]:	2 f449ec65dcbc041b6ae5e6a32717d01b 22:59:23 +0000 bellies/p/itmeh4grg AW Bellies Footwear >> Ballerinas > SHOEH4GRSUBJGZXE 999.0 499.0 ["http://img5a.flixcart.com/image/shoe/7/z/z/r] 3 0973b37acd0c664e3de26e97e5571454 2016-03-25 22:59:23 +0000 women-s-c Cycling Shorts Coloring >> Lingerie, Cycling Shorts Sl STEH2F6HUZMQ6SJ 699.0 267.0 ["http://img5a.flixcart.com/image/short/6/2/h/] 4 bc940ea42ee6bef5ac7cea3fb5cfbee7 2016-03-25 22:59:23 +0000 http://www.flipkart.com/sicons-all-purpose-arn Purpose-arn Dog Shampoo Coat Care Coat Care Coat Care Coat Care Coat Care PSOEH3ZYDMSYARJ5 220.0 210.0 ["http://img5a.flixcart.com/image/pet-shampoo/ Shampoo/ Shampoo/ Coat Care Coat Care Coat Care Coat Care Coat Care Shome)
10]:	df .head() The difference of
	22:59:23 +0000 bellies/p/itmeh4grg 3 0973b37acd0c664e3de26e97e5571454 2016-03-25 22:59:23 +0000 http://www.flipkart.com/alisha-solid-Women's Cycling Shorts 4 bc940ea42ee6bef5ac7cea3fb5cfbee7 2016-03-25 22:59:23 +0000 http://www.flipkart.com/sicons-all-purpose-arm Counting top 10 products in main category column Alisha Solid Women's Clothing >> Women's Clothing >> Women's Clothing >> Lingerie, SRTEH2F6HUZMQ6SJ 699.0 267.0 ["http://img5a.flixcart.com/image/short/6/2/h/ SRTEH2F6HUZMQ6SJ 699.0 267.0 ["http://img5a.flixcart.com/ima
[11]: [11]:	df['main_category'] 0
[12]:	<pre>n = 10 top_products = pd.DataFrame(df['main_category'].value_counts()[:n]).reset_index() print("Before renaming:") print(top_products) Before renaming:</pre>
[13]:	<pre>7 Home Furnishing 700 8 Kitchen & Dining 647 9 Computers 578 # Rename columns top_products.columns = ['Top_Products', 'Total_count'] print("\nAfter renaming:") print(top_products) After renaming:</pre>
[14]:	4 Automotive 1012 5 Home Decor & Festive Needs 929 6 Beauty and Personal Care 710 7 Home Furnishing 700 8 Kitchen & Dining 647 9 Computers 578 Top 10 main brands being purchased df['brand'] 0 Alisha FabHomeDecor
15]:	AW 3 Alisha 4 Sicons 19997 Elite Collection 19998 Elite Collection 19999 Elite Collection 20000 NaN 20001 NaN Name: brand, Length: 20002, dtype: object #Top 10 main brands being purchased n = 10 top_brands=pd.DataFrame(df['brand'].value_counts()[:n]).reset_index()
	print("Before renaming:") print(top_brands) Before renaming:
16]:	<pre># Rename columns top_brands.columns = ['Top Brands', 'Total count'] print("\nAfter renaming:") print(top_brands) After renaming: Top Brands Total count 0 Allure Auto</pre>
17]:	<pre>7 White 155 8 DailyObjects 144 9 Speedwav 141 from plotly.subplots import make_subplots label1=top_products['Top_Products'] value1=top_products['Total_count'] label2=top_brands['Top Brands'] value2=top_brands['Total_count']</pre>
	<pre>#Create subplots fig_both = make_subplots(rows=1, cols=2, specs=[[{'type': 'domain'}, {'type': 'domain'}]]) fig_both.add_trace(go.Pie(labels=label1, values=value1, name="Top Products", pull=[0.3, 0, 0, 0]),1,1) fig_both.add_trace(go.Pie(labels=label2, values=value2, name="Total Brands", pull=[0.3, 0, 0, 0]), 1, 2) #use hole to create a donut-like pie chart fig_both.update_traces(hole=.4, hoverinfo="label+percent+name") #fig_both.update_traces (hoverinfo="Label+percent+name")</pre>
	fig_both.update_layout(title_text="Top products and brands distribution", #Add annotations in the center of the donut pies annotations=[dict(text='Product', x=0.18, y=0.5, font_size=20, showarrow=False),
	21.2% 19.4% 19.4% 19.4% 19.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10.4% 10
	6.09% 6.91% 9.48% 6.91% Slim TheLostPuppy Karatcraft Black White
	<pre># Filter for high discounts (discount_percentage > 90) df_discount = df.query('discount_percentage > 90') # Drop rows with missing values df_discount = df_discount.dropna() # Correct spelling errors in the 'brand' column df_discount["brand"].replace('FashBlush', 'Fash Blush', inplace=True) # Calculate average discount percentage by brand max_discount = (df_discount</pre>
	.sort_values(by='discount_percentage', ascending=False) .reset_index()) print(max_discount) brand discount_percentage 0 Rajcrafts 96.533333 1 Bling 94.548458 2 Fash Blush 92.711714 3 Mydress Mystyle 91.991992 4 Soulful Threads 91.952663 5 Instella 91.719745 6 Bond Beatz 91.596639 7 Fashblush 91.132525 8 Black 90.681676
19]:	<pre>9 KazamaKraft 90.565618 10 Zaicus 90.143281 11 CUBA 90.045023 12 SDZ 90.045023 13 Gia 90.020004 # Create a bar plot with enhancements fig = px.bar(max_discount, x='brand', y='discount_percentage', color='brand', color='brand', color_discrete_sequence=px.colors.qualitative.Plotly, title='Average Discount Percentage by Brand', labels={'discount_percentage': 'Average Discount Percentage', 'brand': 'Brand'}</pre>
	# Update layout for better readability fig.update_layout(xaxis_title='Brand', yaxis_title='Average Discount Percentage') # Show the plot fig.show() Average Discount Percentage by Brand
	Brand Rajcrafts Bling Fash Blush Mydress Mystyle Soulful Threads Instella Bond Beatz Fashblush Black KazamaKraft Zaicus CUBA
20]:	# Group by customer and calculate total spending df_customer = df.groupby("uniq_id")[["discounted_price"]].sum().sort_values(by=['discounted_price'], ascending=[False])
	<pre># Select the top 20 customers with the highest spending top_20_customers = df_customer.head(20) # Ensure 'uniq_id' is a column in the DataFrame top_20_customers = top_20_customers.reset_index() # Reset index to make 'uniq_id' a column # Create a bar plot fig = px.bar(top_20_customers, x='uniq_id', y='discounted_price', color='discounted_price', color='discounted_price', color=continuous_scale=px.colors.diverging.Portland, # Customize the color scale if needed title='Top_20_Customers by Spending',</pre>
	labels={'discounted_price': 'Total Spending', 'uniq_id': 'Customer ID'} # Update layout for better readability fig.update_layout(
	Total Spending Total Spending Total Spending 400k 400k 300k 200k
	100k 100k 100k 100k
21]:	<pre># Filter for 5-star products rating_5 = df[df['product_rating'] == '5'] # Count of main categories with 5-star ratings top_product_type = rating_5['main_category'].value_counts() # Count of brands with 5-star ratings top_brand_type = rating_5['brand'].value_counts() # Top 5 product categories df_top_product = pd.DataFrame(top_product_type.head(5).reset_index()) df_top_product.columns = ['top_prod', 'count']</pre>
	<pre># Top 5 brands df_top_brand = pd.DataFrame(top_brand_type.head(5).reset_index()) df_top_brand.columns = ['top_brands', 'count'] # Display the DataFrames print(df_top_product) print("-"*50) print(df_top_brand) top_prod count 0 Clothing 232 1 Jewellery 70 2 Footwear 47 3 Watches 47</pre>
22]:	top_brands count Regular 21 Slim 13 Black 9 Black 9 JDX 8 Top-rated products and brands # Concatenate the DataFrames horizontally df_product_brand_rate5 = pd.concat([df_top_product, df_top_brand], axis=1)
	<pre># Remove rows with 'No rating available' df.drop(df.index[df['product_rating'] == 'No rating available'], inplace=True) # Count of each rating ratings = pd.DataFrame(df['product_rating'].value_counts().reset_index()) ratings.columns = ['Ratings', 'Counts'] # Convert 'Ratings' to float for numerical sorting ratings['Ratings'] = ratings['Ratings'].astype(float) # Sort by 'Ratings' in descending order ratings = ratings.sort_values(by=['Ratings'], ascending=[False]) # Plot the result</pre>
	<pre>x = ratings['Ratings'] y = ratings['Counts'] figdot2 = go.Figure() figdot2.add_trace(go.Scatter(</pre>
	xaxis_title="Katings", yaxis_title="Count",) figdot2.update_xaxes(showline=True, linewidth=1, linecolor='black', mirror=True) figdot2.update_yaxes(showline=True, linewidth=1, linecolor='black', mirror=True) figdot2.show() Ratings vs Count
	600 500 400 200
	Visualize the trend of average retail and discounted prices over time using an area plot.
23]:	<pre># Group and aggregate data df_date_retail = pd.DataFrame(df.groupby("date")[["retail_price"]].mean().reset_index()) df_date_discount = pd.DataFrame(df.groupby("date")[["discounted_price"]].mean().reset_index()) # Concatenate DataFrames df_date_price = pd.concat([df_date_retail, df_date_discount], axis=1) # Remove duplicate columns df_date_price = df_date_price.loc[:, ~df_date_price.columns.duplicated()] # Prepare data for plotting x = df_date_price['date'] y1 = df_date_price['retail_price']</pre>
	<pre>y2 = df_date_price['discounted_price'] # Create and customize plot fig_area2 = go.Figure() fig_area2.add_trace(go.Scatter(</pre>
	<pre>y=y2, fill='tozeroy', name='Discount Price', line=dict(width=0.5, color='darkslategray'))) fig_area2.update_layout(xaxis_title="Dates", yaxis_title="Price (in 1000s)", plot_bgcolor='white') fig_area2.update_xaxes(showline=True, linewidth=1, linecolor='black', mirror=True) fig_area2.update_yaxes(showline=True, linewidth=1, linecolor='black', mirror=True)</pre>
	fig_area2.show() Retail Price Discount Price
	5000 (800) 4000 2000 1000
24]:	Dates Visualize how the number of clicks on a product varies over time # Create scatter plot scat2 = px. scatter(x=df['Time'].sort_values(ascending=True), y=df['product_url']) # Update layout
	<pre>scat2.update_layout(title='No. of clicks vs time', # Title of the plot xaxis_title='Time', # X-axis label yaxis_title='No. of Clicks' # Y-axis label) # Update axes scat2.update_xaxes(showline=True, linewidth=1, linecolor='black', mirror=True) scat2.update_yaxes(showline=True, linewidth=1, linecolor='black', mirror=True) # Hide y-axis tick labels scat2.update_yaxes(showticklabels=False) # Display plot</pre>
	No. of clicks vs time No. of clicks vs time
	No. of Click
25]:	# Calculate total products total_prod = len(df['pid']) # Calculate total rated products total_prod = len(df['pid']) # Calculate total rated products total_prod = len(df['pid']) # Calculate total rated products total_prod = len(df['pid'])
	<pre>total_ratings = len(df[df['product_rating'] != 'No rating available']) # Calculate 5-star rated products top_ratings = len(df[df['product_rating'] == '5']) # Prepare data for funnel plot df_funnel_1 = pd.DataFrame({ 'number': [total_prod, total_ratings, top_ratings], 'stage': ["Total Products", "Products with Ratings", "Products with 5-Star Rating"] }) # Create and display funnel plot funnel_1_fig = px.funnel(df_funnel_1, x='number', y='stage') funnel_1_fig.show()</pre>
	Total Products 1851
	Products with S-Star Rating Products with 5-Star Rating 620
[]:	