<pre>In [1]: In [2]:</pre>	<pre>import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns</pre>
In [3]: In [4]:	<pre>df_purchase = pd.read_csv("C:/Users/sona3/Downloads/QVI_purchase_behaviour (1).csv") df_transaction.head(10)</pre>
Out[4]: In [5]: Out[5]:	0 43390 1 1000 1 5 Natural Chip Compny SeaSalt175g 2 6.0 1 43599 1 1307 348 66 CCS Nacho Cheese 175g 3 6.3 2 43605 1 1343 383 61 Smiths Crinkle Cut Chips Chicken 170g 2 2.9 3 43329 2 2373 974 69 Smiths Chip Thinly S/Cream&Onion 175g 5 15.0 4 43330 2 2426 1038 108 Kettle Tortilla ChpsHny&Ilpno Chili 150g 3 13.8 5 43604 4 4074 2982 57 Old El Paso Salsa Dip Tonato Mild 300g 1 5.1 6 43601 4 4149 3333 16 Smiths Crinkle Chips Salt & Vinegar 330g 1 5.7 7 43601 4 4196 3539 24 Grain Waves Sweet Chill 210g 1 3.6 8 43332 5 5026 4525 42 Doritos Corn Chip Mexican Jalapeno 150g 1 3.9 9 43330 7 7150 6900 52 Grain Waves Sour Cream&Chives 210G 2 7.2 Data Summary
<pre>In [6]: Out[6]: In [7]:</pre>	DATE int64 STORE_NBR int64 LYLTY_CARD_NBR int64 LYLTY_CARD_NBR int64 PROD_NBR int64 PROD_NBR int64 PROD_OTY int64 TOT_SALES float64 dtype: object df_transaction.describe()
Out[7]: In [8]:	count 264836.00000 264836.00000 2.64836.00000 264836.00000 264836.00000 mean 43464.036260 135.08011 1.355495e+05 1.351583e+05 56.583157 1.907309 7.304200 std 105.389282 76.78418 8.057998e+04 7.813303e+04 32.826638 0.643654 3.083226 min 43282.000000 1.00000 1.000000e+03 1.000000e+00 1.000000 1.500000 25% 43373.00000 70.0000 7.002100e+04 6.760150e+04 28.000000 2.000000 5.400000 50% 43464.000000 130.00000 1.303575e+05 1.351375e+05 56.000000 2.000000 7.400000 75% 43555.000000 203.00000 2.030942e+05 2.027012e+05 85.000000 200.00000 650.000000
Out[8]: In [9]: Out[9]:	STORE_NBR 0 LYLTY_CARD_NBR 0 TXN_ID 0 PROD_NBR 0 PROD_NAME 0 PROD_OTY 0 TOT_SALES 0 dtype: int64 LYLTY_CARD_NBR LIFESTAGE PREMIUM_CUSTOMER
	01000YOUNG SINGLES/COUPLESPremium11002YOUNG SINGLES/COUPLESMainstream21003YOUNG FAMILIESBudget31004OLDER SINGLES/COUPLESMainstream41005MIDAGE SINGLES/COUPLESMainstream51007YOUNG SINGLES/COUPLESBudget61009NEW FAMILIESPremium71010YOUNG SINGLES/COUPLESMainstream81011OLDER SINGLES/COUPLESMainstream91012OLDER FAMILIESMainstream
In [10]: Out[10]: In [11]:	(72637, 3)
Out[11]: In [12]: Out[12]:	LIFESTAGE object PREMIUM_CUSTOMER object dtype: object df_purchase.describe()
	count 7.263700e+04 mean 1.361859e+05 std 8.989293e+04 min 1.000000e+03 25% 6.620200e+04 50% 1.340400e+05
In [13]: Out[13]:	75% 2.033750e+05 max 2.373711e+06 df_purchase.isnull().sum()
In [14]: In [15]:	<pre>changing Date Format df_transaction['DATE'] = pd.to_datetime(df_transaction['DATE'], unit='d', origin='1899-12-30')</pre>
Out[15]:	DATE STORE_NBR LYLTY_CARD_NBR TXN_ID PROD_NAME PROD_QTY TOT_SALES 0 2018-10-17 1 1000 1 5 Natural Chip Compny SeaSalt175g 2 6.0 1 2019-05-14 1 1307 348 66 CCs Nacho Cheese 175g 3 6.3 2 2019-05-20 1 1343 383 61 Smiths Crinkle Cut Chips Chicken 170g 2 2.9 3 2018-08-17 2 2373 974 69 Smiths Chip Thinly S/Cream&Onion 175g 5 15.0 4 2018-08-18 2 2426 1038 108 Kettle Tortilla ChpsHny&Jlpno Chili 150g 3 13.8
In [16]: In [17]:	<pre>Extraction df_transaction['PACKET_SIZE'] = df_transaction['PROD_NAME'].str.extract(r'(\d+)\s*[gG]?')</pre>
Out[17]: In [18]: In [19]:	0 2018-10-17 1 1000 1 5 Natural Chip Compny SeaSalt175g 2 6.0 175 1 2019-05-14 1 1307 348 66 CCs Nacho Cheese 175g 3 6.3 175 2 2019-05-20 1 1343 383 61 Smiths Crinkle Cut Chips Chicken 170g 2 2.9 170 3 2018-08-17 2 2373 974 69 Smiths Chip Thinly S/Cream&Onion 175g 5 15.0 175 4 2018-08-18 2 2426 1038 108 Kettle Tortilla ChpsHny&Jlpno Chili 150g 3 13.8 150 df_transaction['PROD_NAME'] . str.extract(r'^(\w+)')
Out[19]: In [75]:	0 2018-10-17 1 1000 1 5 Natural Chip Compny SeaSalt175g 2 6.0 175 Natural 1 2019-05-14 1 1307 348 66 CCs Nacho Cheese 175g 3 6.3 175 CCs 2 2019-05-20 1 1343 383 61 Smiths Crinkle Cut Chips Chicken 170g 2 2.9 170 Smiths 3 2018-08-17 2 2373 974 69 Smiths Chip Thinly S/Cream&Onion 175g 5 15.0 175 Smiths 4 2018-08-18 2 2426 1038 108 Kettle Tortilla ChpsHny&Jlpno Chili 150g 3 13.8 150 Kettle Merging
In [76]: Out[76]: In [77]:	DATE STORE_NBR LYLTY_CARO_NBR TXN_ID PROD_NBR
In [78]: In [79]:	dr_cleaned = merged_dr[(merged_dr['101_SALES'] >= lower_bound) & (merged_dr['101_SALES'] <= upper_bound)]
In [80]: In [81]: In [82]:	<pre># Total Sales Total_sales = merged_df('TOT_SALES').sum() print(f"Total Sales : (Total_sales :,.0f) ") Total Sales : 1,934,415 # Total number of customers Total_customers = merged_df('LYLTY_CARD_NBR'].nunique() print(f"Total Customers : {Total_customers :,.0f} ") Total Customers : 72,637 # average number of transactions per customer avg_transactions = merged_df.shape[0] / merged_df('LYLTY_CARD_NBR'].nunique() print(f"Average transactions per customer: (avg_transactions:.0f)") Average transactions per customer: 4</pre>
	plt. ylabel ('Total Sales') plt. tight (layout()) plt. show() Total Sales by Brand Total Sales by Brand Total Sales by Brand Total Sales by Brand Brand
In [84]:	# Which customer_segment should target customer_segment = merged_df.groupby('PREMIUM_CUSTOMER')['TOT_SALES'].sum().sort_values(ascending=False) # Plotting bar graph plt.figure(figsize=(10, 6)) customer_segment.plot(kind='bar', color='blue') plt.title('Total Sales by Customer_segment') plt.ylabel('Total Sales') plt.xlabel('Customer_segment') plt.xlabel('Total Sales') plt.tight_layout() plt.show() Total Sales by Customer_segment Total Sales by Customer_segment
In [85]:	500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000000 - 500000 - 500000 - 500000 - 500000 - 500000 - 500000 - 5000
	# Which age group should target age_group = merged_df.groupby('LIFESTAGE')['TOT_SALES'].sum().sort_values(ascending=False) # Plotting bar graph plt.figure(figsize=(10, 6)) age_group.plot(kind='bar', color='blue') plt.title('Total Sales by Age Group') plt.xlabel('Age Group') plt.xlabel('Total Sales') plt.tight_layout() plt.show() Total Sales by Age Group Total Sales by Age Group
In [86]: In [87]:	merged_ar['PACKE1_S1ZE'] = pa.to_numeric(merged_ar['PACKE1_S1ZE'], errors='coerce')
Out [87]:	mergea_ar.atypes
In [88]: In [89]:	BRAND object LIFESTAGE object PREMIUM_CUSTOMER object dtype: object # Are there any relation between Total Sales and Packet size ? correlation = merged_df['PACKET_SIZE'].corr(merged_df['TOT_SALES'])
	<pre>correlation = merged_ar['PACKEI_SIZE'].corr(merged_ar['IOI_SALES']) print(f"Correlation between packet size and total sales: {correlation:.2f}") # Step 4: Visualize plt.figure(figsize=(8, 6)) sns.scatterplot(data=merged_df, x='PACKET_SIZE', y='TOT_SALES') plt.title('Packet Size vs Total Sales') plt.xlabel('Packet Size (in grams)') plt.ylabel('Total Sales') plt.grid(True) plt.tight_layout()</pre>
	plt.show() Correlation between packet size and total sales: 0.31 Packet Size vs Total Sales 600
In [93]:	<pre>#Which Product Sizes are most popular overall # Convert PACKET_SIZE to string to handle it as a category merged_df['PACKET_SIZE'] = merged_df['PACKET_SIZE'].astype(str) plt.figure(figsize=(10, 6)) sns.countplot(data=merged_df, x='PACKET_SIZE', order=merged_df['PACKET_SIZE'].value_counts().index[:10], color = 'blue') plt.title('Top 10 Most Popular Pack Sizes') plt.xlabel('Pack Size') plt.ylabel('Number of Transactions') plt.ylabel('Number of Transactions') plt.ticks(rotation=45) plt.tight_layout() plt.show()</pre>
In [91]:	Top 10 Most Popular Pack Sizes 50000 50000 10000 10000 Pack Size
In [91]: Out[91]: In [92]:	merged_df.groupby('STORE_NBR')['TOT_SALES'].sum().sort_values(ascending=False).head(1) STORE_NBR 226 18905.45 Name: TOT_SALES, dtype: float64
	<pre># Total sales by Month # Extract month from date merged_df['MONTH'] = merged_df['DATE'].dt.month # Group sales by month monthly_sales = merged_df.groupby('MONTH')['TOT_SALES'].sum().reset_index() # Create a DataFrame with all 12 months all_months = pd.DataFrame({'MONTH': range(1, 13)}) # Merge and fill missing months with 0 monthly_sales_full = pd.merge(all_months, monthly_sales, on='MONTH', how='left').fillna(0) # Plot plt.figure(figsize=(10, 5)) plt.plot(monthly_sales_full('MONTH'), monthly_sales_full('TOT_SALES'), marker='o', color='blue', linestyle='-') plt.xileb('Total Sales by Month') plt.xileb('Total Sales by Month') plt.ylabel('Total Sales') plt.xticks(ticks=range(1, 13)) plt.xticks(ticks=range(1, 13))</pre>
	plt.tight_layout() plt.show() Total Sales by Month 167500 162500
In []:	15000 15000 1 2 3 4 5 6 7 8 9 10 11 12