

We need to present a strategic recommendation to our client that is supported by data which she can trust. So, we need to analyse the data to understand the current purchasing trends and behaviours. The client is interested in understanding the current chip purchasing behaviour. Consider what metrics would help describe the customers' purchasing behaviour.

You will also want to derive extra features such as pack size and brand name from the data and define metrics of interest to enable you to draw insights on who spends on chips and what drives spends for each customer segment. Remember, our end goal is to form a strategy based on the findings to provide a clear recommendation to our client the Category Manager so make sure your insights can have a commercial application.

Solution:

Import libraries and Loading datasets

```
# for data visualization
import matplotlib.pyplot as plt
import matplotlib inline
import seaborn as sns
```

Checking the datasets

```

# Filtering the merged data set only include top products
top_products_sales = merged_dataset[merged_dataset['PROD_NAME'].isin(top_products['PROD_NAME'])]

# Aggregating monthly sales for the top products
heatmap_data = top_products_sales.groupby(['Month_Year', 'PROD_NAME'])['TOT_SALES'].sum().unstack(fill_value=0)

# Plotting the heatmap
plt.figure(figsize=(12, 8))

# Create a heatmap
sns.heatmap(heatmap_data, cmap='YlGnBu', annot=True, fmt='.1f', linewidths=.5)

plt.title('Top 15 Performing Products - Total Sales Heatmap')
plt.xlabel('Products')
plt.ylabel('Month-Year')
plt.xticks(rotation=45)
plt.yticks(rotation=45)

plt.tight_layout()
plt.show()

```

Month-Year	Cheezitos Cheese 1200g	Doritos Corn Chip	Doritos Cheese	Kettle Honey Soy	Kettle Mozzarella	Kettle Sea Salt	Kettle Sweet Chili And Sour Cream 175g	Old El Paso Salsa	Old El Paso Salsa Dip Chkly Bm Hotdog	Old El Paso Salsa Dip Tomato Med 200g	Smalls Crinkle	Smalls Crinkle Chps-Slt	Smalls Crinkle Chip	Original 200g	Original 300g
2018-07	2462.4	3347.5	2291.4	2646.0	2554.2	2910.6	2651.4	2359.8	2548.8	2402.1	2249.1	2351.1	2479.5	2650.5	2997.2
2018-08	2394.0	2886.6	2490.9	2413.8	2500.2	2764.8	2527.2	2548.8	2511.0	2340.9	2713.2	2346.0	2576.4	3032.4	2814.3
2018-09	2889.9	3172.0	3237.6	2710.8	2856.6	2862.0	2867.4	2835.0	2478.6	2677.5	2346.0	2340.9	2474.7	2701.6	2955.9
2018-10	3163.5	3471.0	2810.1	2478.6	3024.0	2705.4	2635.2	2824.2	2943.0	2764.2	2493.9	2713.2	2958.3	2821.5	3469.2
2018-11	2650.5	3445.0	2610.6	2462.4	2494.8	2646.0	2797.2	2527.2	2862.0	2555.1	2478.6	2391.9	2610.6	2889.9	2885.1
2018-12	3072.3	3081.0	2935.5	2667.6	2802.6	3272.4	2635.2	3007.8	2721.6	2483.7	2825.4	2733.6	3015.3	2878.5	3056.2
2019-01	2761.6	3289.0	2673.3	2781.0	2683.6	2970.0	2462.4	2619.0	2646.0	2937.6	2616.3	2539.8	2941.2	3095.1	3062.1
2019-02	2644.8	3263.0	2587.8	2532.6	2511.0	2413.8	2359.8	2478.6	2435.4	2473.5	2422.5	2366.4	2867.1	2889.9	2938.2
2019-03	2998.2	3503.5	2969.7	2538.0	2894.4	3024.0	3072.6	2840.4	2867.4	2193.0	2636.7	2636.7	2804.4	2775.9	3156.5
2019-04	2907.0	3100.5	2775.9	2397.6	2775.6	2894.4	2781.0	2802.6	2953.8	2478.6	2289.9	2565.3	3163.5	2901.3	3091.6
2019-05	2815.8	2713.8	2587.8	2516.4	2262.8	2208.6	2667.6	2478.6	2527.2	2325.6	2244.0	2289.5	2850.0	2924.1	2655.0
2019-06	2815.8	3432.0	2889.9	2613.6	2748.6	3088.8	2597.4	2910.6	3164.4	2432.7	2606.1	2305.2	2867.1	2832.9	2590.1
2019-07	262.2	214.5	153.9	253.8	162.0	172.8	162.0	199.8	145.8	219.3	102.0	183.6	159.6	176.7	236.0

```

In [25]: # Group by LIFESTAGE and sum the total sales
lifestage_sales = merged_dataset.groupby(['LIFESTAGE'])['TOT_SALES'].sum().reset_index()

# Create a pie chart
plt.figure(figsize=(10, 8))
plt.pie(lifestage_sales['TOT_SALES'], labels=lifestage_sales['LIFESTAGE'], autopct='%1.1f%%', startangle=140)
plt.title('Sales Distribution by LIFESTAGE')
plt.axis('equal') # Equal aspect ratio ensures that pie is drawn as a circle.
plt.tight_layout()
plt.show()

```

LIFESTAGE	TOT_SALES	Percentage
YOUNG SINGLES/COUPLES	13.5%	13.5%
YOUNG FAMILIES	16.3%	16.3%
MIDDLE SINGLES/COUPLES	16.3%	16.3%
Other	16.3%	16.3%



OLDER SINGLES/COUPLES

```

In [26]: merged_dataset['Month_Year'] = merged_dataset['DATE'].dt.to_period('W')

# Aggregating monthly sales by PREMIUM_CUSTOMER category
monthly_sales_comparative = merged_dataset.groupby(['Month_Year', 'PREMIUM_CUSTOMER'])['TOT_SALES'].sum().unstack(fill_value=0)

# Reset the index for plotting
monthly_sales_comparative = monthly_sales_comparative.reset_index()

# Set a dark grid style
sns.set(style='darkgrid')

# Plotting the results

```

```
plt.xticks(rotation=45)
plt.legend(title='Premium Customer Category', fontsize=12)
plt.grid(True)
plt.tight_layout()

# Show the plot
plt.show()
```

Key Findings for Each Analysis:

Data Cleaning and Preparation:

1. Data Integrity:
No null values were present in the datasets. Duplicates were minimal (only one duplicate in the transaction dataset).
2. Outlier Removal:
Identified and removed outliers from the TOT_SALES column using the IQR method.
3. Data Merging:
Datasets were merged on LYLTY_CARD_NBR, resulting in 264,258 entries.
4. Feature Engineering:
Converted dates from numerical format to datetime and categorized lifetime and premium customer columns.

various.

4. Monthly Sales Trends:

Steady sales trends across 2018-2019, with peaks during specific months such as December (likely due to seasonal demand).

5. Sales Distribution by Lifestyle:

Older demographic groups formed the bulk of sales, aligning with the trend observed in total sales by lifestyle.

6. Monthly Sales Comparison by Customer Category:

Budget and Mainstream categories showed more consistent performance compared to Premium customers, which had sporadic peaks.

7. Top Product Sales Heatmap:

Certain products showed seasonal trends, peaking during festive periods.

Summary/Highlights:

Older demographics are the primary contributors to chip sales, especially in the Mainstream and Budget segments. Product preference varies, with a clear bifurcation toward popular chip brands like Natural Chn. Co. and GGA. Seasonal trends suggest the importance of strategic promotions during high-demand months.