

✓ Sales Analysis: Proj 2 Col Rakesh Pedram

✓ Data Wrangling and Import

Using Google Colabs with files in google drive.

```
1
2 import pandas as pd
3 from google.colab import drive
4 drive.mount('/content/drive')
5 import numpy as np
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8
9 # Get the file path from Google Drive
10 file_path = '/content/drive/My Drive/Dataset/apparel.csv'
11
12 # Read the file into a Pandas DataFrame
13 df = pd.read_csv(file_path)
14 print(df.head())
15 print(df.info())
16 # Describe numerical variables to check for missing values, outliers, etc.
17 print(df.describe())
18 #df['Dt_Customer'] = pd.to_datetime(df['Dt_Customer'], format='%d-%m-%Y')
19 print(df.info())
```

Mounted at /content/drive

	Date	Time	State	Group	Unit	Sales
0	01-Oct-20	Morning	WA	Kids	8	20000
1	01-Oct-20	Morning	WA	Men	8	20000
2	01-Oct-20	Morning	WA	Women	4	10000
3	01-Oct-20	Morning	WA	Seniors	15	37500
4	01-Oct-20	Afternoon	WA	Kids	3	7500

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 7560 entries, 0 to 7559

Data columns (total 6 columns):

#	Column	Non-Null	Count	Dtype
0	Date	7560	non-null	object
1	Time	7560	non-null	object
2	State	7560	non-null	object
3	Group	7560	non-null	object
4	Unit	7560	non-null	int64
5	Sales	7560	non-null	int64

dtypes: int64(2), object(4)

memory usage: 354.5+ KB

None

	Unit	Sales
count	7560.000000	7560.000000
mean	18.005423	45013.558201
std	12.901403	32253.506944
min	2.000000	5000.000000
25%	8.000000	20000.000000
50%	14.000000	35000.000000
75%	26.000000	65000.000000
max	65.000000	162500.000000

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 7560 entries, 0 to 7559

Data columns (total 6 columns):

#	Column	Non-Null	Count	Dtype
0	Date	7560	non-null	object
1	Time	7560	non-null	object
2	State	7560	non-null	object
3	Group	7560	non-null	object
4	Unit	7560	non-null	int64
5	Sales	7560	non-null	int64

dtypes: int64(2), object(4)

memory usage: 354.5+ KB

None

Check for outliers for sales and unit numeric field. Data is clean. Unit and sales are integers only date field needs conversion

```

1 condition = df['Sales'] > 6500 # Sales more than 75 percentile to check outliers
2 df_sorted = df[condition].sort_values(by='Sales')
3
4 # Display the sorted DataFrame
5 print(df_sorted['Sales'])
6 condition = df['Unit'] < 8 # Sales more than 75 percentile to check outliers
7 df_sorted = df[condition].sort_values(by='Unit')
8
9 # Display the sorted DataFrame
10 print(df_sorted['Unit'].count())

1643

```

There are no outliers in the Sales and unit. Will run a histogram to see my numeric data is normalised.

```

1 # Convert date field to date obj

1 df['Date']=pd.to_datetime(df['Date'])
2 print (df['Date'].dtype)
3 #Converted to date time field now can be used in timeline charts
4 #Rest is string data

datetime64[ns]

1

datetime64[ns]

1 # Check no nulls in date field
2 df[df.notna()]
3 df[df.isna()]
4 df[df.notna()].count()
5 #df[df.isna()].count()
6 print(df.isna().sum())

Date      0
Time      0
State     0
Group     0
Unit      0
Sales     0
dtype: int64

```

I have no nulls data is clean. no negative values.

✓ Normalisation of data

I studied the problem and there is no analysis for which I might need a normalised data that will be required to be fed in a model. However sticking to the directions. Used a Min-Max Scaler to normalise my sales and Unit data just for learning. Using built in sklearn lib function.

```

1 from sklearn.preprocessing import MinMaxScaler
2 # Create a min-max scaler object
3 scaler = MinMaxScaler()
4
5 # Fit and transform the dataframe
6 df_norm = scaler.fit_transform(df[['Unit', 'Sales']])
7
8 # Convert the numpy array to a dataframe
9 df_norm = pd.DataFrame(df_norm, columns=['Unit_norm', 'Sales_norm'])
10
11 # Print the dataframe
12 print(df_norm)

      Unit_norm  Sales_norm
0      0.095238    0.095238
1      0.095238    0.095238
2      0.031746    0.031746
3      0.206349    0.206349
4      0.015873    0.015873
...         ...         ...
7555  0.190476    0.190476
7556  0.206349    0.206349
7557  0.206349    0.206349
7558  0.142857    0.142857
7559  0.174603    0.174603

```

[7560 rows x 2 columns]

Using the Group by function to make analysis

```
1 # Group data by state or other relevant columns
2 grouped_data = df.groupby("State")
3 print(grouped_data.first())
4 #Group by Time
5 grouped_data = df.groupby("Time")
6 print(grouped_data.first())
```

	Time	Group	Unit	Sales	
State					
NSW	Morning	Kids	39	97500	
NT	Morning	Kids	13	32500	
QLD	Morning	Kids	20	50000	
SA	Morning	Kids	12	30000	
TAS	Morning	Kids	13	32500	
VIC	Morning	Kids	49	122500	
WA	Morning	Kids	8	20000	
Time					
	Afternoon	WA	Kids	3	7500
	Evening	WA	Kids	15	37500
	Morning	WA	Kids	8	20000

✚ pivot tables

The pivot tables are ideal for this data

```
1 # Example 1: Total sales by group and state
2 pivot_table1 = df.pivot_table(index="Group", columns="State", values="Sales", aggfunc="sum")
3 # Example 2: Sum of sales for time of day
4 pivot_table2 = df.pivot_table(index="Time", values="Sales", aggfunc="sum")
5 print(pivot_table1)
6 print(pivot_table2)
```

State	NSW	NT	QLD	SA	TAS	VIC	WA
Group							
Kids	18587500	5700000	8510000	14515000	5775000	26360000	5625000
Men	19022500	5762500	8392500	14655000	5757500	26407500	5752500
Seniors	18187500	5465000	8190000	14717500	5650000	26315000	5512500
Women	19172500	5652500	8325000	14970000	5577500	26482500	5262500
Sales							
Time							
Afternoon	114007500						
Evening	112087500						
Morning	114207500						

✚ Data Analysis

```

1 # Descriptive statistics
2 print(df[['Unit', 'Sales']].describe())
3
4 # Top and bottom sales groups/states
5 top_group = grouped_data["Sales"].mean().idxmax()
6 bottom_group = grouped_data["Sales"].mean().idxmin()
7
8 top_state = df[df["Sales"] == df["Sales"].max()]["State"].values[0]
9 bottom_state = df[df["Sales"] == df["Sales"].min()]["State"].values[0]
10 # Which group is generating highest sales
11 # Total sales by state
12 pivot_table3 = df.pivot_table(index="State", values="Sales", aggfunc="sum")
13 print('Table of state wise revenue')
14 print(pivot_table3)
15 # Which group is generating highest sale
16 pivot_table4 = df.pivot_table(index="Group", values="Sales", aggfunc="sum")
17 print('Table of group wise revenue')
18 print(pivot_table4)
19
20 # Time analysis done over timeline in Date
21 # To do timeline analysis set the date as index and then use the resample function
22 df = df.set_index("Date")
23
24 # Weekly, monthly, quarterly reports
25 weekly_sales = df.resample("W-Sun")["Sales"].sum()
26 monthly_sales = df.resample("M")["Sales"].sum()
27 quarterly_sales = df.resample("Q")["Sales"].sum()
28 print(f'The weekly sales are /n {weekly_sales}')
29 print(f'The monthly sales are /n {monthly_sales}') # Should give just 3 entries
30 print(f'The quarterly sales are /n {quarterly_sales}') # This data is only for one quarter

```

	Unit	Sales
count	7560.000000	7560.000000
mean	18.005423	45013.558201
std	12.901403	32253.506944
min	2.000000	5000.000000
25%	8.000000	20000.000000
50%	14.000000	35000.000000
75%	26.000000	65000.000000
max	65.000000	162500.000000

Table of state wise revenue

	Sales
--	-------

State	
-------	--

NSW	74970000
NT	22580000
QLD	33417500
SA	58857500
TAS	22760000
VIC	105565000
WA	22152500

Table of group wise revenue

	Sales
--	-------

Group	
-------	--

Kids	85072500
Men	85750000
Seniors	84037500
Women	85442500

The weekly sales are /n Date

2020-10-04	15045000
2020-10-11	27002500
2020-10-18	26640000
2020-10-25	26815000
2020-11-01	21807500
2020-11-08	20865000
2020-11-15	21172500
2020-11-22	21112500
2020-11-29	21477500
2020-12-06	29622500
2020-12-13	31525000
2020-12-20	31655000
2020-12-27	31770000
2021-01-03	13792500

Freq: W-SUN, Name: Sales, dtype: int64

The monthly sales are /n Date

2020-10-31	114290000
2020-11-30	90682500
2020-12-31	135330000

Freq: M, Name: Sales, dtype: int64

The quarterly sales are /n Date

2020-12-31	340302500
------------	-----------

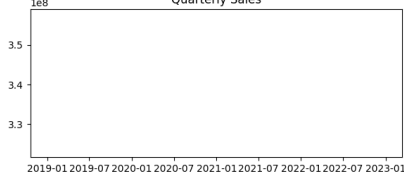
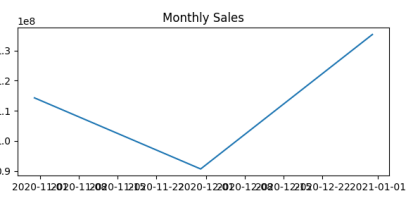
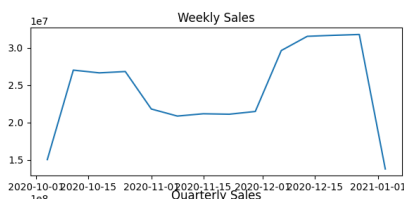
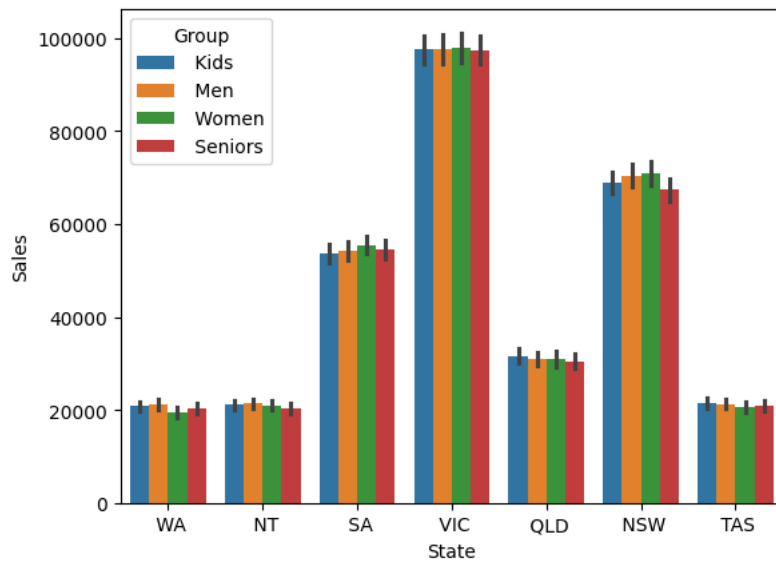
Freq: Q-DEC, Name: Sales, dtype: int64

✓ Data Visualisation

- State-wise sales analysis for different groups (kids, women, men, and seniors)
- Group-wise sales analysis (kids, women, men, and seniors) across different states.
- Time-of-the-day analysis: during which time of the day are sales the highest, and during which time are sales the lowest?

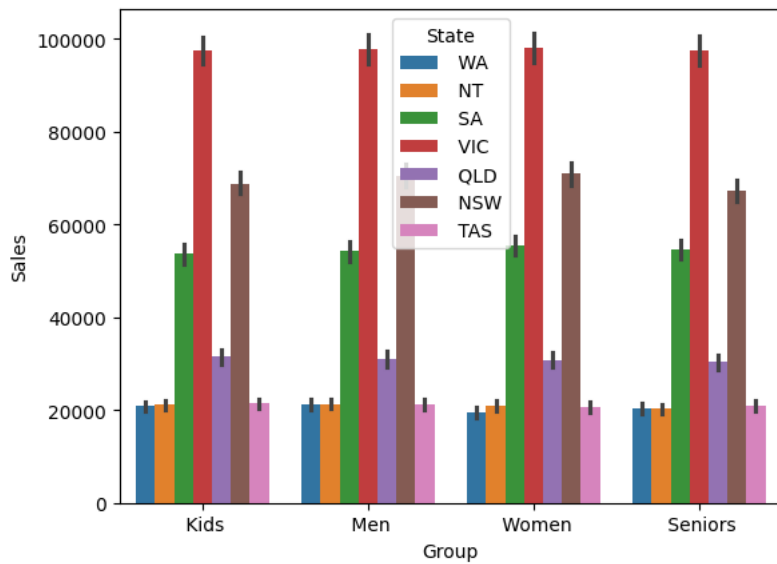
```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # State-wise sales analysis for different groups
5 sns.barplot(x="State", y="Sales", hue="Group", data=df)
6
7
8
9 # Time-of-day analysis
10 # ... use appropriate libraries and data transformation for analysis
11
12 # Daily, weekly, monthly, quarterly charts
13 plt.figure(figsize=(15, 6))
14 plt.subplot(2, 2, 1)
15 plt.plot(weekly_sales)
16 plt.title("Weekly Sales")
17 plt.subplot(2, 2, 2)
18 plt.plot(monthly_sales)
19 plt.title("Monthly Sales")
20 plt.subplot(2, 2, 3)
21 plt.plot(quarterly_sales)
22 plt.title("Quarterly Sales")# Single value wont show
23
24
```

Text(0.5, 1.0, 'Quarterly Sales')



✓ Group wise sale

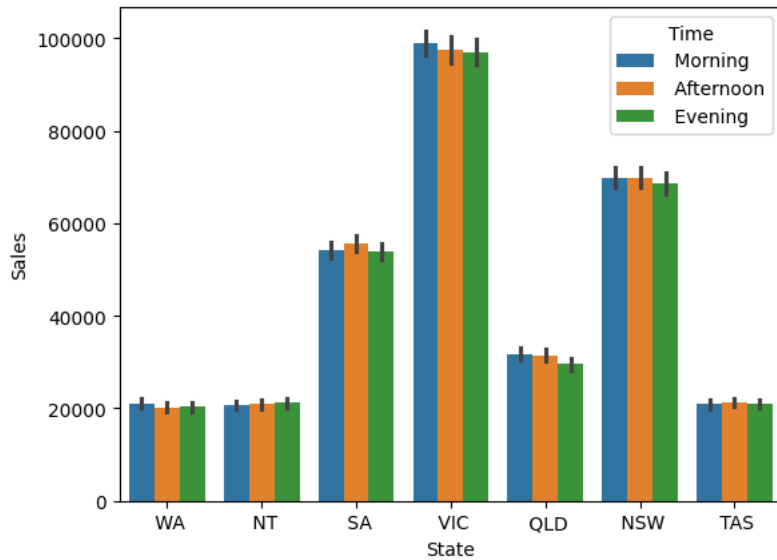
```
1 # Group-wise sales analysis across different states
2 plot1=sns.barplot(x="Group", y="Sales", hue="State", data=df)
```



✓ Impact of time on sales

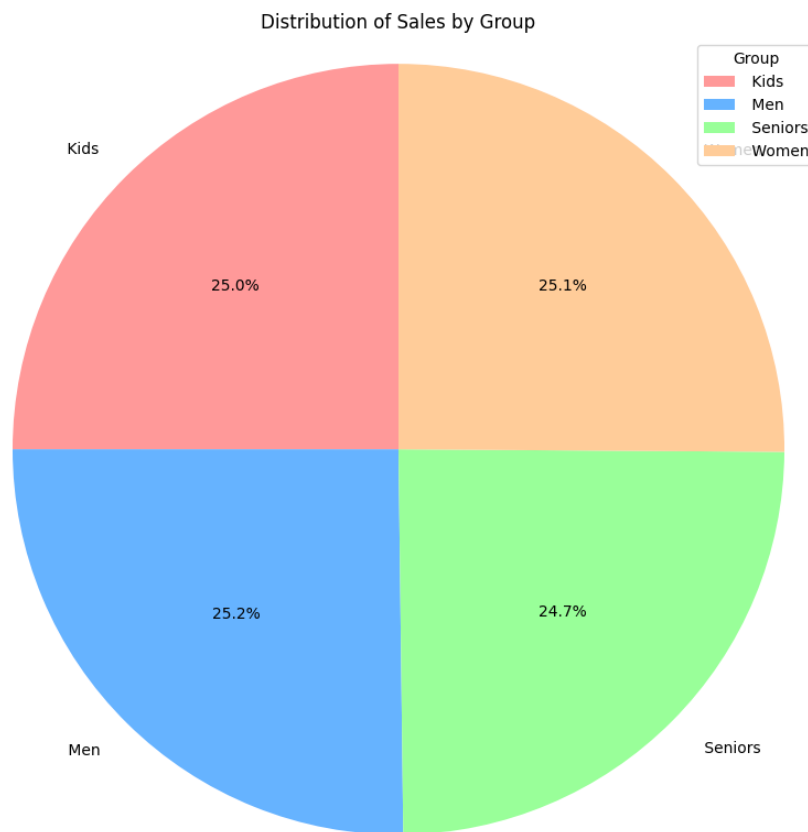
```
1 sns.barplot(x="State", y="Sales", hue="Time", data=df)
```

<Axes: xlabel='State', ylabel='Sales'>



Group sales Pie chart

```
1 # Create a DataFrame showing total sales per group
2 group_sales = df.groupby("Group")["Sales"].sum().reset_index()
3
4 # Create the pie chart
5 # Define colors for each group (optional)
6 colors = ['#ff9999', '#66b3ff', '#99ff99', '#ffcc99']
7
8 # Create the pie chart
9 plt.figure(figsize=(8, 8))
10 plt.pie(group_sales["Sales"], labels=group_sales["Group"], autopct="%1.1f%%", startangle=90, colors=colors)
11 plt.title("Distribution of Sales by Group")
12 plt.axis("equal") # Equal aspect ratio ensures a circular pie chart
13
14 # Customize the chart (optional)
15 plt.legend(title="Group")
16 plt.tight_layout()
17
18 # Show the chart
19 plt.show()
```



✓ Report Generation

- Libraries for visualisation I am using seaborn lib for bar charts and matplotlib for the line plots
- Data Analysis State(VIC) contributes to more than 30% revenue and state WA is about 6%
- Unit sales reflect the total sales pattern
- There is no impact of time on the overall sales. All are almost equal.
- Histogram of sales show max sales are in the range 4000-5000. This is a right skewed data. Which is standard for sales
- There is insignificant variation between the sales figures of the group

1

✓ Visualisation of descriptive data

The code generates descriptive stats for All numeric data in the df


```

1 # This is a code sinp from my previous project repurposed for analysis
2 import scipy.stats as stats
3 df.columns
4
5 # Identify numerical columns for analysis
6 #Exclude your binaries . they got reset to zero with automated Winsorization
7
8 #exclude_list = ['AcceptedCmp3', 'AcceptedCmp4', 'AcceptedCmp5', 'AcceptedCmp1',
9 #               'AcceptedCmp2', 'Response', 'Complain']
10 exclude_list=[]
11 numerical_columns = df.select_dtypes(include=['int64', 'float64']).columns.difference(exclude_list)
12
13 #numerical_columns = df.select_dtypes(include=['int64', 'float64']) ---- this didnt work
14 # Careful with outlier check on binary columns. It will set them to zero
15 # Create a function for outlier detection and treatment
16 def detect_and_treat_outliers(column):
17     """
18     Detects outliers in a numerical column using IQR and removes or winsorizes them.
19
20     Args:
21         column: The numerical column to analyze.
22
23     Returns:
24         The cleaned column with outliers removed or winsorized.
25     """
26
27     # Calculate IQR and quantiles
28     q1 = column.quantile(0.25)
29     q3 = column.quantile(0.75)
30     iqr = q3 - q1
31
32     # Identify potential outliers (values beyond 1.5 IQR from quartiles)
33     lower_bound = q1 - 1.5 * iqr
34     upper_bound = q3 + 1.5 * iqr
35     outliers = column[(column < lower_bound) | (column > upper_bound)]
36
37
38     # 1. Winsorize outliers (replace with nearest non-outlier value)
39     # Choosing winsorize with clip method as I want the shape of my data frame unchanged with no nulls.
40     # Another method we can use the matstats library the command is slick!
41     #from scipy.stats import mstats
42     #cleaned_column = mstats.winsorize(df[column_to_winsorize], limits=(lower_limit, upper_limit))
43
44     cleaned_column = column.clip(lower_bound, upper_bound,axis=0)
45
46     return cleaned_column
47 clean_df=pd.DataFrame()
48 # Iterate through numerical columns and visualize distributions
49 #Basic plot for every column descriptive anaysis.
50 #Eyeballing the data is much better I feel and faster. Gimmick
51 for column in numerical_columns:
52     print(f"\nAnalyzing column: {column}")
53
54     # Create box plot
55     plt.figure(figsize=(8, 6))
56     df[column].plot(kind='box')
57     plt.title(f"Box Plot for {column}")
58     plt.show()
59
60     # Create histogram
61     plt.figure(figsize=(8, 6))
62     df[column].hist(bins=20)
63     plt.title(f"Histogram for {column}")
64     plt.show()
65
66     # Detect and treat outliers (use the chosen method from the function)
67
68     clean_df[column] = detect_and_treat_outliers(df[column])
69     #df[column] = detect_and_treat_outliers(df[column])
70 clean_df.info()
71 clean_df.describe()
72 df.describe()
73
74 # Save the cleaned DataFrame if needed
75 # cleaned_df.to_csv('cleaned_data.csv', index=False)
76
77

```

Analyzing column: Unit

