

Import libraries

```
In [1]: import pandas as pd  
In [2]: trnx = pd.read_csv("Retail_Data_Transactions.csv")  
response = pd.read_csv("Retail_Data_Response.csv")  
In [3]: trnx.head(5)  
Out[3]: customer_id trans_date tran_amount  
0 CS5295 2013-02-11 35  
1 CS4768 2015-03-15 39  
2 CS2122 2013-02-26 52  
3 CS1217 2011-11-16 99  
4 CS1850 2013-11-20 78
```

```
In [4]: response.head(5)
```

```
Out[4]: customer_id response  
0 CS1112 0  
1 CS1113 0  
2 CS1114 1  
3 CS1115 1  
4 CS1116 1
```

Merging the tables

```
In [5]: df = trnx.merge(response,on="customer_id",how="left")
```

```
In [6]: df.head(3)
```

```
Out[6]: customer_id trans_date tran_amount response  
0 CS5295 2013-02-11 35 1.0  
1 CS4768 2015-03-15 39 1.0  
2 CS2122 2013-02-26 52 0.0
```

Check nulls values

```
In [7]: df.isnull().sum()
```

```
Out[7]: customer_id 0  
trans_date 0  
tran_amount 0  
response 31  
dtype: int64
```

Check Data types

```
In [8]: df.dtypes
```

```
Out[8]: customer_id object  
trans_date object  
tran_amount int64  
response float64  
dtype: object
```

```
In [9]: df["trans_date"] = pd.to_datetime(df["trans_date"])
```

```
In [10]: df["response"] = df["response"].astype("object")
```

```
In [11]: df.dtypes
```

```
Out[11]: customer_id object  
trans_date datetime64[ns]  
tran_amount int64  
response object  
dtype: object
```

```
In [12]: df["response"] = df["response"].fillna(0).astype("int64")
```

```
C:\Users\HP\AppData\Local\Temp\ipykernel_6472\101564169.py:1: FutureWarning: Downcasting object dtype arrays on .fillna, .ffill, .bfill is deprecated and will change in a future version. Call result.infer_objects(copy=False) instead. To opt-in to the future behavior, set pd.set_option('future.no_silent_downcasting', True)
```

```
df["response"] = df.fillna(0).astype("int64")
```

```
In [13]: df.dtypes
```

```
Out[13]: customer_id object  
trans_date datetime64[ns]  
tran_amount int64  
response int64  
dtype: object
```

Checking Outliers using z_score

```
In [14]: import numpy as np  
from scipy import stats
```

```
In [15]: z_score = np.abs(stats.zscore(df["tran_amount"]))
```

```
threshold = 3
```

```
outliers = z_score > threshold
```

```
print(outliers)
```

```
[False False False ... False False False]
```

```
In [16]: z_score = np.abs(stats.zscore(df["response"]))
```

```
threshold = 3
```

```
outliers = z_score > threshold
```

```
print(outliers)
```

```
[False False False ... False False False]
```

Create new columns

```
In [17]: df["month"] = df["trans_date"].dt.month_name()
```

```
df["year"] = df["trans_date"].dt.year
```

```
In [18]: df.head(3)
```

```
Out[18]: customer_id trans_date tran_amount response month year
```

```
0 CS5295 2013-02-11 35 1 February 2013
```

```
1 CS4768 2015-03-15 39 1 March 2015
```

```
2 CS2122 2013-02-26 52 0 February 2013
```

```
In [19]: import plotly.express as px
```

Which three months have the highest amount of transactions

```
In [20]: month_sales = df.groupby("month")["tran_amount"].sum().sort_values(ascending=False).reset_index()
```

```
In [21]: month_sales.head(3)
```

```
Out[21]: month tran_amount
```

```
0 August 726921
```

```
1 October 725320
```

```
2 January 724107
```

```
In [22]: fig = px.bar(month_sales,
```

```
x="month",
y="tran_amount",
title="Highest Month Transaction",
color="tran_amount",
color_continuous_scale="Cividis",
)
```

```
fig.update_layout(width=700,height=400)
fig.show()
```

Top 5 customer have highest co of orders

```
In [23]: customer_count = df["customer_id"].value_counts().reset_index()
```

```
customer_count.head(5)
```

```
Out[23]: customer_id count
```

```
0 CS4424 39
```

```
1 CS4320 38
```

```
2 CS3799 36
```

```
3 CS5109 35
```

```
4 CS2620 35
```

```
In [24]: fig = px.bar(customer_count.head(10),
x="customer_id",
y="count",
title="Top 5 Customers By Order Count",
color="count",
color_continuous_scale="Cividis",
)
fig.update_layout(width=500,height=400)
fig.show()
```

Top 10 customer having highest sales

```
In [25]: customer_sales = df.groupby("customer_id")["tran_amount"].sum().reset_index().sort_values(by="tran_amount",ascending=False)
```

```
Out[25]: customer_id tran_amount
```

```
3312 CS5424 2933
```

```
3208 CS4320 2647
```

```
4640 CS5752 2612
```

```
3548 CS4660 2527
```

```
2687 CS3799 2513
```

```
3997 CS5109 2506
```

```
2962 CS4074 2462
```

```
2693 CS3805 2453
```

```
3496 CS4608 2449
```

```
4443 CS5555 2439
```

```
In [26]: fig = px.bar(customer_sales.head(10),
x="customer_id",
y="tran_amount",
title="Highest Customers Sales",
color="tran_amount",
color_continuous_scale="Cividis",
)
fig.update_layout(width=700,height=400)
fig.show()
```

Which year has the highest number of transactions

```
In [27]: yearly_sales = df.groupby("year")["tran_amount"].sum().reset_index().sort_values(by="tran_amount",ascending=False)
```

```
In [28]: yearly_sales
```

```
Out[28]: year tran_amount
```

```
2 2013 213768
```

```
1 2012 2116599
```

```
3 2014 2094508
```

```
0 2011 1340339
```

```
4 2015 435175
```

```
In [29]: fig = px.bar(yearly_sales,
x="year",
y="tran_amount",
title="Monthly Transaction Trend",
color="tran_amount",
color_discrete_sequence="Cividis",
)
fig.update_layout(width=700,height=400)
fig.show()
```

Advanced Analytics

Times Series Analysis

```
In [30]: df["month_year"] = df["trans_date"].dt.to_period("M").dt.to_timestamp()
```

```
month_year_sales = df.groupby("month_year")["tran_amount"]
```

```
.reset_index()
```

```
.sort_values("month_year")
```

```
)
```

```
fig = px.line(month_year_sales,
x="month_year",
y="tran_amount",
title="Monthly Transaction Trend",
color_discrete_sequence="olive",
markers=True
)
```

```
fig.update_xaxes(
    tick="M",
    tickformat="%b-%Y",
    tickangle=45
)
fig.show()
```

```
Monthly Transaction Trend
```



```
In [32]: df.head(3)
```

```
Out[32]: customer_id trans_date tran_amount response month year month_year
```

```
0 CS5295 2013-02-11 35 1 February 2013 2013-02-01
```

```
1 CS4768 2015-03-15 39 1 March 2015 2015-03-01
```

```
2 CS2122 2013-02-26 52 0 February 2013 2013-02-01
```

```
Customer Segmentation
```

```
In [33]: # Recency - most recent orders
```

```
recency = df.groupby("customer_id")["trans_date"].max()
```

```
# Frequency - no. of orders
```

```
frequency = df.groupby("customer_id")["trans_date"].count()
```

```
# Monetory - sum of amount
```

```
monetary = df.groupby("customer_id")["tran_amount"].sum()
```

```
# combine all into DataFrame
```

```
rfm = pd.DataFrame([{"recency":recency,
"frequency":frequency,
"monetary":monetary}])
```

```
fig = px.bar(rfm,
x="customer_id",
y="monetary",
title="Top 5 Customers By Order Count",
color="monetary",
color_continuous_scale="Cividis",
)
fig.update_layout(width=500,height=400)
fig.show()
```

```
Top 5 customer have highest co of orders
```

```
In [34]: customer_count = df["customer_id"].value_counts().reset_index()
```

```
customer_count.head(5)
```

```
Out[34]: customer_id count
```

```
0 CS4424 39
```

```
1 CS4320 38
```

```
2 CS3799 36
```

```
3 CS5109 35
```

```
4 CS2620 35
```

```
In [35]: fig = px.bar(customer_count.head(10),
x="customer_id",
y="count",
title="Top 5 Customers By Order Count",
color="count",
color_continuous_scale="Cividis",
)
fig.update_layout(width=500,height=400)
fig.show()
```

```
Top 5 customer have highest co of orders
```

```
In [36]: customer_sales = df.groupby("customer_id")["tran_amount"].sum().reset_index()
```

```
customer_sales.head(5)
```

```
Out[36]: customer_id count
```

```
3312 CS5424 2933
```

```
3208 CS4320 2647
```

```
4640 CS5752 2612
```

```
3548 CS4660 2527
```

```
2687 CS3799 2513
```

```
3997 CS5109 2506
```

```
2962 CS4074 2462
```

```
2693 CS3805 2453
```

```
3496 CS4608 2449
```

```
4443 CS5555 2439
```

```
In [37]: fig = px.bar(customer_sales.head(10),
x="customer_id",
y="count",
title="Top 5 Customers By Order Count",
color="count",
color_continuous_scale="Cividis",
)
fig.update_layout(width=500,height=400)
fig.show()
```

```
Top 5 customer have highest co of orders
```

```
In [38]: rfm["segment"] = rfm.apply(seg_customer,axis=1)
```

```
Out[38]: rfm
```

```
Out[38]: recency frequency monetary segment
```

```
customer_id
```

```
CS1112 2015-01-04 15 1012 P0
```

```
CS1113 2015-02-09 20 1490 P0
```

```
CS1114 2015-02-12 19 1432 P0
```

```
CS1115 2015-03-05 22 1659 P0
```

```
CS1116 2014-08-25 13 857 P1
```

```
...
```

```
CS
```

Top 5 Customers Monthly Sales Trend



In []: