

Data Analysis Final Project

Shopping Tendencies Analysis

By:
Sahar Elersawi

Introduction

- ❑ Our project is talking about the analysis of sales in 10 malls in Turkey.
- ❑ We used different analysis methodologies and techniques to analyze our dataset like we used
- ❑ Python coding and MS Power Bi visualization and reports.
- ❑ We set 4 questions to use its answer for good analysis to be shown clearly to the stakeholders :

1- Which gender frequency buy more ? And What is the average age?

2- What is the sales volume of each category? Which category has the majority of sales volume?

3- What is the sales volume for each mall ? Which mall has the majority sales volume?

4- What is the most common payment method?

##loading data

#read the dataset csv file

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns
```

```
import numpy as np
```

```
df=pd.read_csv("/Users/saharsayed/Downloads/customer_shopping_data  
2.csv")
```

```
# Display the first few rows of the dataset
```

```
#print(df.head())
```

```
df.head()
```

	invoice_no	customer_id	gender	age	category	quantity	price	payment_method	invoice_date	shopping_mall
0	I138884	C241288	Female	28	Clothing	5	1500.40	Credit Card	5/8/2022	Kanyon
1	I317333	C111565	Male	21	Shoes	3	1800.51	Debit Card	12/12/2021	Forum Istanbul
2	I127801	C266599	Male	20	Clothing	1	300.08	Cash	9/11/2021	Metrocity
3	I173702	C988172	Female	66	Shoes	5	3000.85	Credit Card	16/05/2021	Metropol AVM
4	I337046	C189076	Female	53	Books	4	60.60	Cash	24/10/2021	Kanyon

#Clean the data

##finding missing null

```
valuesdf.isnull().sum()
```

Drop rows with missing values

```
df=df.dropna()
```

Drop duplicate

```
rowsdf = df.drop_duplicates()
```

Understanding the Dataset

#understanding the dataset

get first 5 rows

rowsdf.head()

##data processing and # check the content

df.shapedf.describe()

	age	quantity	price
count	99457.000000	99457.000000	99457.000000
mean	43.427089	3.003429	689.256321
std	14.990054	1.413025	941.184567
min	18.000000	1.000000	5.230000
25%	30.000000	2.000000	45.450000
50%	43.000000	3.000000	203.300000
75%	56.000000	4.000000	1200.320000
max	69.000000	5.000000	5250.000000

Checking number of unique entries

Checking number of unique entries

df.nunique()

invoice_no.	99457
customer_id	99457
gender	2
age	52
category	8
quantity	5
price	40
payment_method	3
invoice_date	797
shopping_mall	10
year	3
month	12 dtype: int64

1- Which gender buy more ?

Most frequent entries and their frequencies for 'Gender'

```
gender_counts = df['gender'].value_counts()  
print("Gender Counts:")  
print(gender_counts)
```

Gender Counts:

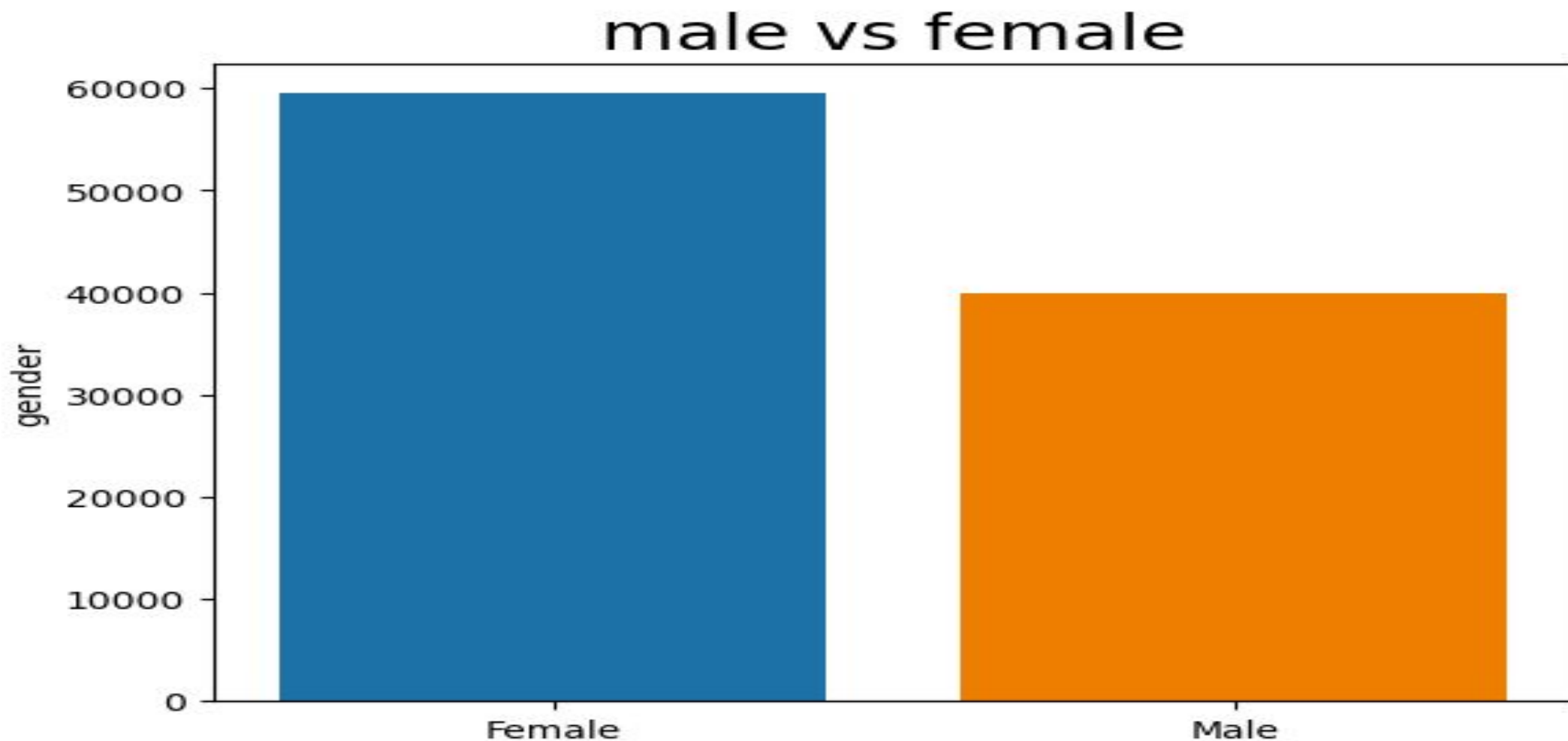
Female 59482

Male 39975

Name: gender, dtype: int64

- ❑ **Females buy with about 20% more than males.**


```
In [24]: gen_num=df["gender"].value_counts()  
sns.barplot(y = gen_num, x = gen_num.index, data = df)  
plt.title("male vs female", size=20)  
  
Out[24]: Text(0.5, 1.0, 'male vs female')
```



```
In [25]: ##Majority of customers are female
```

FileHomeInsertModelingViewOptimizeHelpFormatData / Drill

PasteCutCopyFormat painterClipboard

Get dataExcel workbookData hubSQL ServerEnter dataData

Transform dataRefresh dataQueries

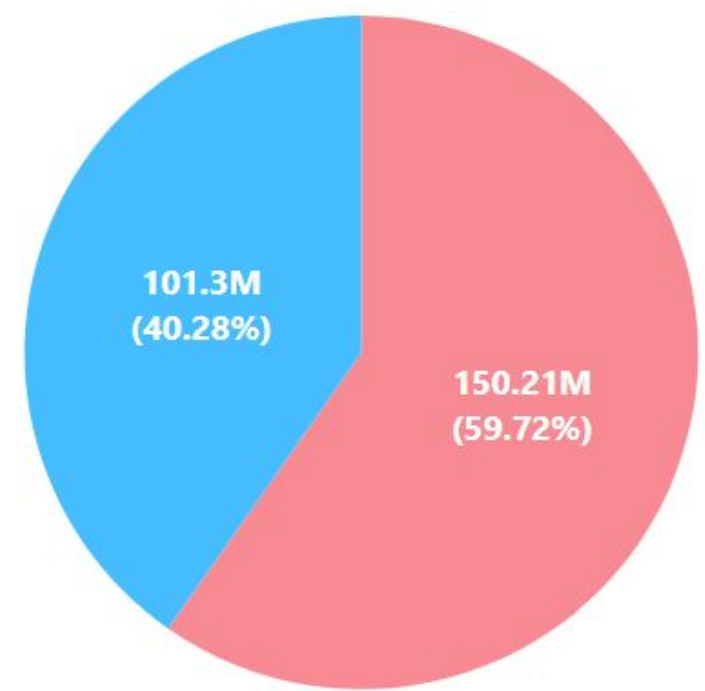
New visualText boxMore visualsInsert

New measureQuick measureCalculations

SensitivitySensitivity

PublishShare

Back to report | SUM OF TOTAL_PRICE BY GENDER



gender
● Female
● Male

What is the average age?

```
# Calculate average age
```

```
average_age = df['age'].mean()
```

```
print('Average age')
```

```
print (average_age)
```

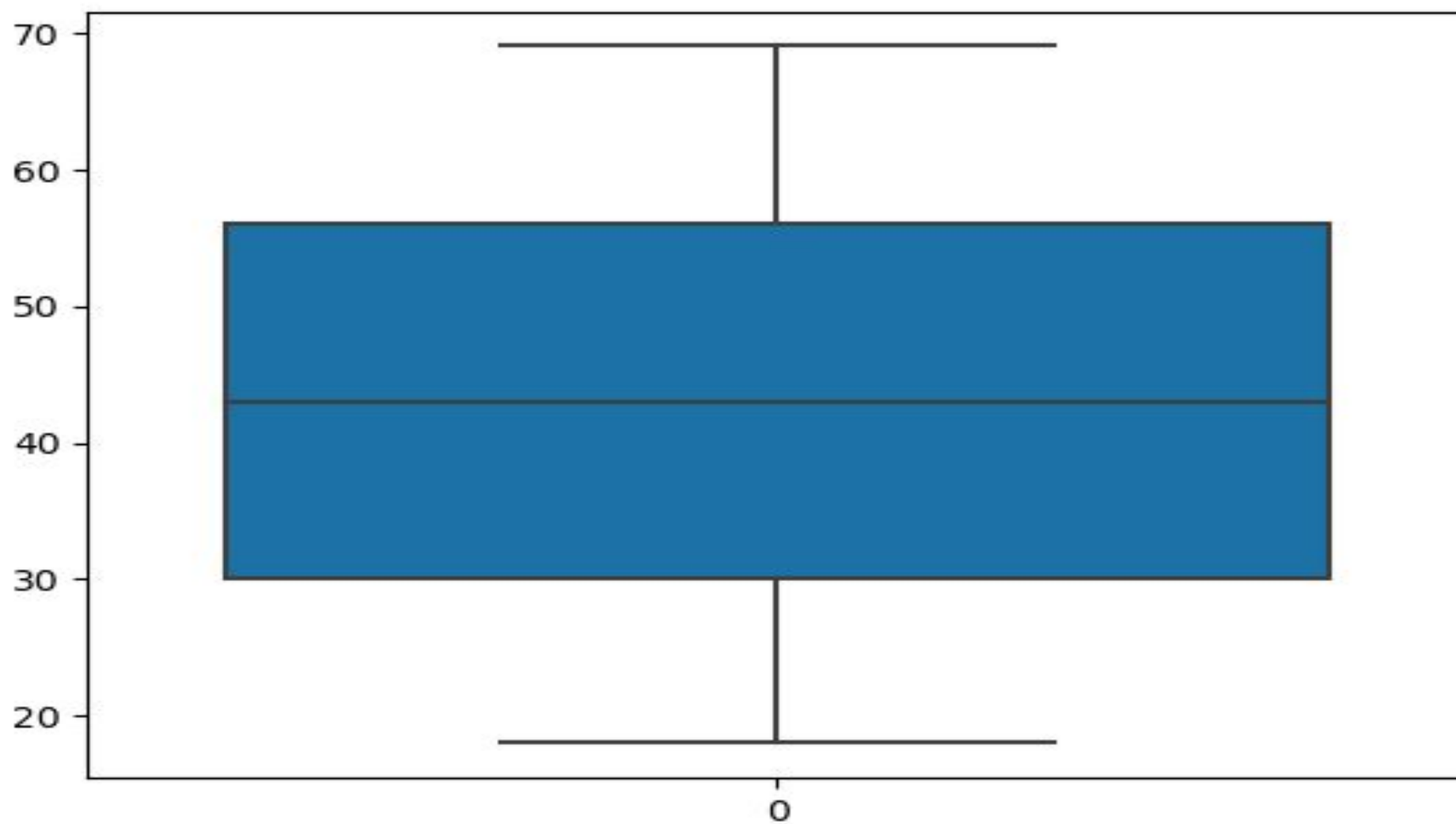
Average age

43.42708909377922

❑ **Average age (mean) is 43 years old.**

```
In [22]: sns.boxplot(df['age'])
```

```
Out[22]: <Axes: >
```



```
In [23]: #So, average customer age is between 30-55 years
```

2- What is the sales volume of each category? Which category has the majority of sales volume?

Group by category and calculate sales volume

```
category_sales = df.groupby("category")[["quantity", "price"]].sum()
category_sales["sales_volume"] = category_sales["quantity"] * category_sales["price"]
print("SALES VOLUME OF EACH CATEGORY")
print(category_sales["sales_volume"])
```

Identify category with the majority sales volume

```
majority_category = category_sales["sales_volume"].idxmax()
print("Category with the majority sales volume:", majority_category)
```

SALES VOLUME OF EACH CATEGORY

category

Books	3.400574e+09
Clothing	3.218136e+12
Cosmetics	8.404691e+10
Food & Beverage	1.025317e+10
Shoes	5.479955e+11
Souvenir	2.594050e+09
Technology	2.369120e+11
Toys	3.294997e+10

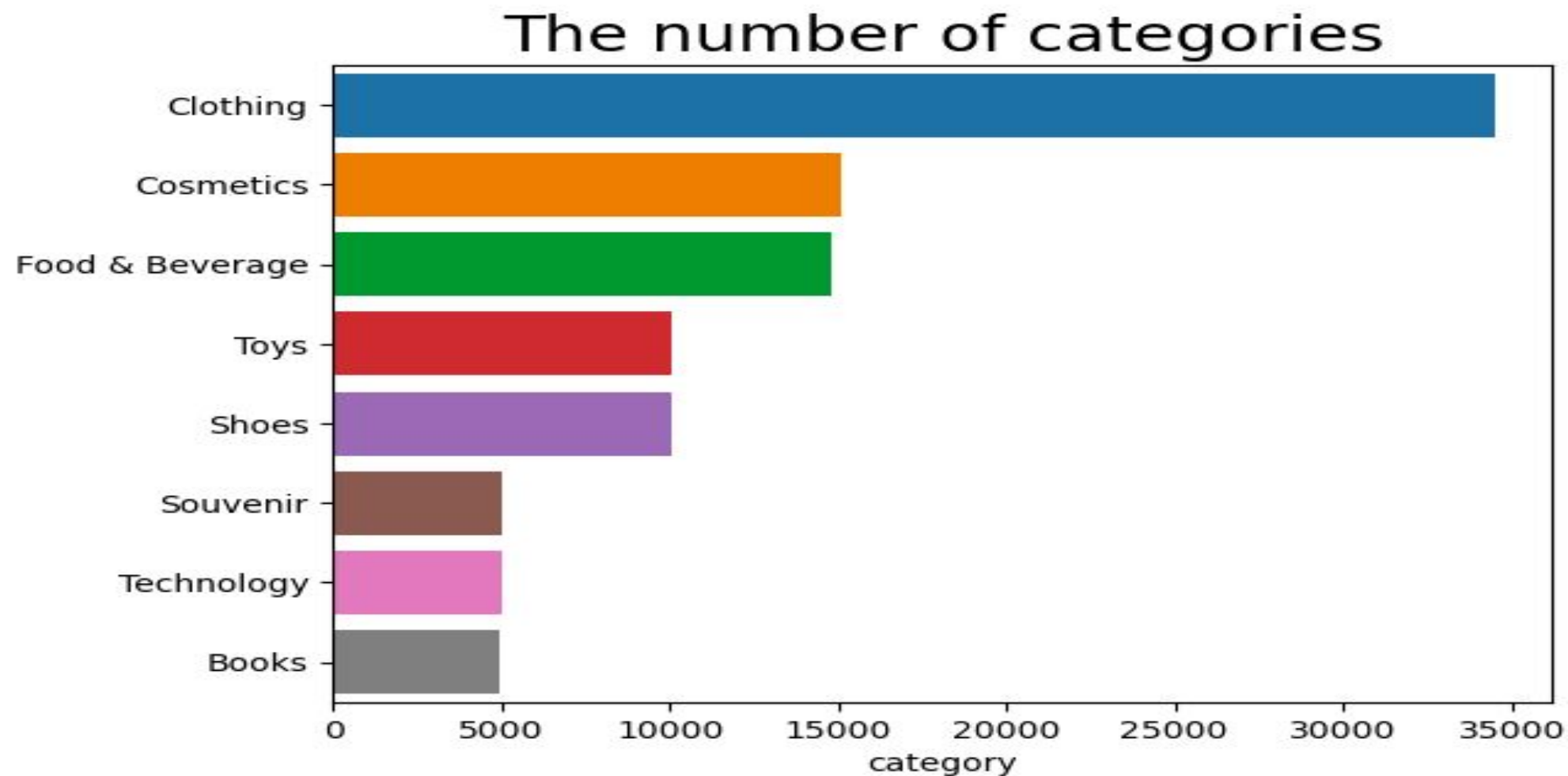
Name: sales_volume, dtype: float64

Category with the majority sales volume: Clothing

- ❑ Clothes and shoes have the majority of sales volume and books and Souvenir have the minority of sales volume.

```
In [23]: cat_num = df["category"].value_counts()  
sns.barplot(x = cat_num, y = cat_num.index, data = df)  
plt.title("The number of categories", size=20)
```

```
Out[23]: Text(0.5, 1.0, 'The number of categories')
```



```
In [24]: ##There are total 8 categories of products and maximum sales by quantity is clothes .
```

FileHomeInsertModelingViewOptimizeHelpFormatData / Drill

PasteCutCopyFormat painterClipboard

Get dataExcel workbookData hubSQL ServerEnter dataDataverseRecent sourcesData

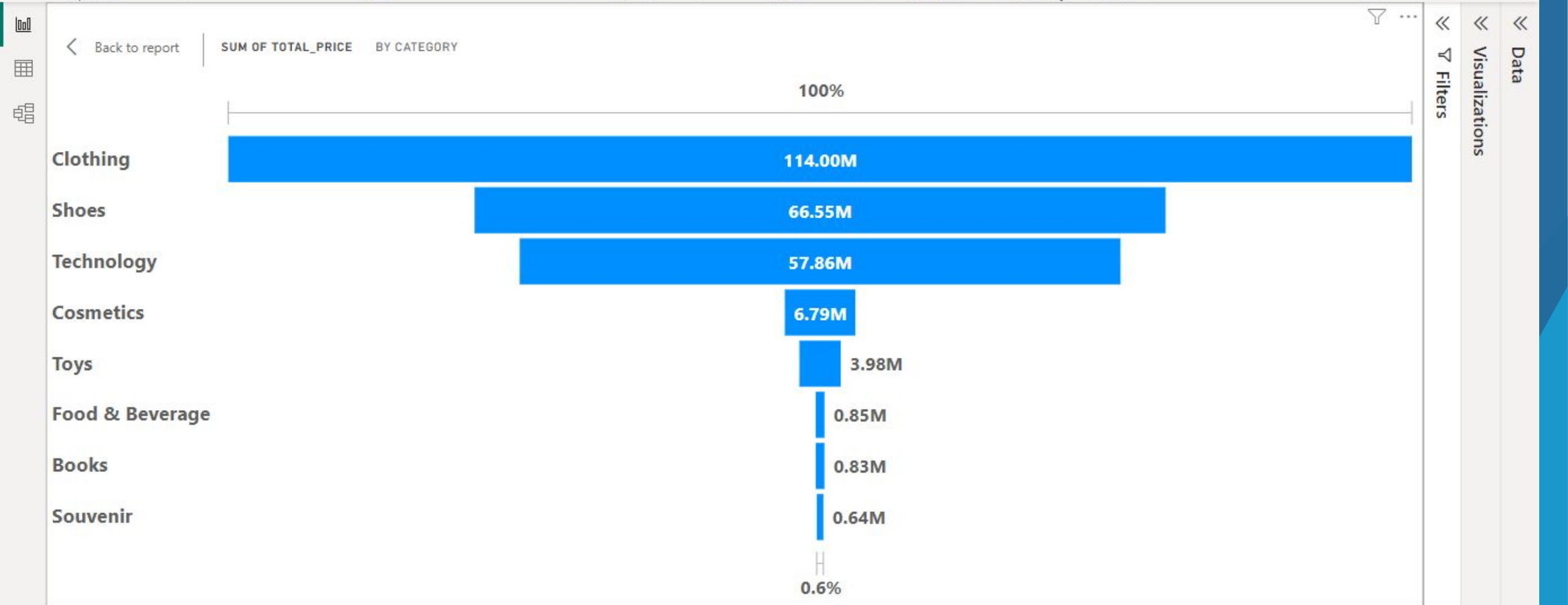
Transform dataRefresh dataQueries

New visualText boxMore visualsInsert

New measureQuick measureCalculations

SensitivitySensitivity

PublishShare



3- What is the sales volume for each mall ? Which mall has the majority sales volume?

Group by shopping mall and calculate sales volume

```
mall_sales = df.groupby("shopping_mall")[["quantity", "price"]].sum()  
mall_sales["sales_volume"] = mall_sales["quantity"] * mall_sales["price"]  
print("SALES VOLUME OF EACH MALL:")  
print(mall_sales["sales_volume"])
```

#identify mall with the majority sales volume

```
mall_max_sales = mall_sales["sales_volume"] . idxmax()  
print ("MALL WITH MAJORITY SALES:")  
print (mall_max_sales)
```

SALES VOLUME OF EACH MALL:

shopping_mall

Cevahir AVM	5.132996e+10
Emaar Square Mall	4.916431e+10
Forum Istanbul	4.954737e+10
Istinye Park	1.979187e+11
Kanyon	8.152004e+11
Mall of Istanbul	8.326834e+11
Metrocity	4.601626e+11
Metropol AVM	2.118169e+11
Viaport Outlet	5.024071e+10
Zorlu Center	5.346599e+10

Name: sales_volume, dtype: float64

MALL WITH MAJORITY SALES:

Mall of Istanbul

- ☐ Mall of Istanbul has the majority of sales.
- ☐ Emaar square mall has the minority of sales.

FileHomeInsertModelingViewOptimizeHelpFormatData / Drill

PasteCutCopyFormat painterClipboard

Get dataExcel workbookData hubSQL ServerEnter dataDataaverseRecent sourcesData

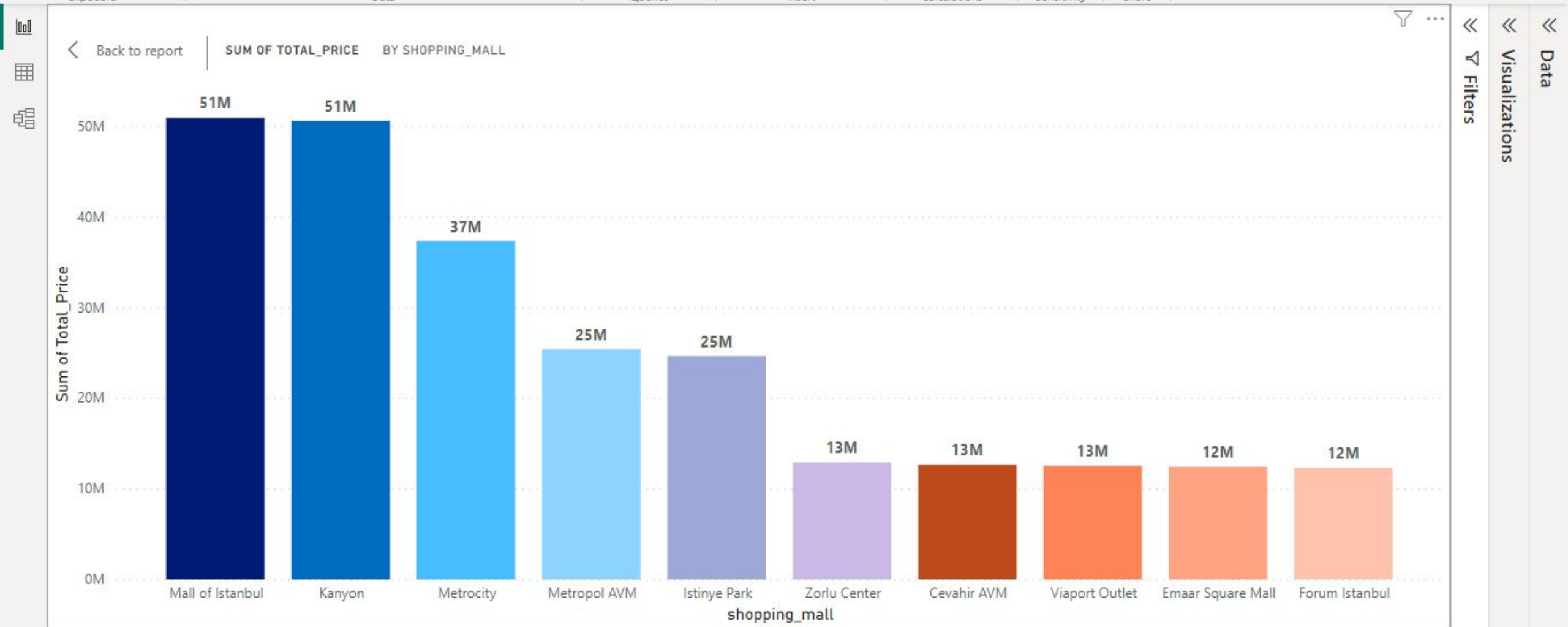
Transform dataRefresh dataQueries

New visualText boxMore visualsInsert

New measureQuick measureCalculations

SensitivitySensitivity

PublishShare



4- What is the most common payment method?

```
#### Most frequent entries and their frequencies for 'Payment_Method'
payment_counts = df['payment_method'].value_counts()
print("\nPayment Method Counts:")
print(payment_counts)

## identify the common payment method
most_common_payment_method = payment_method_frequency.idxmax()
print('MOST COMMON PAYMENT METHOD :')
print(MOST_COMMON_PAYMENT_METHOD)
```

Payment Method Counts:

Cash 44447

Credit Card 34931

Debit Card 20079

Name: payment_method, dtype: int64

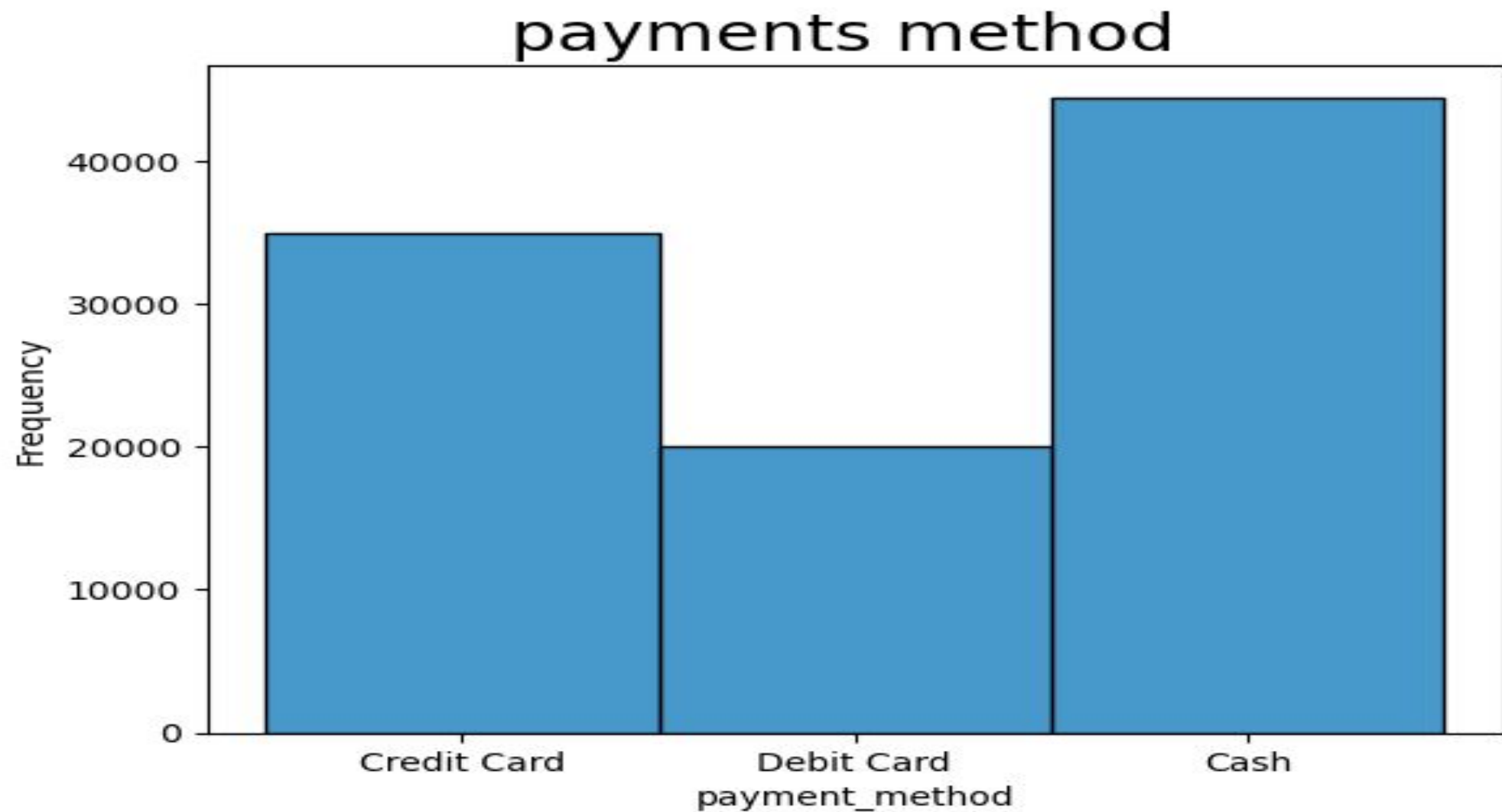
MOST COMMON PAYMENT METHOD:

CASH

☐ **Cash is the common payment method**

```
In [26]: sns.histplot(df["payment_method"], kde = False, stat='frequency')  
plt.title("payments method ", size=20,)
```

```
Out[26]: Text(0.5, 1.0, 'payments method ')
```



```
In [27]: ## Cash is the most prefferd mode of transaction followed
```

FileHomeInsertModelingViewOptimizeHelpFormatData / Drill

PasteCutCopyFormat painterClipboard

Get dataExcel workbookData hubSQL ServerEnter dataDataaverseRecent sourcesData

Transform dataRefresh dataQueries

New visualText boxMore visualsInsert

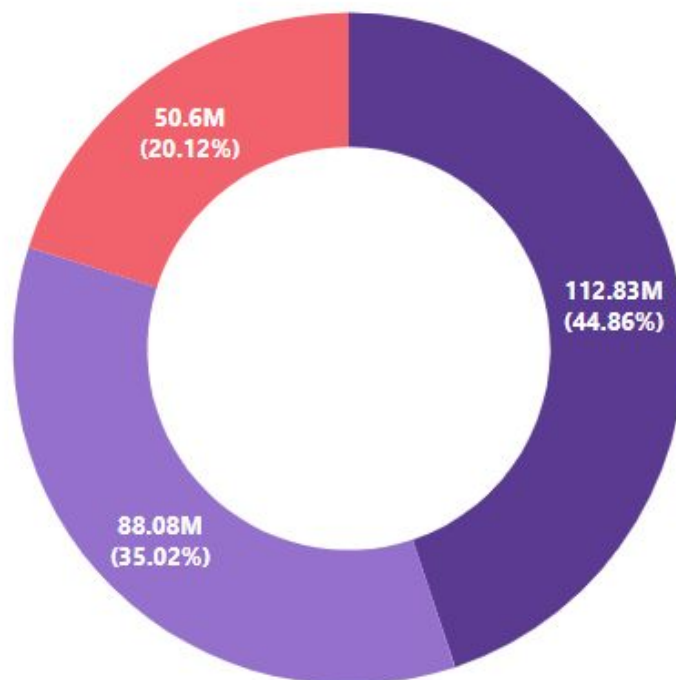
New measureQuick measureCalculations

SensitivitySensitivity

PublishShare

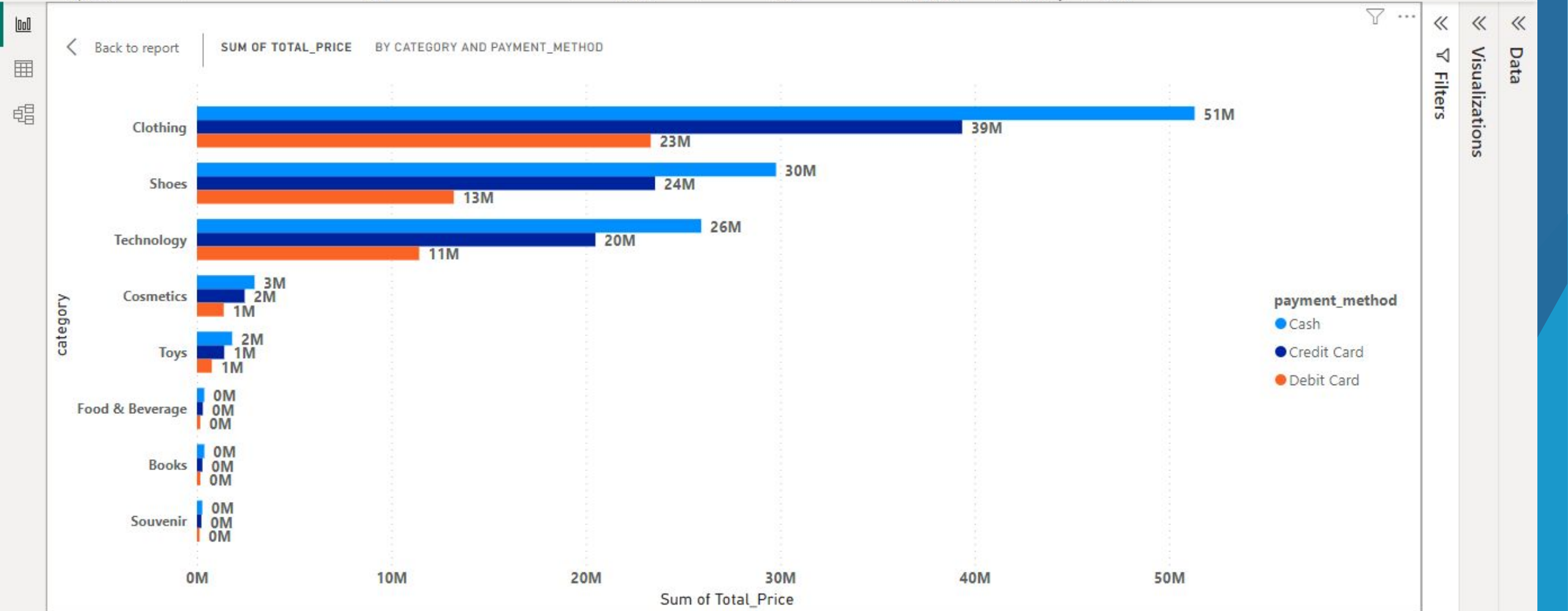
Back to report

SUM OF TOTAL_PRICE BY PAYMENT_METHOD



payment_method

- Cash
- Credit Card
- Debit Card



FileHomeInsertModelingViewOptimizeHelpFormatData / Drill

PasteCutCopyFormat painterClipboard

Get dataExcel workbookData hubSQL ServerEnter dataDataaverseRecent sourcesData

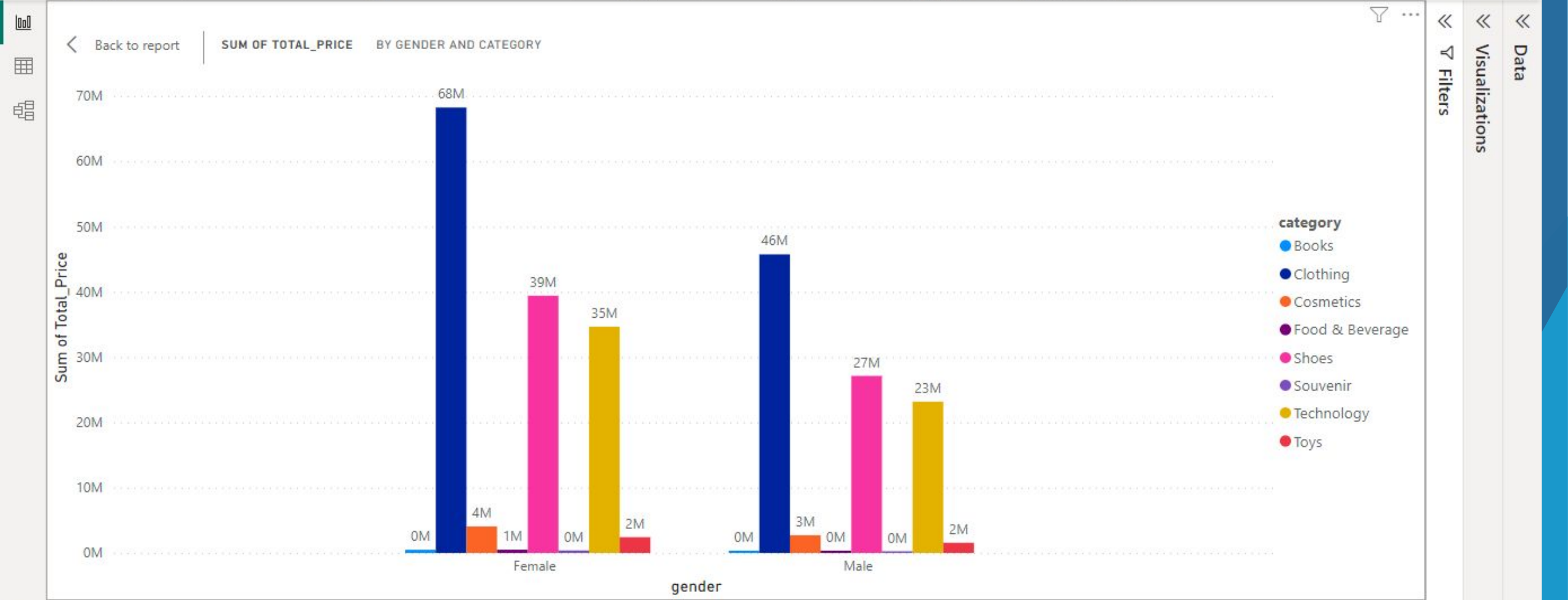
Transform dataRefreshQueries

New visualText boxMore visualsInsert

New measureQuick measureCalculations

SensitivitySensitivity

PublishShare



FileHomeInsertModelingViewOptimizeHelp

PasteCutCopyFormat painterClipboard

Get dataExcelDataSQLServerEnter dataDatawarehouseRecent sourcesData

Transform dataRefreshQueries

New visualText boxMore visualsInsert

New measureQuick measureCalculations

SensitivitySensitivity

PublishShare



Conclusion

Data Analytics recommendations

- ❑ The malls of minority of sales should make some marketing activities to increase its sales.
- ❑ Make promotions for the categories that has minority of sales to attract the customers and increase its sales.
- ❑ We recommend to make application to get the customer feedback in each mall for each category to enhance the services and the quality of the products.
- ❑ To improve the using of the credit cards the banks who owned the credit cards can make points program for each transaction paid by it and get free gifts from the minority categories of sales.

Thank You