

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')

In [2]: df = pd.read_excel("SuperStoreCanada (2).xlsx")
df

Out[2]:
```

	RefID	OrderDate	OrderID	Priority	ShipDate	ShipMode	CustID	CustomerSegment	CustName	PostalCode	ProdCategory	ProdSubCat	ProName	ProContainer	Quantity	SalePrice	CostPrice	Discount	ShipCost	Status
0	100001.0	2009-04-01	10001.0	Specified	Not Specified	Express Air	NA#N	Home Office	Max Customer	NA#N	Office Supplies	Storage & Organization	Safco Industrial Wire Shelving	Large Box	10.0	95.99	57.5940	0.08	35.00	OK
1	100002.0	2009-04-01	10002.0	High	2009-04-03	Regular Air	NA#N	Small Business	Jessica Myers	NA#N	Office Supplies	Storage & Organization	Penna STOR-ALL™ Hanging File Box, 13 1/8"W x 1...	Small Box	36.0	5.98	1.9136	0.10	4.69	OK
...
8394	103965.0	2013-03-30	15494.0	Not Specified	Not Specified	Regular Air	NA#N	Small Business	Jim Epp	NA#N	Furniture	Office Furnishings	DAX Wood Document Frame	Whip Bag	65.0	13.73	6.5904	0.08	8.85	OK
8395	103965.0	2013-03-30	15495.0	High	2013-04-01	Regular Air	NA#N	Home Office	Maribeth Vedwab	NA#N	Technology	Computer Peripherals	Bekins 105-Key Black Keyboard	Small Box	68.0	13.98	5.9504	0.00	4.00	OK
8396	103970.0	2013-03-30	15495.0	Low	2013-04-01	Regular Air	NA#N	Home Office	Maribeth Vedwab	NA#N	Office Supplies	Scissors, Rulers and Trimmers	Marin Vale Chaffers Opener Electric Letter Op...	Medium Box	1.0	832.81	199.8744	0.09	24.49	OK
8397	103980.0	2013-03-30	15496.0	Not Specified	Not Specified	Delivery Truck	NA#N	Corporate	Tony Molinari	NA#N	Furniture	Chairs & Chaimans	Fabric Task Chair	Jumbo Drum	8.0	60.98	26.2214	0.06	30.00	OK
8398	103989.0	2013-03-30	15496.0	Not Specified	Not Specified	Express Air	NA#N	Corporate	Tony Molinari	NA#N	Office Supplies	Storage & Organization	Series Personal Project File with Scope Front D...	Small Box	44.0	13.48	5.7964	0.10	4.51	OK

8399 rows x 20 columns

Exploratory Data Analysis - EDA

In [3]: df.size

Out[3]: 218374

In [4]: df.shape

Out[4]: (8399, 20)

In [5]: df.columns

Out[5]: Index(['RefID', 'OrderDate', 'OrderID', 'Priority', 'ShipDate', 'ShipMode', 'CustID', 'CustomerSegment', 'CustName', 'PostalCode', 'City', 'State', 'Country', 'Region', 'Market', 'ProdID', 'ProdCategory', 'ProdSubCat', 'ProName', 'ProContainer', 'Quantity', 'SalePrice', 'CostPrice', 'Discount', 'ShipCost', 'Status'], dtype='object')

In [6]: df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8399 entries, 0 to 8398
Data columns (total 20 columns):
 #   Column                                Non-Null Count  Dtype
---  --
 0   RefID                                8399 non-null    float64
 1   OrderDate                            8399 non-null    datetime64[ns]
 2   OrderID                              8399 non-null    float64
 3   Priority                              8399 non-null    object
 4   ShipDate                             8399 non-null    datetime64[ns]
 5   ShipMode                             8399 non-null    float64
 6   CustID                               8399 non-null    float64
 7   CustomerSegment                      8399 non-null    object
 8   CustName                             8399 non-null    object
 9   PostalCode                           8 non-null       float64
10   City                                 8 non-null       float64
11   State                                8399 non-null    object
12   Country                              8 non-null       float64
13   Region                              8399 non-null    float64
14   Market                               8 non-null       float64
15   ProdID                               8 non-null       float64
16   ProdCategory                         8399 non-null    object
17   ProdSubCat                           8399 non-null    object
18   ProName                              8399 non-null    object
19   ProContainer                         8399 non-null    object
20   Quantity                             8399 non-null    float64
21   SalePrice                            8399 non-null    float64
22   CostPrice                            8399 non-null    float64
23   Discount                             8399 non-null    float64
24   ShipCost                             8399 non-null    float64
25   Status                               8399 non-null    object
dtypes: datetime64[ns](2), float64(13), object(11)
memory usage: 1.7+ MB
```

In [7]: df.describe()

	RefID	OrderID	CustID	PostalCode	City	Country	Market	ProdID	Quantity	SalePrice	CostPrice	Discount	ShipCost
count	8399.000000	8399.000000	0.0	0.0	0.0	0.0	0.0	0.0	8399.000000	8399.000000	8399.000000	8399.000000	8399.000000
mean	10420.000000	1726.322778	NA#N	NA#N	NA#N	NA#N	NA#N	25.259999	89.346259	43.356775	0.049671	12.838657	32.092148
std	204.727698	1082.927977	NA#N	NA#N	NA#N	NA#N	NA#N	20.146174	290.354283	140.709955	0.032603	17.346652	106.000000
min	100001.000000	10001.000000	NA#N	NA#N	NA#N	NA#N	NA#N	1.000000	0.995000	0.182400	0.000000	0.000000	0.000000
25%	102101.500000	11359.500000	NA#N	NA#N	NA#N	NA#N	NA#N	16.000000	6.480000	3.168000	0.020000	0.340000	0.000000
50%	104200.000000	1726.000000	NA#N	NA#N	NA#N	NA#N	NA#N	32.000000	20.990000	10.146000	0.050000	6.070000	0.000000
75%	106299.500000	14105.000000	NA#N	NA#N	NA#N	NA#N	NA#N	51.000000	85.990000	40.639200	0.080000	13.990000	0.000000
max	108399.000000	15496.000000	NA#N	NA#N	NA#N	NA#N	NA#N	111.000000	6783.020000	4137.642200	0.250000	164.730000	0.000000

In [8]: df.isnull().sum()

	RefID	OrderDate	OrderID	Priority	ShipDate	ShipMode	CustID	CustomerSegment	CustName	PostalCode	City	State	Country	Region	Market	ProdID	ProdCategory	ProdSubCat	ProName	ProContainer	Quantity	SalePrice	CostPrice	Discount	ShipCost	Status
count	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
RefID	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
OrderDate	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
OrderID	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
Priority	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
ShipDate	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
ShipMode	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
CustID	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
CustomerSegment	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
CustName	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
PostalCode	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
City	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
State	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
Country	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
Region	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
Market	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
ProdID	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
ProdCategory	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
ProdSubCat	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
ProName	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
ProContainer	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
Quantity	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
SalePrice	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
CostPrice	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
Discount	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
ShipCost	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
Status	8399	8399	8399	8399	8399	8399	8	8399	8399	8	8	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399	8399
dtype: object	int64	int64	int64	int64	int64	int64	object	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64

8399 rows x 18 columns

In [9]: df1.columns

Out[9]: Index(['OrderDate', 'OrderID', 'Priority', 'ShipDate', 'ShipMode', 'CustomerSegment', 'State', 'Region', 'ProdCategory', 'ProdSubCat', 'ProName', 'ProContainer', 'Quantity', 'SalePrice', 'CostPrice', 'Discount', 'ShipCost', 'Status', 'TotalSales', 'TotalCost', 'discount_perc', 'Discount_amount', 'GrossSales', 'GrossCostPrice'], dtype='object')

In [11]: #Create a new column total cost
df1['TotalSales'] = df1['SalePrice'] * df1['Quantity']
df1.head(3)

Out[11]:

OrderDate	OrderID	Priority	ShipDate	ShipMode	CustomerSegment	State	Region	ProdCategory	ProdSubCat	ProName	ProContainer	Quantity	SalePrice	CostPrice	Discount	ShipCost	Status	TotalSales
2009-04-01	10001.0	Specified	Not Specified	Express Air	Home Office	Quebec	Quebec	Office Supplies	Storage & Organization	Safco Industrial Wire Shelving	Large Box	10.0	95.99	57.5940	0.08	35.00	OK	959.90
2009-04-01	10002.0	High	2009-04-03	Regular Air	Small Business	Ontario	Ontario	Office Supplies	Storage & Organization	Penna STOR-ALL™ Hanging File Box, 13 1/8"W x 1...	Small Box	36.0	5.98	1.9136	0.10	4.69	OK	215.28

In [12]: #Create a new column total cost
df1['TotalCost'] = df1['CostPrice'] * df1['Quantity']
df1.head(3)

Out[12]:

OrderDate	OrderID	Priority	ShipDate	ShipMode	CustomerSegment	State	Region	ProdCategory	ProdSubCat	ProName	ProContainer	Quantity	SalePrice	CostPrice	Discount	ShipCost	Status	TotalSales	TotalCost
2009-04-01	10001.0	Specified	Not Specified	Express Air	Home Office	Quebec	Quebec	Office Supplies	Storage & Organization	Safco Industrial Wire Shelving	Large Box	10.0	95.99	57.5940	0.08	35.00	OK	959.90	575.940
2009-04-01	10002.0	High	2009-04-03	Regular Air	Small Business	Ontario	Ontario	Office Supplies	Storage & Organization	Penna STOR-ALL™ Hanging File Box, 13 1/8"W x 1...	Small Box	36.0	5.98	1.9136	0.10	4.69	OK	215.28	68.886

In [13]: df1['Discount'].unique()

Out[13]: array([0.08, 0.1, 0.06, 0., ..., 0.07, 0.05, 0.09, 0.03, 0.04, 0.01, 0.02, 0.1, 0.17, 0.11, 0.16, 0.25])

In [14]: #Calculating discount amount
df1['Discount_perc'] = df1['Discount'] * 100
df1.head(3)

Out[14]:

OrderDate	OrderID	Priority	ShipDate	ShipMode	CustomerSegment	State	Region	ProdCategory	ProdSubCat	ProName	ProContainer	Quantity	SalePrice	CostPrice	Discount	ShipCost	Status	TotalSales	TotalCost	discount_perc
2009-04-01	10001.0	Specified	Not Specified	Express Air	Home Office	Quebec	Quebec	Office Supplies	Storage & Organization	Safco Industrial Wire Shelving	Large Box	10.0	95.99	57.5940	0.08	35.00	OK	959.90	575.940	8.0
2009-04-01	10002.0	High	2009-04-03	Regular Air	Small Business	Ontario	Ontario	Office Supplies	Storage & Organization	Penna STOR-ALL™ Hanging File Box, 13 1/8"W x 1...	Small Box	36.0	5.98	1.9136	0.10	4.69	OK	215.28	68.886	10.0

2 rows x 22 columns

In [15]: #Sales after Discount
df1['GrossSales'] = df1['TotalSales'] - df1['Discount_amount']
df1.head(3)

Out[15]: