

project3

April 19, 2024

1 Import Lib

```
[1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import statsmodels.api as sm
from sklearn.impute import KNNImputer
from sklearn.feature_selection import SelectKBest, f_regression
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error, r2_score, mean_absolute_error
from sklearn.ensemble import RandomForestRegressor
from sklearn.neighbors import KNeighborsRegressor
from sklearn.neural_network import MLPRegressor
```

2 Reading Dataset

```
[2]: df=pd.read_excel(r"C:\Users\ganesh\Desktop\Sample - Superstore (1).xlsx") #_
      ↪reading excel file
pd.set_option("display.max_column",22)
df.head()      #to display top 5 rows
```

```
[2]:
```

	Row ID	Order ID	Order Date	Ship Date	Ship Mode	Customer ID	\
0	1.0	CA-2016-152156	2016-11-08	2016-11-11	Second Class	CG-12520	
1	2.0	CA-2016-152156	2016-11-08	2016-11-11	Second Class	CG-12520	
2	3.0	CA-2016-138688	2016-06-12	2016-06-16	Second Class	DV-13045	
3	4.0	US-2015-108966	2015-10-11	2015-10-18	Standard Class	SO-20335	
4	5.0	US-2015-108966	2015-10-11	2015-10-18	Standard Class	SO-20335	

	Customer Name	Segment	Country	City	State	\
0	Claire Gute	Consumer	United States	Henderson	Kentucky	
1	Claire Gute	Consumer	United States	Henderson	Kentucky	
2	Darrin Van Huff	Corporate	United States	Los Angeles	California	
3	Sean O'Donnell	Consumer	United States	Fort Lauderdale	Florida	
4	Sean O'Donnell	Consumer	United States	Fort Lauderdale	Florida	

	Postal Code	Region	Product ID	Category	Sub-Category	\
0	42420.0	South	FUR-BO-10001798	Furniture	Bookcases	
1	42420.0	South	FUR-CH-10000454	Furniture	Chairs	
2	90036.0	West	OFF-LA-10000240	Office Supplies	Labels	
3	33311.0	South	FUR-TA-10000577	Furniture	Tables	
4	33311.0	South	OFF-ST-10000760	Office Supplies	Storage	

	Product Name	Sales	Quantity	\
0	Bush Somerset Collection Bookcase	261.96	2.0	
1	Hon Deluxe Fabric Upholstered Stacking Chairs,...	731.94	3.0	
2	Self-Adhesive Address Labels for Typewriters b...	14.62	2.0	
3	Bretford CR4500 Series Slim Rectangular Table	957.5775	5.0	
4	Eldon Fold 'N Roll Cart System	22.368	2.0	

	Discount	Profit
0	0.00	41.9136
1	0.00	219.5820
2	0.00	6.8714
3	0.45	-383.0310
4	0.20	2.5164

```
[3]: df.shape      #shape of dataframe
```

```
[3]: (9994, 21)
```

```
[4]: df.info()     #information of dataframe
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9994 entries, 0 to 9993
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Row ID                9985 non-null   float64
1   Order ID              9981 non-null   object
2   Order Date            9981 non-null   datetime64[ns]
3   Ship Date             9979 non-null   datetime64[ns]
4   Ship Mode             9975 non-null   object
5   Customer ID           9968 non-null   object
6   Customer Name         9937 non-null   object
7   Segment               9942 non-null   object
8   Country               9930 non-null   object
9   City                  9949 non-null   object
10  State                 9937 non-null   object
11  Postal Code           9957 non-null   float64
12  Region                9954 non-null   object
13  Product ID            9956 non-null   object
14  Category              9965 non-null   object
```

```

15 Sub-Category    9952 non-null    object
16 Product Name    9936 non-null    object
17 Sales           9932 non-null    object
18 Quantity        9948 non-null    float64
19 Discount        9957 non-null    float64
20 Profit          9944 non-null    float64
dtypes: datetime64[ns](2), float64(5), object(14)
memory usage: 1.6+ MB

```

```
[5]: df.isnull().sum() #finding number of null values in each column
```

```

[5]: Row ID          9
     Order ID       13
     Order Date     13
     Ship Date      15
     Ship Mode      19
     Customer ID    26
     Customer Name  57
     Segment       52
     Country       64
     City          45
     State         57
     Postal Code   37
     Region        40
     Product ID    38
     Category      29
     Sub-Category  42
     Product Name  58
     Sales         62
     Quantity      46
     Discount      37
     Profit        50
     dtype: int64

```

```
[6]: df.isnull().sum()/df.shape[0]*100 #percentage of missing values in each
      ↪ columns
```

```

[6]: Row ID          0.090054
     Order ID       0.130078
     Order Date     0.130078
     Ship Date      0.150090
     Ship Mode      0.190114
     Customer ID    0.260156
     Customer Name  0.570342
     Segment       0.520312
     Country       0.640384
     City          0.450270

```

```

State          0.570342
Postal Code    0.370222
Region         0.400240
Product ID     0.380228
Category       0.290174
Sub-Category   0.420252
Product Name   0.580348
Sales          0.620372
Quantity       0.460276
Discount       0.370222
Profit         0.500300
dtype: float64

```

```
[7]: df.duplicated().sum()  #duplicate rows
```

```
[7]: 10
```

```
[8]: for i in df.select_dtypes(include="object").columns:  #counts number of unique
      ↪ values
      print(df[i].value_counts())
      print("***"*10)
```

```

CA-2017-100111    14
CA-2017-157987    12
US-2016-108504    11
CA-2016-165330    11
US-2015-126977    10
..
US-2017-107636    1
US-2014-165862    1
CA-2016-101448    1
US-2017-117331    1
CA-2017-119914    1
Name: Order ID, Length: 5000, dtype: int64
*****
Standard Class    5958
Second Class      1940
First Class       1534
Same Day          543
Name: Ship Mode, dtype: int64
*****
WB-21850          37
AP-10915          35
PP-18955          34
JL-15835          34
MA-17560          34
..

```

LD-16855	1
AO-10810	1
CJ-11875	1
RE-19405	1
JR-15700	1

Name: Customer ID, Length: 793, dtype: int64

William Brown	37
Arthur Prichep	35
John Lee	34
Paul Prost	34
Matt Abelman	34
	..
Lela Donovan	1
Jocasta Rupert	1
Carl Jackson	1
Anthony O'Donnell	1
Ricardo Emerson	1

Name: Customer Name, Length: 792, dtype: int64

Consumer	5171
Corporate	2997
Home Office	1774

Name: Segment, dtype: int64

United States	9930
---------------	------

Name: Country, dtype: int64

New York City	913
Los Angeles	744
Philadelphia	529
San Francisco	507
Seattle	428
	...
Cedar Rapids	1
Palatine	1
Jefferson City	1
Waterloo	1
Iowa City	1

Name: City, Length: 530, dtype: int64

California	1991
New York	1122
Texas	975
Pennsylvania	578
Washington	506
	...
31.744	1

```

18.180000000000003      1
471.92                  1
89.584                  1
12.624                  1
Name: State, Length: 68, dtype: int64
*****
West      3197
East      2828
Central   2305
South     1605
0          11
0.2        6
0.7         1
0.1         1
Name: Region, dtype: int64
*****
OFF-PA-10001970    19
TEC-AC-10003832    18
FUR-FU-10004270    16
FUR-CH-10001146    15
TEC-AC-10002049    15
..
TEC-MA-10002937     1
TEC-PH-10003535     1
OFF-AP-10002734     1
TEC-MA-10003353     1
OFF-ST-10001627     1
Name: Product ID, Length: 1878, dtype: int64
*****
Office Supplies    6003
Furniture          2116
Technology         1846
Name: Category, dtype: int64
*****
Binders      1512
Paper        1367
Furnishings   956
Phones       890
Storage       843
Art           787
Accessories   774
Chairs        612
Appliances    464
Labels        363
Tables        317
Envelopes     252
Bookcases     228
Fasteners     217

```

```

Supplies          190
Machines          114
Copiers           66
Name: Sub-Category, dtype: int64
*****
Staple envelope          48
Easy-staple paper       46
Staples                 46
Avery Non-Stick Binders 20
Staples in misc. colors 18
..
Eldon File Chest Portable File 1
Hewlett-Packard Deskjet D4360 Printer 1
Jiffy Padded Mailers with Self-Seal Closure 1
Hunt BOSTON Model 1606 High-Volume Electric Pencil Sharpener, Beige 1
Eldon Jumbo ProFile Portable File Boxes Graphite/Black 1
Name: Product Name, Length: 1857, dtype: int64
*****
12.960    54
15.552    39
19.440    39
10.368    35
32.400    28
..
487.920    1
25.920     1
95.736     1
3.392      1
275.880     1
Name: Sales, Length: 6128, dtype: int64
*****

```

#Converting DataType of sales from object to float (datatype was object because there were some string values in column,this process convert string to nan)

```

[9]: df["Sales"]=pd.to_numeric(df["Sales"],errors='coerce')
     pro=df[df["Sales"].isna()]
     print(pro)

```

	Row ID	Order ID	Order Date	Ship Date	Ship Mode \
193	194.0	CA-2015-102281	2015-10-12	2015-10-14	First Class
194	195.0	CA-2015-131457	2015-10-31	2015-11-06	Standard Class
195	196.0	CA-2014-140004	2014-03-21	2014-03-25	Standard Class
196	197.0	CA-2014-140004	2014-03-21	2014-03-25	Standard Class
197	198.0	CA-2017-107720	2017-11-06	2017-11-13	Standard Class
...
7177	7187.0	CA-2017-133102	2017-08-17	2017-08-24	Standard Class
7178	7188.0	CA-2017-133102	2017-08-17	2017-08-24	Standard Class

7179	7189.0	CA-2017-133102	2017-08-17	2017-08-24	Standard Class
7180	7190.0	CA-2016-164399	2016-11-12	2016-11-15	First Class
7181	7191.0	CA-2016-164399	2016-11-12	2016-11-15	First Class

	Customer ID	Customer Name	Segment	Country	\
193	MP-17470	Mark Packer	Home Office	United States	
194	MZ-17515	Mary Zewe	Corporate	United States	
195	CB-12025	Cassandra Brandow	Consumer	United States	
196	CB-12025	Cassandra Brandow	Consumer	United States	
197	VM-21685	Valerie Mitchum	Home Office	United States	
...	
7177	ED-13885	Emily Ducich	Home Office	United States	
7178	ED-13885	Emily Ducich	Home Office	United States	
7179	ED-13885	Emily Ducich	Home Office	United States	
7180	DW-13480	Dianna Wilson	Home Office	United States	
7181	DW-13480	Dianna Wilson	Home Office	United States	

	City	State	Postal Code	Region	Product ID	\
193	New York City	New York	10035.0	East	NaN	
194	Redlands	California	92374.0	West	NaN	
195	Hamilton	Ohio	45011.0	East	NaN	
196	Hamilton	Ohio	45011.0	East	NaN	
197	Westfield	New Jersey	7090.0	East	NaN	
...	
7177	Houston	Texas	77095.0	Central	FUR-CH-10002017	
7178	Houston	Texas	77095.0	Central	FUR-FU-10003247	
7179	Houston	Texas	77095.0	Central	OFF-AP-10001563	
7180	San Diego	California	92024.0	West	TEC-PH-10004908	
7181	San Diego	California	92024.0	West	FUR-TA-10003392	

	Category	Sub-Category	Product Name	Sales	Quantity	Discount	\
193	NaN	NaN	NaN	NaN	NaN	NaN	
194	NaN	NaN	NaN	NaN	NaN	NaN	
195	NaN	NaN	NaN	NaN	NaN	NaN	
196	NaN	NaN	NaN	NaN	NaN	NaN	
197	NaN	NaN	NaN	NaN	NaN	NaN	
...	
7177	Furniture	Chairs	NaN	NaN	NaN	NaN	
7178	Furniture	Furnishings	NaN	NaN	NaN	NaN	
7179	Office Supplies	Appliances	NaN	NaN	NaN	NaN	
7180	Technology	Phones	NaN	NaN	NaN	NaN	
7181	Furniture	Tables	NaN	NaN	NaN	NaN	

	Profit
193	NaN
194	NaN
195	NaN
196	NaN


```

197      NaN
...
7177      NaN
7178      NaN
7179      NaN
7180      NaN
7181      NaN

```

[83 rows x 21 columns]

```
[10]: df.isnull().sum() #checking number of null values for eah column
```

```

[10]: Row ID          9
      Order ID       13
      Order Date     13
      Ship Date      15
      Ship Mode      19
      Customer ID    26
      Customer Name   57
      Segment        52
      Country        64
      City           45
      State          57
      Postal Code     37
      Region         40
      Product ID     38
      Category       29
      Sub-Category   42
      Product Name   58
      Sales          83
      Quantity       46
      Discount       37
      Profit         50
      dtype: int64

```

```
[11]: df[df["Row ID"]==6858] #checking row that has changed from string to nan
```

```

[11]:      Row ID      Order ID Order Date  Ship Date      Ship Mode \
6857  6858.0  CA-2017-128965 2017-04-17 2017-04-22  Standard Class

      Customer ID Customer Name  Segment      Country      City \
6857    PS-18760  Pamela Stobb  Consumer  United States  Los Angeles

      State  Postal Code Region      Product ID      Category \
6857  California      90008.0    West  OFF-PA-10004911  Office Supplies

      Sub-Category Product Name  Sales  Quantity  Discount  Profit

```

6857 Paper Paper NaN 28.14 3.0 0.0

```
[12]: df.dropna(subset=["Sales"], inplace=True) #dropping null values of sales(we cant  
      ↳impute data in sales(dependent variable) because it creates a arbitrary  
      ↳values whih affects model prediction)
```

```
[13]: df.shape #shape of dataframe
```

```
[13]: (9911, 21)
```

```
[14]: df["Sales"].dtypes #cheking dataframe of sales column
```

```
[14]: dtype('float64')
```

```
[15]: df.isnull().sum() #checking number of null values
```

```
[15]: Row ID          9  
      Order ID       13  
      Order Date     13  
      Ship Date      15  
      Ship Mode       19  
      Customer ID    26  
      Customer Name  46  
      Segment        41  
      Country        53  
      City           45  
      State          51  
      Postal Code    31  
      Region         31  
      Product ID     11  
      Category        2  
      Sub-Category   4  
      Product Name    4  
      Sales           0  
      Quantity        0  
      Discount        0  
      Profit         13  
      dtype: int64
```

```
[16]: df["City"].isnull().sum() #cheking number of null values in city
```

```
[16]: 45
```

```
[17]: df["City"].unique() #hecking unique values in city column
```

```
[17]: array(['Henderson', 'Los Angeles', 'Fort Lauderdale', 'Concord',  
          'Seattle', 'Fort Worth', 'Madison', 'West Jordan', 'San Francisco',
```

'Fremont', 'Philadelphia', 'Orem', 'Houston', 'Richardson',
 'Naperville', 'Melbourne', 'Eagan', 'Westland', 'Dover',
 'New Albany', 'New York City', 'Troy', 'Chicago', 'Gilbert',
 'Springfield', 'Jackson', 'Memphis', 'Decatur', 'Durham',
 'Columbia', 'Rochester', nan, 'Aurora', 'Charlotte', 'Orland Park',
 'Urbandale', 'Columbus', 'Bristol', 'Wilmington', 'Bloomington',
 'Phoenix', 'Roseville', 'Independence', 'Pasadena', 'Newark',
 'Franklin', 'Scottsdale', 'San Jose', 'Edmond', 'Carlsbad',
 'San Antonio', 'Monroe', 'Fairfield', 'Grand Prairie', 'Denver',
 'Dallas', 'Whittier', 'Saginaw', 'Medina', 'Detroit', 'Tampa',
 'Santa Clara', 'Lakeville', 'San Diego', 'Brentwood',
 'Chapel Hill', 'Morristown', 'Cincinnati', 'Inglewood', 'Portland',
 'Tamarac', 'Colorado Springs', 'Belleville', 'Taylor', 'Lakewood',
 'Arlington', 'Arvada', 'Hackensack', 'Saint Petersburg',
 'Long Beach', 'Hesperia', 'Murfreesboro', 'Austin', 'Lowell',
 'Manchester', 'Harlingen', 'Tucson', 'Quincy', 'Pembroke Pines',
 'Des Moines', 'Peoria', 'Las Vegas', 'Warwick', 'Miami',
 'Huntington Beach', 'Richmond', 'Louisville', 'Lawrence', 'Canton',
 'New Rochelle', 'Gastonia', 'Jacksonville', 'Auburn', 'Akron',
 'Norman', 'Park Ridge', 'Amarillo', 'Lindenhurst', 'Huntsville',
 'Fayetteville', 'Costa Mesa', 'Parker', 'Atlanta', 'Gladstone',
 'Great Falls', 'Montgomery', 'Mesa', 'Green Bay', 'Anaheim',
 'Marysville', 'Salem', 'Laredo', 'Grove City', 'Dearborn',
 'Warner Robins', 'Vallejo', 'Minneapolis', 'Mission Viejo',
 'Rochester Hills', 'Plainfield', 'Sierra Vista', 'Vancouver',
 'Cleveland', 'Tyler', 'Burlington', 'Waynesboro', 'Chester',
 'Cary', 'Palm Coast', 'Mount Vernon', 'Hialeah', 'Oceanside',
 'Evanston', 'Trenton', 'Cottage Grove', 'Bossier City',
 'Lancaster', 'Asheville', 'Lake Elsinore', 'Omaha', 'Edmonds',
 'Santa Ana', 'Milwaukee', 'Florence', 'Lorain', 'Linden',
 'Salinas', 'New Brunswick', 'Garland', 'Norwich', 'Alexandria',
 'Toledo', 'Farmington', 'Riverside', 'Torrance', 'Round Rock',
 'Boca Raton', 'Virginia Beach', 'Murrieta', 'Olympia',
 'Washington', 'Jefferson City', 'Saint Peters', 'Rockford',
 'Brownsville', 'Yonkers', 'Oakland', 'Clinton', 'Encinitas',
 'Roswell', 'Jonesboro', 'Antioch', 'Homestead', 'La Porte',
 'Lansing', 'Cuyahoga Falls', 'Reno', 'Harrisonburg', 'Escondido',
 'Royal Oak', 'Rockville', 'Coral Springs', 'Buffalo',
 'Boynton Beach', 'Gulfport', 'Fresno', 'Greenville', 'Macon',
 'Cedar Rapids', 'Providence', 'Pueblo', 'Saint Paul', 'Deltona',
 'Murray', 'Middletown', 'Freeport', 'Pico Rivera', 'Provo',
 'Pleasant Grove', 'Smyrna', 'Parma', 'Mobile', 'New Bedford',
 'Irving', 'Vineland', 'Glendale', 'Niagara Falls', 'Thomasville',
 'Westminster', 'Coppell', 'Pomona', 'North Las Vegas', 'Allentown',
 'Tempe', 'Laguna Niguel', 'Bridgeton', 'Everett', 'Watertown',
 'Appleton', 'Bellevue', 'Allen', 'El Paso', 'Grapevine',
 'Carrollton', 'Kent', 'Lafayette', 'Tigard', 'Skokie', 'Plano',

'Suffolk', 'Indianapolis', 'Bayonne', 'Dublin', 'Greensboro',
'Baltimore', 'Kenosha', 'Olathe', 'Tulsa', 'Redmond', 'Raleigh',
'Muskogee', 'Meriden', 'Bowling Green', 'South Bend', 'Spokane',
'Keller', 'Port Orange', 'Medford', 'Charlottesville', 'Missoula',
'Apopka', 'Reading', 'Broomfield', 'Paterson', 'Oklahoma City',
'Chesapeake', 'Lubbock', 'Johnson City', 'San Bernardino',
'Leominster', 'Bozeman', 'Perth Amboy', 'Ontario',
'Rancho Cucamonga', 'Moorhead', 'Mesquite', 'Redlands', 'Stockton',
'Ormond Beach', 'Sunnyvale', 'York', 'College Station',
'Saint Louis', 'Manteca', 'San Angelo', 'Salt Lake City',
'Knoxville', 'Little Rock', 'Lincoln Park', 'Marion', 'Littleton',
'Bangor', 'Southaven', 'New Castle', 'Midland', 'Sioux Falls',
'Fort Collins', 'Clarksville', 'Sacramento', 'Thousand Oaks',
'Malden', 'Holyoke', 'Albuquerque', 'Sparks', 'Coachella',
'Elmhurst', 'Passaic', 'North Charleston', 'Newport News',
'Jamestown', 'Mishawaka', 'Westfield', 'La Quinta', 'Tallahassee',
'Nashville', 'Bellingham', 'Woodstock', 'Haltom City', 'Wheeling',
'Summerville', 'Hot Springs', 'Englewood', 'Las Cruces', 'Hoover',
'Frisco', 'Vacaville', 'Waukesha', 'Bakersfield', 'Pompano Beach',
'Corpus Christi', 'Redondo Beach', 'Orlando', 'Orange',
'Lake Charles', 'Highland Park', 'Hempstead', 'Noblesville',
'Apple Valley', 'Mount Pleasant', 'Sterling Heights', 'Eau Claire',
'Pharr', 'Billings', 'Gresham', 'Chattanooga', 'Meridian',
'Bolingbrook', 'Lakeland', 'Maple Grove', 'Woodland',
'Missouri City', 'Pearland', 'San Mateo', 'Grand Rapids',
'Visalia', 'Overland Park', 'Temecula', 'Yucaipa', 'Revere',
'Conroe', 'Tinley Park', 'Dubuque', 'Dearborn Heights', 'Santa Fe',
'Hickory', 'Carol Stream', 'Saint Cloud', 'North Miami',
'Plantation', 'Port Saint Lucie', 'Rock Hill', 'Odessa',
'West Allis', 'Chula Vista', 'Manhattan', 'Altoona', 'Thornton',
'Champaign', 'Texarkana', 'Edinburg', 'Baytown', 'Greenwood',
'Woonsocket', 'Superior', 'Bedford', 'Covington', 'Broken Arrow',
'Miramar', 'Hollywood', 'Deer Park', 'Wichita', 'McAllen',
'Iowa City', 'Boise', 'Cranston', 'Port Arthur', 'Citrus Heights',
'The Colony', 'Daytona Beach', 'Bullhead City', 'Portage', 'Fargo',
'Elkhart', 'San Gabriel', 'Hamilton', 'Margate', 'Sandy Springs',
'Mentor', 'Lawton', 'Hampton', 'Rome', 'La Crosse', 'Lewiston',
'Hattiesburg', 'Danville', 'Logan', 'Waterbury', 'Athens',
'Avondale', 'Marietta', 'Yuma', 'Wausau', 'Pasco', 'Oak Park',
'Pensacola', 'League City', 'Gaithersburg', 'Lehi', 'Tuscaloosa',
'Moreno Valley', 'Georgetown', 'Loveland', 'Chandler', 'Helena',
'Kirkwood', 'Waco', 'Frankfort', 'Bethlehem', 'Grand Island',
'Woodbury', 'Rogers', 'Clovis', 'Jupiter', 'Santa Barbara',
'Cedar Hill', 'Norfolk', 'Draper', 'Ann Arbor', 'La Mesa',
'Pocatello', 'Holland', 'Milford', 'Buffalo Grove', 'Lake Forest',
'Redding', 'Chico', 'Utica', 'Conway', 'Cheyenne', 'Owensboro',
'Caldwell', 'Kenner', 'Nashua', 'Bartlett', 'Redwood City',

```
'Lebanon', 'Santa Maria', 'Des Plaines', 'Longview',
'Hendersonville', 'Waterloo', 'Cambridge', 'Palatine', 'Beverly',
'Eugene', 'Oxnard', 'Renton', 'Glenview', 'Delray Beach',
'Commerce City', 'Texas City', 'Wilson', 'Rio Rancho', 'Goldsboro',
'Montebello', 'El Cajon', 'West Palm Beach', 'Abilene', 'Normal',
'Saint Charles', 'Camarillo', 'Hillsboro', 'Burbank', 'Modesto',
'Garden City', 'Atlantic City', 'Longmont', 'Davis', 'Morgan Hill',
'Clifton', 'Sheboygan', 'East Point', 'Rapid City', 'Andover',
'Kissimmee', 'Shelton', 'Danbury', 'Sanford', 'San Marcos',
'Greeley', 'Mansfield', 'Elyria', 'Twin Falls', 'Coral Gables',
'Romeoville', 'Marlborough', 'Laurel', 'Bryan', 'Pine Bluff',
'Aberdeen', 'Hagerstown', 'East Orange', 'Arlington Heights',
'Oswego', 'Beaumont', 'Coon Rapids', 'San Clemente',
'San Luis Obispo', 'Springdale', 'Lodi', 'Mason'], dtype=object)
```

```
[18]: df.dropna(subset=["City"],inplace=True) #dropped nan values in city
```

```
[19]: df.isnull().sum() #checking null values
```

```
[19]: Row ID          5
      Order ID      9
      Order Date    9
      Ship Date     11
      Ship Mode      8
      Customer ID    5
      Customer Name  23
      Segment       18
      Country       21
      City          0
      State         6
      Postal Code    0
      Region        0
      Product ID     2
      Category       2
      Sub-Category   4
      Product Name    4
      Sales          0
      Quantity       0
      Discount       0
      Profit        13
      dtype: int64
```

```
[20]: df.shape #shape of dataframe
```

```
[20]: (9866, 21)
```

```
[21]: df["State"].isnull().sum() #checking null values in state
```

[21]: 6

```
[22]: df["State"].unique() #finding unique values in state column
```

```
[22]: array(['Kentucky', 'California', 'Florida', 'North Carolina',  
        'Washington', 'Texas', 'Wisconsin', 'Utah', 'Nebraska',  
        'Pennsylvania', 'Illinois', 'Minnesota', 'Michigan', 'Delaware',  
        'Indiana', 'New York', 'Arizona', 'Virginia', 'Tennessee',  
        'Alabama', 'South Carolina', 'Colorado', 'Iowa', 'Ohio',  
        'Missouri', 'Oklahoma', 'New Mexico', 'Louisiana', 'Connecticut',  
        'New Jersey', 'Oregon', 'Massachusetts', 'Georgia', 'Nevada',  
        'Rhode Island', nan, 'Mississippi', 'Arkansas', 'Montana',  
        'New Hampshire', 'Maryland', 'District of Columbia', 'Kansas',  
        'Vermont', 'Maine', 'South Dakota', 'Idaho', 'North Dakota',  
        'Wyoming', 24.849999999999998, 12.624, 89.584, 471.92,  
        18.180000000000003, 31.744, 5.904, 621.7600000000001, 59.98, 48.87,  
        154.9, 5.92, 30.18, 24.1, 8.78, 376.74, 29.52, 11.96,  
        26.400000000000002, 'West Virginia'], dtype=object)
```

```
[23]: state_city=df[["State","City"]] #created dataframe for state, city  
unique_state_city=state_city.drop_duplicates() #drop dupliates  
unique_state_city.head(10) #hecking respective city name for state to change  
↪errors in state column
```

```
[23]:
```

	State	City
0	Kentucky	Henderson
2	California	Los Angeles
3	Florida	Fort Lauderdale
12	North Carolina	Concord
13	Washington	Seattle
14	Texas	Fort Worth
16	Wisconsin	Madison
17	Utah	West Jordan
18	California	San Francisco
21	Nebraska	Fremont

```
[24]: a=[24.849999999999998, 12.624, 89.584, 471.92,  
        18.180000000000003, 31.744, 5.904, 621.7600000000001, 59.98, 48.87,  
        154.9, 5.92, 30.18, 24.1, 8.78, 376.74, 29.52, 11.96,  
        26.400000000000002] #values are errors in state column  
for i in a:  
    print(df[df["State"]==i])
```

Row ID	Order ID	Order Date	Ship Date	Ship Mode	\
6972	6973.0	CA-2017-153822	2017-09-19	2017-09-25	Standard Class

Customer ID	Customer Name	Segment	Country	City	State	\
-------------	---------------	---------	---------	------	-------	---

6972	AB-10105	Adrian Barton	Consumer	United States	Phoenix	24.85
------	----------	---------------	----------	---------------	---------	-------

	Postal Code	Region	Product ID	Category	Sub-Category	\
6972	5.0	0	7.7035	Technology	Phones	

	Product Name	Sales	Quantity	Discount	\
6972	Polycom VoiceStation 500 Conference phone	471.92	2.0	0.2	

	Profit
6972	29.495

	Row ID	Order ID	Order Date	Ship Date	Ship Mode	\
6973	6974.0	CA-2017-153822	2017-09-19	2017-09-25	Standard Class	

	Customer ID	Customer Name	Segment	Country	City	State	\
6973	AB-10105	Adrian Barton	Consumer	United States	Phoenix	12.624	

	Postal Code	Region	Product ID	Category	Sub-Category	\
6973	2.0	0.2	-2.5248	Office Supplies	Binders	

	Product Name	Sales	Quantity	Discount	Profit
6973	Plastic Binding Combs	18.18	4.0	0.7	-13.938

	Row ID	Order ID	Order Date	Ship Date	Ship Mode	\
6974	6975.0	CA-2017-146185	2017-09-15	2017-09-19	Standard Class	

	Customer ID	Customer Name	Segment	Country	City	State	\
6974	CC-12145	Charles Crestani	Consumer	United States	Houston	89.584	

	Postal Code	Region	Product ID	Category	Sub-Category	\
6974	2.0	0.2	4.4792	Office Supplies	Art	

	Product Name	Sales	Quantity	Discount	Profit
6974	Prismacolor Color Pencil Set	31.744	2.0	0.2	8.3328

	Row ID	Order ID	Order Date	Ship Date	Ship Mode	\
6975	6976.0	CA-2015-112144	2015-06-28	2015-07-02	Standard Class	

	Customer ID	Customer Name	Segment	Country	City	State	\
6975	CY-12745	Craig Yedwab	Corporate	United States	Gilbert	471.92	

	Postal Code	Region	Product ID	Category	Sub-Category	\
6975	2.0	0.2	29.495	Office Supplies	Labels	

	Product Name	Sales	Quantity	Discount	Profit
6975	Avery 501	5.904	2.0	0.2	1.9926

	Row ID	Order ID	Order Date	Ship Date	Ship Mode	\
6976	6977.0	CA-2015-112144	2015-06-28	2015-07-02	Standard Class	

	Customer ID	Customer Name	Segment	Country	City	State	\
6976	CY-12745	Craig Yedwab	Corporate	United States	Gilbert	18.18	

Postal Code	Region	Product ID	Category	Sub-Category	\
6976	4.0	0.7	-13.938	Furniture	Furnishings

Product Name	Sales	Quantity	Discount	Profit
6976 Electrix Halogen Magnifier Lamp	621.76	4.0	0.2	46.632

Row ID	Order ID	Order Date	Ship Date	Ship Mode	Customer ID	\
6977	6978.0	US-2016-119298	2016-11-25	2016-11-28	First Class	EP-13915

Customer Name	Segment	Country	City	State	Postal Code	\
6977 Emily Phan	Consumer	United States	Jonesboro	31.744	2.0	

Region	Product ID	Category	Sub-Category	\
6977	0.2	8.3328	Technology	Phones

Product Name	Sales	Quantity	\
6977 OtterBox Defender Series Case - Samsung Galaxy S4	59.98	2.0	

Discount	Profit
6977	0.0 17.994

Row ID	Order ID	Order Date	Ship Date	Ship Mode	Customer ID	\
6978	6979.0	CA-2017-155159	2017-11-25	2017-11-29	Second Class	DL-13315

Customer Name	Segment	Country	City	State	Postal Code	\
6978 Delfina Latchford	Consumer	United States	Atlanta	5.904	2.0	

Region	Product ID	Category	Sub-Category	\
6978	0.2	1.9926	Office Supplies	Paper

Product Name	Sales	Quantity	Discount	Profit
6978 Wirebound Message Book, 4 per Page	48.87	9.0	0.0	23.9463

Row ID	Order ID	Order Date	Ship Date	Ship Mode	\
6979	6980.0	CA-2017-149076	2017-01-14	2017-01-19	Standard Class

Customer ID	Customer Name	Segment	Country	City	\
6979	SO-20335	Sean O'Donnell	Consumer	United States	Los Angeles

State	Postal Code	Region	Product ID	Category	Sub-Category	\
6979	621.76	4.0	0.2	46.632	Office Supplies	Paper

Product Name	Sales	Quantity	Discount	Profit
6979 Xerox 19	154.9	5.0	0.0	69.705

Row ID	Order ID	Order Date	Ship Date	Ship Mode	Customer ID	\
6980	6981.0	CA-2014-146990	2014-11-07	2014-11-08	First Class	BP-11095

Customer Name	Segment	Country	City	State	\
6980 Bart Pistole	Corporate	United States	New York City	59.98	

	Postal Code	Region	Product ID	Category	Sub-Category	\
6980	2.0	0	17.994	Office Supplies	Fasteners	
	Product Name	Sales	Quantity	Discount	Profit	
6980	Binder Clips by OIC	5.92	4.0	0.0	2.8416	
	Row ID	Order ID	Order Date	Ship Date	Ship Mode	Customer ID \
6981	6982.0	CA-2014-146990	2014-11-07	2014-11-08	First Class	BP-11095
	Customer Name	Segment	Country	City	State	\
6981	Bart Pistole	Corporate	United States	New York City	48.87	
	Postal Code	Region	Product ID	Category	Sub-Category	\
6981	9.0	0	23.9463	Office Supplies	Paper	
	Product Name	Sales	Quantity	Discount	\	
6981	Riverleaf Stik-Withit Designer Note Cubes	30.18	3.0	0.0		
	Profit					
6981	13.8828					
	Row ID	Order ID	Order Date	Ship Date	Ship Mode	\
6982	6983.0	CA-2016-116526	2016-09-01	2016-09-05	Standard Class	
	Customer ID	Customer Name	Segment	Country	City	State \
6982	JA-15970	Joseph Airdo	Consumer	United States	Detroit	154.9
	Postal Code	Region	Product ID	Category	Sub-Category	\
6982	5.0	0	69.705	Office Supplies	Binders	
	Product Name	Sales	Quantity	\		
6982	Wilson Jones Turn Tabs Binder Tool for Ring Bi...	24.1	5.0			
	Discount	Profit				
6982	0.0	11.086				
	Row ID	Order ID	Order Date	Ship Date	Ship Mode	\
6983	6984.0	CA-2016-116526	2016-09-01	2016-09-05	Standard Class	
	Customer ID	Customer Name	Segment	Country	City	State \
6983	JA-15970	Joseph Airdo	Consumer	United States	Detroit	5.92
	Postal Code	Region	Product ID	Category	Sub-Category	\
6983	4.0	0	2.8416	Technology	Phones	
	Product Name	Sales	Quantity	\		
6983	Belkin Grip Candy Sheer Case / Cover for iPhon...	8.78	1.0			
	Discount	Profit				
6983	0.0	2.2828				
	Row ID	Order ID	Order Date	Ship Date	Ship Mode	\

6984	6985.0	CA-2016-116526	2016-09-01	2016-09-05	Standard Class	
------	--------	----------------	------------	------------	----------------	--

	Customer ID	Customer Name	Segment	Country	City	State \
6984	JA-15970	Joseph Airdo	Consumer	United States	Detroit	30.18

	Postal Code	Region	Product ID	Category	Sub-Category \
6984	3.0	0	13.8828	Office Supplies	Appliances

	Product Name	Sales	Quantity \
6984	Eureka The Boss Plus 12-Amp Hard Box Upright V...	376.74	4.0

	Discount	Profit
6984	0.1	71.162

	Row ID	Order ID	Order Date	Ship Date	Ship Mode \
6985	6986.0	CA-2016-116526	2016-09-01	2016-09-05	Standard Class

	Customer ID	Customer Name	Segment	Country	City	State \
6985	JA-15970	Joseph Airdo	Consumer	United States	Detroit	24.1

	Postal Code	Region	Product ID	Category	Sub-Category \
6985	5.0	0	11.086	Office Supplies	Binders

	Product Name	Sales	Quantity	Discount	Profit
6985	GBC Plastic Binding Combs	29.52	4.0	0.0	14.4648

	Row ID	Order ID	Order Date	Ship Date	Ship Mode \
6986	6987.0	CA-2016-116526	2016-09-01	2016-09-05	Standard Class

	Customer ID	Customer Name	Segment	Country	City	State \
6986	JA-15970	Joseph Airdo	Consumer	United States	Detroit	8.78

	Postal Code	Region	Product ID	Category	Sub-Category \
6986	1.0	0	2.2828	Office Supplies	Art

	Product Name	Sales	Quantity	Discount	Profit
6986	Newell 315	11.96	2.0	0.0	2.99

	Row ID	Order ID	Order Date	Ship Date	Ship Mode \
6987	6988.0	CA-2016-116526	2016-09-01	2016-09-05	Standard Class

	Customer ID	Customer Name	Segment	Country	City	State \
6987	JA-15970	Joseph Airdo	Consumer	United States	Detroit	376.74

	Postal Code	Region	Product ID	Category	Sub-Category \
6987	4.0	0.1	71.162	Office Supplies	Binders

	Product Name	Sales	Quantity \
6987	Wilson Jones 1" Hanging DublLock Ring Binders	26.4	5.0

	Discount	Profit

```

6987      0.0  12.672
      Row ID      Order ID Order Date  Ship Date      Ship Mode Customer ID \
6988  6989.0  CA-2017-158561 2017-11-11 2017-11-16  Second Class  BB-11545

      Customer Name  Segment      Country      City State \
6988  Brenda Bowman  Corporate  United States  Fort Lauderdale  29.52

      Postal Code Region Product ID      Category Sub-Category \
6988      4.0      0      14.4648  Office Supplies  Appliances

      Product Name      Sales  Quantity  Discount \
6988  Hoover Upright Vacuum With Dirt Cup  1158.12      5.0      0.2

      Profit
6988  130.2885
      Row ID      Order ID Order Date  Ship Date      Ship Mode Customer ID \
6989  6990.0  CA-2017-165099 2017-12-11 2017-12-13  First Class  DK-13375

      Customer Name  Segment      Country      City State Postal Code \
6989  Dennis Kane  Consumer  United States  Abilene  11.96      2.0

      Region Product ID      Category Sub-Category \
6989      0      2.99  Office Supplies  Appliances

      Product Name      Sales  Quantity  Discount \
6989  Hoover Commercial Lightweight Upright Vacuum  1.392      2.0      0.8

      Profit
6989 -3.7584
      Row ID      Order ID Order Date  Ship Date      Ship Mode Customer ID \
6990  6991.0  CA-2015-109386 2015-11-08 2015-11-13  Second Class  RH-19600

      Customer Name  Segment      Country      City State Postal Code \
6990  Rob Haberlin  Consumer  United States  Hampton  26.4      5.0

      Region Product ID      Category Sub-Category \
6990      0      12.672  Office Supplies  Appliances

      Product Name      Sales  Quantity \
6990  Holmes Replacement Filter for HEPA Air Cleaner...  44.43      3.0

      Discount  Profit
6990      0.0  18.6606

```

```
[25]: df[df["State"].isnull()] #checking nul values in state
```

```
[25]:
```

	Row ID	Order ID	Order Date	Ship Date	Ship Mode	Customer ID	\
416	417.0	CA-2017-122105	2017-06-24	2017-06-28	Standard Class	CJ-12010	
423	424.0	CA-2017-125388	2017-10-19	2017-10-23	NaN	MP-17965	
428	429.0	CA-2017-152275	2017-10-01	2017-10-08	Standard Class	KH-16630	
430	431.0	US-2016-123750	2016-04-15	2016-04-21	Standard Class	RB-19795	
486	487.0	CA-2017-140963	2017-06-10	2017-06-13	First Class	MT-18070	
659	660.0	CA-2015-146563	2015-08-24	2015-08-28	Standard Class	CB-12025	

	Customer Name	Segment	Country	City	State	\
416	Caroline Jumper	Consumer	United States	Huntington Beach	NaN	
423	Michael Paige	Corporate	United States	Lawrence	NaN	
428	Ken Heidel	Corporate	United States	San Antonio	NaN	
430	Ross Baird	Home Office	United States	Gastonia	NaN	
486	Michelle Tran	Home Office	United States	Los Angeles	NaN	
659	NaN	Consumer	United States	Arlington	NaN	

	Postal Code	Region	Product ID	Category	Sub-Category	\
416	92646.0	West	OFF-AR-10004344	NaN	Art	
423	1841.0	East	OFF-ST-10000918	Office Supplies	Storage	
428	78207.0	Central	OFF-AR-10000369	NaN	Art	
430	28052.0	South	TEC-AC-10004659	Technology	NaN	
486	90045.0	West	TEC-PH-10001924	Technology	Phones	
659	76017.0	Central	OFF-ST-10001490	Office Supplies	NaN	

	Product Name	Sales	Quantity	\
416	Bulldog Vacuum Base Pencil Sharpener	95.920	8.0	
423	Crate-A-Files	32.700	3.0	
428	Design Ebony Sketching Pencil	6.672	6.0	
430	Imation Secure+ Hardware Encrypted USB 2.0 Fla...	408.744	7.0	
486	NaN	279.960	5.0	
659	Hot File 7-Pocket, Floor Stand	999.432	7.0	

	Discount	Profit
416	0.0	25.8984
423	0.0	8.5020
428	0.2	0.5004
430	0.2	76.6395
486	0.2	17.4975
659	0.2	124.9290

```
[26]: df[df["City"]=="Huntington Beach"].head(2) #checkin city name to replace nan
↳with respective state name
```

```
[26]:
```

	Row ID	Order ID	Order Date	Ship Date	Ship Mode	\
416	417.0	CA-2017-122105	2017-06-24	2017-06-28	Standard Class	
1890	1891.0	CA-2014-157623	2014-03-14	2014-03-18	Standard Class	

	Customer ID	Customer Name	Segment	Country	City \
416	CJ-12010	Caroline Jumper	Consumer	United States	Huntington Beach
1890	DK-13225	Dean Katz	Corporate	United States	Huntington Beach

	State	Postal Code	Region	Product ID	Category \
416	NaN	92646.0	West	OFF-AR-10004344	NaN
1890	California	92646.0	West	OFF-PA-10001204	Office Supplies

	Sub-Category	Product Name	Sales	Quantity \
416	Art	Bulldog Vacuum Base Pencil Sharpener	95.92	8.0
1890	Paper	Xerox 1972	10.56	2.0

	Discount	Profit
416	0.0	25.8984
1890	0.0	4.7520

```
[27]: df[df["State"]=="Maryland"].head(3)
```

```
[27]:
```

	Row ID	Order ID	Order Date	Ship Date	Ship Mode \
887	888.0	CA-2017-150707	2017-10-14	2017-10-19	Standard Class
1093	1094.0	CA-2015-165085	2015-12-27	2015-12-31	Standard Class
1094	1095.0	CA-2015-165085	2015-12-27	2015-12-31	Standard Class

	Customer ID	Customer Name	Segment	Country	City \
887	EL-13735	Ed Ludwig	Home Office	United States	Columbia
1093	BT-11485	Brad Thomas	Home Office	United States	Clinton
1094	BT-11485	Brad Thomas	Home Office	United States	Clinton

	State	Postal Code	Region	Product ID	Category \
887	Maryland	21044.0	East	OFF-BI-10001078	Office Supplies
1093	Maryland	20735.0	East	OFF-PA-10000605	Office Supplies
1094	Maryland	20735.0	East	OFF-AP-10002518	Office Supplies

	Sub-Category	Product Name	Sales \
887	Binders	Acco PRESSTEX Data Binder with Storage Hooks, ...	37.66
1093	Paper	Xerox 1950	28.90
1094	Appliances	Kensington 7 Outlet MasterPiece Power Center	355.96

	Quantity	Discount	Profit
887	7.0	0.0	18.4534
1093	5.0	0.0	14.1610
1094	2.0	0.0	103.2284

```
[28]: city_to_state = {
    'Phoenix': 'Arizona',
    'Houston': 'Texas',
    'Gilbert': 'Arizona',
```

```

'Jonesboro': 'Arkansas',
'Atlanta': 'Georgia',
'Los Angeles': 'California',
'New York City': 'New York',
'Detroit': 'Michigan',
'Fort Lauderdale': 'Florida',
'Hampton': 'Virginia',
'Arlington': 'Virginia',
'Gastonia': 'North Carolina',
'San Antonio': 'Texas',
'Lawrence': 'Massachusetts',
'Huntington Beach': 'California'
}

# Update 'State' column for cities in city_to_state dictionary
df['State'] = df['City'].map(city_to_state).fillna(df['State'])

print(df)

```

	Row ID	Order ID	Order Date	Ship Date	Ship Mode \
0	1.0	CA-2016-152156	2016-11-08	2016-11-11	Second Class
1	2.0	CA-2016-152156	2016-11-08	2016-11-11	Second Class
2	3.0	CA-2016-138688	2016-06-12	2016-06-16	Second Class
3	4.0	US-2015-108966	2015-10-11	2015-10-18	Standard Class
4	5.0	US-2015-108966	2015-10-11	2015-10-18	Standard Class
...
9989	NaN	NaN	NaT	NaT	Standard Class
9990	NaN	NaN	NaT	NaT	Second Class
9991	NaN	NaN	NaT	NaT	Standard Class
9992	NaN	NaN	NaT	NaT	Standard Class
9993	NaN	NaN	NaT	NaT	Second Class

	Customer ID	Customer Name	Segment	Country \
0	CG-12520	Claire Gute	Consumer	United States
1	CG-12520	Claire Gute	Consumer	United States
2	DV-13045	Darrin Van Huff	Corporate	United States
3	SO-20335	Sean O'Donnell	Consumer	United States
4	SO-20335	Sean O'Donnell	Consumer	United States
...
9989	SR-20425	Sharelle Roach	Home Office	United States
9990	AG-10330	Alex Grayson	Consumer	United States
9991	BP-11095	NaN	NaN	NaN
9992	JW-16075	NaN	NaN	NaN
9993	LH-16900	NaN	NaN	NaN

	City	State	Postal Code	Region	Product ID \
0	Henderson	Kentucky	42420.0	South	FUR-BO-10001798

1	Henderson	Kentucky	42420.0	South	FUR-CH-10000454
2	Los Angeles	California	90036.0	West	OFF-LA-10000240
3	Fort Lauderdale	Florida	33311.0	South	FUR-TA-10000577
4	Fort Lauderdale	Florida	33311.0	South	OFF-ST-10000760
...
9989	Tuscaloosa	Alabama	35401.0	South	FUR-CH-10002647
9990	Mesa	Arizona	85204.0	West	FUR-TA-10003008
9991	Jacksonville	North Carolina	28540.0	South	OFF-PA-10004071
9992	Chicago	Illinois	60610.0	Central	OFF-AP-10004980
9993	Columbus	Georgia	31907.0	South	FUR-FU-10000747

	Category	Sub-Category	\
0	Furniture	Bookcases	
1	Furniture	Chairs	
2	Office Supplies	Labels	
3	Furniture	Tables	
4	Office Supplies	Storage	
...	
9989	Furniture	Chairs	
9990	Furniture	Tables	
9991	Office Supplies	Paper	
9992	Office Supplies	Appliances	
9993	Furniture	Furnishings	

	Product Name	Sales	Quantity	\
0	Bush Somerset Collection Bookcase	261.9600	2.0	
1	Hon Deluxe Fabric Upholstered Stacking Chairs,...	731.9400	3.0	
2	Self-Adhesive Address Labels for Typewriters b...	14.6200	2.0	
3	Bretford CR4500 Series Slim Rectangular Table	957.5775	5.0	
4	Eldon Fold 'N Roll Cart System	22.3680	2.0	
...	
9989	Situations Contoured Folding Chairs, 4/Set	141.9600	2.0	
9990	Lesro Round Back Collection Coffee Table, End ...	182.5500	2.0	
9991	Eaton Premium Continuous-Feed Paper, 25% Cotto...	88.7680	2.0	
9992	3M Replacement Filter for Office Air Cleaner f...	53.0880	7.0	
9993	Tenex B1-RE Series Chair Mats for Low Pile Car...	275.8800	6.0	

	Discount	Profit
0	0.00	41.9136
1	0.00	219.5820
2	0.00	6.8714
3	0.45	-383.0310
4	0.20	2.5164
...
9989	0.00	35.4900
9990	0.50	-135.0870
9991	0.20	31.0688
9992	0.80	-108.8304

9993 0.00 46.8996

[9866 rows x 21 columns]

```
[29]: #checking changes
df[df["City"]=="Huntington Beach"].head(2)
```

```
[29]:      Row ID      Order ID Order Date  Ship Date      Ship Mode \
416    417.0  CA-2017-122105 2017-06-24 2017-06-28  Standard Class
1890   1891.0  CA-2014-157623 2014-03-14 2014-03-18  Standard Class

      Customer ID  Customer Name  Segment      Country      City \
416      CJ-12010  Caroline Jumper  Consumer  United States  Huntington Beach
1890      DK-13225      Dean Katz  Corporate  United States  Huntington Beach

      State  Postal Code Region      Product ID      Category \
416  California      92646.0  West  OFF-AR-10004344      NaN
1890  California      92646.0  West  OFF-PA-10001204  Office Supplies

      Sub-Category      Product Name  Sales  Quantity \
416      Art  Bulldog Vacuum Base Pencil Sharpener  95.92      8.0
1890      Paper      Xerox 1972  10.56      2.0

      Discount  Profit
416      0.0  25.8984
1890      0.0   4.7520
```

```
[30]: df["State"].isnull().sum() #cheeking number of nulls in state
```

```
[30]: 0
```

```
[31]: #now null also chnaged
df[df["Row ID"]==107]
```

```
[31]:      Row ID      Order ID Order Date  Ship Date      Ship Mode Customer ID \
106    107.0  CA-2017-119004 2017-11-23 2017-11-28  Standard Class      JM-15250

      Customer Name  Segment      Country      City      State \
106  Janet Martin  Consumer  United States  Charlotte  North Carolina

      Postal Code Region      Product ID      Category Sub-Category \
106      28205.0  South  TEC-AC-10003499  Technology  Accessories

      Product Name  Sales  Quantity \
106  Memorex Mini Travel Drive 8 GB USB 2.0 Flash D...  74.112      8.0

      Discount  Profit
```


106 0.2 17.6016

```
[32]: #drop 11.96 due to no name of state
df.drop(df[df["State"]==11.96].index,axis=0,inplace=True)
```

```
[33]: df["State"].unique() #now rechecking uniques values in state
```

```
[33]: array(['Kentucky', 'California', 'Florida', 'North Carolina',
        'Washington', 'Texas', 'Wisconsin', 'Utah', 'Nebraska',
        'Pennsylvania', 'Illinois', 'Minnesota', 'Michigan', 'Delaware',
        'Indiana', 'New York', 'Arizona', 'Virginia', 'Tennessee',
        'Alabama', 'South Carolina', 'Colorado', 'Iowa', 'Ohio',
        'Missouri', 'Oklahoma', 'New Mexico', 'Louisiana', 'Connecticut',
        'New Jersey', 'Oregon', 'Massachusetts', 'Georgia', 'Nevada',
        'Rhode Island', 'Mississippi', 'Arkansas', 'Montana',
        'New Hampshire', 'Maryland', 'District of Columbia', 'Kansas',
        'Vermont', 'Maine', 'South Dakota', 'Idaho', 'North Dakota',
        'Wyoming', 'West Virginia'], dtype=object)
```

```
[34]: df["Region"].unique() #hecking unique values in region(found some errors)
```

```
[34]: array(['South', 'West', 'Central', 'East', 0, 0.2, 0.7, 0.1], dtype=object)
```

```
[35]: df["Region"].isnull().sum() #checking nan values in region column
```

```
[35]: 0
```

```
[36]: a = [0, 0.2, 0.7, 0.1]
cities_by_region = {}

for i in a:
    filtered_df = df[df["Region"] == i]
    cities_by_region[i] = filtered_df["City"].tolist()

print(cities_by_region)
```

```
{0: ['Phoenix', 'New York City', 'New York City', 'Detroit', 'Detroit',
'Detroit', 'Detroit', 'Detroit', 'Fort Lauderdale', 'Hampton'], 0.2: ['Phoenix',
'Houston', 'Gilbert', 'Jonesboro', 'Atlanta', 'Los Angeles'], 0.7: ['Gilbert'],
0.1: ['Detroit']}
```

```
[37]: df[df["City"]=="Huntington Beach"].head(2) #checking region for city
```

```
[37]:      Row ID      Order ID Order Date  Ship Date      Ship Mode \
416      417.0    CA-2017-122105 2017-06-24 2017-06-28  Standard Class
1890    1891.0    CA-2014-157623 2014-03-14 2014-03-18  Standard Class
```

	Customer ID	Customer Name	Segment	Country	City \
416	CJ-12010	Caroline Jumper	Consumer	United States	Huntington Beach
1890	DK-13225	Dean Katz	Corporate	United States	Huntington Beach

	State	Postal Code	Region	Product ID	Category \
416	California	92646.0	West	OFF-AR-10004344	NaN
1890	California	92646.0	West	OFF-PA-10001204	Office Supplies

	Sub-Category	Product Name	Sales	Quantity \
416	Art	Bulldog Vacuum Base Pencil Sharpener	95.92	8.0
1890	Paper	Xerox 1972	10.56	2.0

	Discount	Profit
416	0.0	25.8984
1890	0.0	4.7520

```
[38]: city_to_region = {
    'Phoenix': 'West',
    'Houston': 'Central',
    'Gilbert': 'West',
    'Jonesboro': 'South',
    'Atlanta': 'South',
    'Los Angeles': 'West',
    'New York City': 'East',
    'Detroit': 'Central',
    'Fort Lauderdale': 'South',
    'Hampton': 'South'
}

# Update 'Region' column permanently
df['Region'] = df['City'].map(city_to_region).fillna(df['Region'])
```

```
[39]: df["Region"].unique() #checking unique values in region
```

```
[39]: array(['South', 'West', 'Central', 'East'], dtype=object)
```

```
[40]: df.shape #checking shape of dataframe
```

```
[40]: (9865, 21)
```

```
[41]: df.dropna(subset=["Order Date", "Ship Date"], inplace=True) #dropping null values
      ↪ in order date, ship date
```

```
[42]: df.isnull().sum() #checking null values in columns
```

```
[42]: Row ID          0
      Order ID       0
```

```

Order Date      0
Ship Date       0
Ship Mode       4
Customer ID     1
Customer Name   19
Segment        15
Country        18
City           0
State          0
Postal Code     0
Region         0
Product ID     2
Category       2
Sub-Category   4
Product Name   4
Sales          0
Quantity       0
Discount       0
Profit         13
dtype: int64

```

```

[43]: #data formating(creating new columns from existing data)
df["Order Date"] = pd.to_datetime(df["Order Date"], format="%d/%m/%Y")
df["Ship Date"] = pd.to_datetime(df["Ship Date"], format="%d/%m/%Y")

df["order_year"]=df["Order Date"].dt.year
df["order_month"]=df["Order Date"].dt.month
df["order_date"]=df["Order Date"].dt.day

df["Ship_year"]=df["Ship Date"].dt.year
df["Ship_month"]=df["Ship Date"].dt.month
df["Ship_date"]=df["Ship Date"].dt.day

```

```

[44]: df["Postal Code"].isnull().sum() #sum of null values in Postal code

```

```

[44]: 0

```

```

[45]: s=df["Postal Code"].astype(int) #unique postal code
s.unique()

```

```

[45]: array([42420, 90036, 33311, 90032, 28027, 98103, 76106, 53711, 84084,
          94109, 68025, 19140, 84057, 90049, 77095, 75080, 77041, 60540,
          32935, 55122, 48185, 19901, 47150, 10024, 12180, 90004, 60610,
          85234, 22153, 10009, 49201, 38109, 77070, 35601, 94122, 27707,
          60623, 29203, 55901, 80013, 28205, 60462, 10035, 50322, 43229,
          37620, 19805, 61701, 85023, 95661, 64055, 91104, 43055, 53132,
          85254, 95123, 98105, 98115, 73034, 90045, 19134, 88220, 78207,

```

77036, 62521, 71203, 6824, 75051, 80219, 75220, 37064, 90604,
 48601, 44256, 48227, 38401, 33614, 95051, 55044, 92037, 77506,
 94513, 27514, 7960, 45231, 94110, 90301, 97206, 33319, 80906,
 7109, 48180, 8701, 22204, 80004, 7601, 33710, 19143, 90805,
 92345, 37130, 78745, 1852, 31907, 6040, 78550, 85705, 62301,
 2038, 33024, 98198, 61604, 89115, 2886, 33180, 28403, 92646,
 40475, 80027, 1841, 39212, 48187, 10801, 28052, 32216, 47201,
 13021, 44312, 73071, 94521, 60068, 79109, 11757, 90008, 92024,
 77340, 14609, 72701, 92627, 80134, 30318, 64118, 59405, 48234,
 36116, 85204, 60653, 54302, 45503, 92804, 98270, 97301, 78041,
 19120, 75217, 43123, 10011, 48126, 31088, 94591, 55407, 92691,
 48307, 7060, 85635, 98661, 60505, 76017, 40214, 75081, 44105,
 75701, 27217, 22980, 19013, 27511, 32137, 10550, 48205, 33012,
 11572, 92105, 60201, 48183, 55016, 71111, 50315, 93534, 23223,
 28806, 92530, 68104, 98026, 92704, 53209, 41042, 44052, 7036,
 93905, 8901, 17602, 3301, 21044, 75043, 6360, 22304, 43615,
 87401, 92503, 90503, 78664, 92054, 33433, 23464, 92563, 28540,
 52601, 98502, 20016, 65109, 63376, 61107, 33142, 78521, 10701,
 94601, 28110, 20735, 30076, 72401, 47374, 94509, 33030, 46350,
 48911, 44221, 89502, 22801, 92025, 48073, 20852, 33065, 14215,
 33437, 39503, 93727, 27834, 11561, 35630, 31204, 52402, 2908,
 81001, 94533, 55106, 32725, 42071, 6457, 11520, 90660, 84604,
 84062, 30080, 24153, 44134, 36608, 2740, 75061, 8360, 85301,
 14304, 27360, 92683, 38301, 75019, 91767, 89031, 18103, 19711,
 85281, 92677, 8302, 2149, 13601, 54915, 98006, 75002, 79907,
 76051, 75007, 37167, 98031, 70506, 97224, 60076, 75023, 23434,
 46203, 7002, 43017, 28314, 27405, 21215, 53142, 66062, 98002,
 74133, 97756, 27604, 74403, 6450, 42104, 46614, 6010, 89015,
 99207, 76248, 45014, 32127, 97504, 22901, 59801, 33178, 29501,
 97477, 32712, 19601, 80020, 65807, 7501, 73120, 23320, 79424,
 65203, 37604, 36830, 92404, 1453, 59715, 85345, 44107, 8861,
 91761, 91730, 56560, 75150, 92374, 95207, 32174, 94086, 3820,
 17403, 77840, 63116, 2169, 95336, 44240, 76903, 84106, 35810,
 37918, 72209, 48146, 43302, 80122, 5408, 4401, 38671, 47362,
 48640, 57103, 80525, 47905, 37042, 95823, 91360, 2148, 1040,
 87105, 89431, 92236, 60126, 7055, 29406, 23602, 14701, 46544,
 43402, 7090, 92253, 32303, 37211, 98226, 60098, 76117, 60090,
 29483, 71901, 80112, 43130, 88001, 35244, 75034, 95687, 84107,
 53186, 93309, 33068, 45373, 78415, 90278, 32839, 7050, 70601,
 60035, 11550, 46060, 55124, 29464, 48310, 54703, 78577, 59102,
 97030, 37421, 83642, 92307, 60440, 33801, 55369, 95695, 77489,
 77581, 94403, 49505, 93277, 66212, 92592, 92399, 2151, 77301,
 60477, 52001, 48127, 87505, 28601, 60188, 56301, 33161, 46226,
 33317, 34952, 29730, 79762, 53214, 91911, 66502, 16602, 80229,
 61821, 47401, 71854, 78539, 77520, 46142, 90712, 2895, 54880,
 76021, 98042, 74012, 33023, 33021, 77536, 67212, 78501, 52240,
 83704, 2920, 61032, 77642, 95610, 75056, 98052, 32114, 86442,

```

46368, 58103, 46514, 91776, 45011, 33063, 30328, 44060, 73505,
23666, 13440, 54601, 83501, 39401, 94526, 48858, 84321, 6708,
30605, 4240, 61832, 85323, 30062, 85364, 54401, 99301, 60302,
32503, 77573, 20877, 84043, 35401, 92553, 40324, 80538, 85224,
59601, 63122, 76706, 48066, 60423, 18018, 55113, 68801, 55125,
48237, 72756, 88101, 33458, 93101, 75104, 68701, 84020, 48104,
91941, 83201, 49423, 6460, 60089, 92630, 96003, 95928, 13501,
72032, 82001, 42301, 83605, 70065, 3060, 38134, 94061, 37087,
93454, 60016, 98632, 37075, 50701, 2138, 60067, 1915, 97405,
93030, 98059, 60025, 33445, 80022, 77590, 27893, 87124, 27534,
98208, 90640, 92020, 33407, 5, 2, 4, 9, 3,
1, 61761, 60174, 93010, 97123, 91505, 95351, 67846, 8401,
80501, 95616, 26003, 95037, 7011, 53081, 30344, 57701, 1810,
34741, 6484, 6810, 52302, 32771, 78666, 80634, 76063, 44035,
83301, 63301, 33134, 60441, 1752, 20707, 77803, 71603, 57401,
21740, 7017, 60004, 60543, 77705, 55433, 92672, 94568, 93405,
72762, 95240, 77571, 45040, 30188])

```

```

[46]: a=[5,2,4,9,3,1] #errors in postal code
      c=df[df["Postal Code"].isin(a)]
      c

```

```

[46]: Row ID      Order ID Order Date  Ship Date      Ship Mode \
6972  6973.0    CA-2017-153822 2017-09-19 2017-09-25  Standard Class
6973  6974.0    CA-2017-153822 2017-09-19 2017-09-25  Standard Class
6974  6975.0    CA-2017-146185 2017-09-15 2017-09-19  Standard Class
6975  6976.0    CA-2015-112144 2015-06-28 2015-07-02  Standard Class
6976  6977.0    CA-2015-112144 2015-06-28 2015-07-02  Standard Class
6977  6978.0    US-2016-119298 2016-11-25 2016-11-28    First Class
6978  6979.0    CA-2017-155159 2017-11-25 2017-11-29    Second Class
6979  6980.0    CA-2017-149076 2017-01-14 2017-01-19  Standard Class
6980  6981.0    CA-2014-146990 2014-11-07 2014-11-08    First Class
6981  6982.0    CA-2014-146990 2014-11-07 2014-11-08    First Class
6982  6983.0    CA-2016-116526 2016-09-01 2016-09-05  Standard Class
6983  6984.0    CA-2016-116526 2016-09-01 2016-09-05  Standard Class
6984  6985.0    CA-2016-116526 2016-09-01 2016-09-05  Standard Class
6985  6986.0    CA-2016-116526 2016-09-01 2016-09-05  Standard Class
6986  6987.0    CA-2016-116526 2016-09-01 2016-09-05  Standard Class
6987  6988.0    CA-2016-116526 2016-09-01 2016-09-05  Standard Class
6988  6989.0    CA-2017-158561 2017-11-11 2017-11-16    Second Class
6990  6991.0    CA-2015-109386 2015-11-08 2015-11-13    Second Class

      Customer ID      Customer Name      Segment      Country \
6972      AB-10105      Adrian Barton      Consumer      United States
6973      AB-10105      Adrian Barton      Consumer      United States
6974      CC-12145      Charles Crestani      Consumer      United States
6975      CY-12745      Craig Yedwab      Corporate      United States

```

6976	CY-12745	Craig Yedwab	Corporate	United States
6977	EP-13915	Emily Phan	Consumer	United States
6978	DL-13315	Delfina Latchford	Consumer	United States
6979	SO-20335	Sean O'Donnell	Consumer	United States
6980	BP-11095	Bart Pistole	Corporate	United States
6981	BP-11095	Bart Pistole	Corporate	United States
6982	JA-15970	Joseph Airdo	Consumer	United States
6983	JA-15970	Joseph Airdo	Consumer	United States
6984	JA-15970	Joseph Airdo	Consumer	United States
6985	JA-15970	Joseph Airdo	Consumer	United States
6986	JA-15970	Joseph Airdo	Consumer	United States
6987	JA-15970	Joseph Airdo	Consumer	United States
6988	BB-11545	Brenda Bowman	Corporate	United States
6990	RH-19600	Rob Haberlin	Consumer	United States

	City	State	...	\
6972	Phoenix	Arizona	...	
6973	Phoenix	Arizona	...	
6974	Houston	Texas	...	
6975	Gilbert	Arizona	...	
6976	Gilbert	Arizona	...	
6977	Jonesboro	Arkansas	...	
6978	Atlanta	Georgia	...	
6979	Los Angeles	California	...	
6980	New York City	New York	...	
6981	New York City	New York	...	
6982	Detroit	Michigan	...	
6983	Detroit	Michigan	...	
6984	Detroit	Michigan	...	
6985	Detroit	Michigan	...	
6986	Detroit	Michigan	...	
6987	Detroit	Michigan	...	
6988	Fort Lauderdale	Florida	...	
6990	Hampton	Virginia	...	

	Product Name	Sales	Quantity	\
6972	Polycom VoiceStation 500 Conference phone	471.920	2.0	
6973	Plastic Binding Combs	18.180	4.0	
6974	Prismacolor Color Pencil Set	31.744	2.0	
6975	Avery 501	5.904	2.0	
6976	Electrix Halogen Magnifier Lamp	621.760	4.0	
6977	OtterBox Defender Series Case - Samsung Galaxy S4	59.980	2.0	
6978	Wirebound Message Book, 4 per Page	48.870	9.0	
6979	Xerox 19	154.900	5.0	
6980	Binder Clips by OIC	5.920	4.0	
6981	Riverleaf Stik-Withit Designer Note Cubes	30.180	3.0	
6982	Wilson Jones Turn Tabs Binder Tool for Ring Bi...	24.100	5.0	

6983	Belkin Grip Candy Sheer Case / Cover for iPhon...	8.780	1.0
6984	Eureka The Boss Plus 12-Amp Hard Box Upright V...	376.740	4.0
6985	GBC Plastic Binding Combs	29.520	4.0
6986	Newell 315	11.960	2.0
6987	Wilson Jones 1" Hanging DublLock Ring Binders	26.400	5.0
6988	Hoover Upright Vacuum With Dirt Cup	1158.120	5.0
6990	Holmes Replacement Filter for HEPA Air Cleaner...	44.430	3.0

	Discount	Profit	order_year	order_month	order_date	Ship_year	\
6972	0.2	29.4950	2017	9	19	2017	
6973	0.7	-13.9380	2017	9	19	2017	
6974	0.2	8.3328	2017	9	15	2017	
6975	0.2	1.9926	2015	6	28	2015	
6976	0.2	46.6320	2015	6	28	2015	
6977	0.0	17.9940	2016	11	25	2016	
6978	0.0	23.9463	2017	11	25	2017	
6979	0.0	69.7050	2017	1	14	2017	
6980	0.0	2.8416	2014	11	7	2014	
6981	0.0	13.8828	2014	11	7	2014	
6982	0.0	11.0860	2016	9	1	2016	
6983	0.0	2.2828	2016	9	1	2016	
6984	0.1	71.1620	2016	9	1	2016	
6985	0.0	14.4648	2016	9	1	2016	
6986	0.0	2.9900	2016	9	1	2016	
6987	0.0	12.6720	2016	9	1	2016	
6988	0.2	130.2885	2017	11	11	2017	
6990	0.0	18.6606	2015	11	8	2015	

	Ship_month	Ship_date
6972	9	25
6973	9	25
6974	9	19
6975	7	2
6976	7	2
6977	11	28
6978	11	29
6979	1	19
6980	11	8
6981	11	8
6982	9	5
6983	9	5
6984	9	5
6985	9	5
6986	9	5
6987	9	5
6988	11	16
6990	11	13

[18 rows x 27 columns]

```
[47]: #finding most frequent postal code for city
cities = ["Phoenix", "Houston", "Gilbert", "Jonesboro", "Atlanta", "Los Angeles", "New York City", "Detroit", "Fort Lauderdale", "Abilene", "Hampton"]

for city in cities:
    filtered_df = df[df["City"] == city]
    if filtered_df.empty:
        print(f"No rows found for city: {city}")
    else:
        postal_code_mode = filtered_df["Postal Code"].mode()[0]
        print(f"Most frequent postal code for city {city}: {postal_code_mode}")
```

```
Most frequent postal code for city Phoenix: 85023.0
Most frequent postal code for city Houston: 77041.0
Most frequent postal code for city Gilbert: 85234.0
Most frequent postal code for city Jonesboro: 72401.0
Most frequent postal code for city Atlanta: 30318.0
Most frequent postal code for city Los Angeles: 90049.0
Most frequent postal code for city New York City: 10035.0
Most frequent postal code for city Detroit: 48227.0
Most frequent postal code for city Fort Lauderdale: 33311.0
No rows found for city: Abilene
Most frequent postal code for city Hampton: 23666.0
```

```
[48]: city_to_region = {
    'Phoenix': 85023.0,
    'Houston': 77041.0,
    'Gilbert': 85234.0,
    'Jonesboro': 72401.0,
    'Atlanta': 30318.0,
    'Los Angeles': 90049.0,
    'New York City': 10035.0,
    'Detroit': 48227.0,
    'Fort Lauderdale': 33311.0,
    'Abilene': 77041.0,
    'Hampton': 23666.0
}

# Update 'Region' column permanently
c['Postal Code'] = c['City'].map(city_to_region).fillna(c['Postal Code'])
```

C:\Users\ganesh\AppData\Local\Temp\ipykernel_49052\1011837133.py:16:
SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.
Try using `.loc[row_indexer,col_indexer] = value` instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
`c['Postal Code'] = c['City'].map(city_to_region).fillna(c['Postal Code'])`

[49]: c

```
[49]:      Row ID      Order ID Order Date  Ship Date      Ship Mode \
6972  6973.0  CA-2017-153822 2017-09-19 2017-09-25  Standard Class
6973  6974.0  CA-2017-153822 2017-09-19 2017-09-25  Standard Class
6974  6975.0  CA-2017-146185 2017-09-15 2017-09-19  Standard Class
6975  6976.0  CA-2015-112144 2015-06-28 2015-07-02  Standard Class
6976  6977.0  CA-2015-112144 2015-06-28 2015-07-02  Standard Class
6977  6978.0  US-2016-119298 2016-11-25 2016-11-28    First Class
6978  6979.0  CA-2017-155159 2017-11-25 2017-11-29    Second Class
6979  6980.0  CA-2017-149076 2017-01-14 2017-01-19  Standard Class
6980  6981.0  CA-2014-146990 2014-11-07 2014-11-08    First Class
6981  6982.0  CA-2014-146990 2014-11-07 2014-11-08    First Class
6982  6983.0  CA-2016-116526 2016-09-01 2016-09-05  Standard Class
6983  6984.0  CA-2016-116526 2016-09-01 2016-09-05  Standard Class
6984  6985.0  CA-2016-116526 2016-09-01 2016-09-05  Standard Class
6985  6986.0  CA-2016-116526 2016-09-01 2016-09-05  Standard Class
6986  6987.0  CA-2016-116526 2016-09-01 2016-09-05  Standard Class
6987  6988.0  CA-2016-116526 2016-09-01 2016-09-05  Standard Class
6988  6989.0  CA-2017-158561 2017-11-11 2017-11-16    Second Class
6990  6991.0  CA-2015-109386 2015-11-08 2015-11-13    Second Class
```

```
      Customer ID      Customer Name      Segment      Country \
6972    AB-10105      Adrian Barton    Consumer    United States
6973    AB-10105      Adrian Barton    Consumer    United States
6974    CC-12145    Charles Crestani    Consumer    United States
6975    CY-12745      Craig Yedwab    Corporate    United States
6976    CY-12745      Craig Yedwab    Corporate    United States
6977    EP-13915      Emily Phan      Consumer    United States
6978    DL-13315    Delfina Latchford    Consumer    United States
6979    SO-20335    Sean O'Donnell      Consumer    United States
6980    BP-11095      Bart Pistole    Corporate    United States
6981    BP-11095      Bart Pistole    Corporate    United States
6982    JA-15970      Joseph Airdo      Consumer    United States
6983    JA-15970      Joseph Airdo      Consumer    United States
6984    JA-15970      Joseph Airdo      Consumer    United States
6985    JA-15970      Joseph Airdo      Consumer    United States
6986    JA-15970      Joseph Airdo      Consumer    United States
6987    JA-15970      Joseph Airdo      Consumer    United States
6988    BB-11545      Brenda Bowman    Corporate    United States
```

6990	RH-19600	Rob Haberlin	Consumer	United States
------	----------	--------------	----------	---------------

	City	State	...	\
6972	Phoenix	Arizona	...	
6973	Phoenix	Arizona	...	
6974	Houston	Texas	...	
6975	Gilbert	Arizona	...	
6976	Gilbert	Arizona	...	
6977	Jonesboro	Arkansas	...	
6978	Atlanta	Georgia	...	
6979	Los Angeles	California	...	
6980	New York City	New York	...	
6981	New York City	New York	...	
6982	Detroit	Michigan	...	
6983	Detroit	Michigan	...	
6984	Detroit	Michigan	...	
6985	Detroit	Michigan	...	
6986	Detroit	Michigan	...	
6987	Detroit	Michigan	...	
6988	Fort Lauderdale	Florida	...	
6990	Hampton	Virginia	...	

	Product Name	Sales	Quantity	\
6972	Polycom VoiceStation 500 Conference phone	471.920	2.0	
6973	Plastic Binding Combs	18.180	4.0	
6974	Prismacolor Color Pencil Set	31.744	2.0	
6975	Avery 501	5.904	2.0	
6976	Electrix Halogen Magnifier Lamp	621.760	4.0	
6977	OtterBox Defender Series Case - Samsung Galaxy S4	59.980	2.0	
6978	Wirebound Message Book, 4 per Page	48.870	9.0	
6979	Xerox 19	154.900	5.0	
6980	Binder Clips by OIC	5.920	4.0	
6981	Riverleaf Stik-Withit Designer Note Cubes	30.180	3.0	
6982	Wilson Jones Turn Tabs Binder Tool for Ring Bi...	24.100	5.0	
6983	Belkin Grip Candy Sheer Case / Cover for iPhon...	8.780	1.0	
6984	Eureka The Boss Plus 12-Amp Hard Box Upright V...	376.740	4.0	
6985	GBC Plastic Binding Combs	29.520	4.0	
6986	Newell 315	11.960	2.0	
6987	Wilson Jones 1" Hanging DublLock Ring Binders	26.400	5.0	
6988	Hoover Upright Vacuum With Dirt Cup	1158.120	5.0	
6990	Holmes Replacement Filter for HEPA Air Cleaner...	44.430	3.0	

	Discount	Profit	order_year	order_month	order_date	Ship_year	\
6972	0.2	29.4950	2017	9	19	2017	
6973	0.7	-13.9380	2017	9	19	2017	
6974	0.2	8.3328	2017	9	15	2017	
6975	0.2	1.9926	2015	6	28	2015	

6976	0.2	46.6320	2015	6	28	2015
6977	0.0	17.9940	2016	11	25	2016
6978	0.0	23.9463	2017	11	25	2017
6979	0.0	69.7050	2017	1	14	2017
6980	0.0	2.8416	2014	11	7	2014
6981	0.0	13.8828	2014	11	7	2014
6982	0.0	11.0860	2016	9	1	2016
6983	0.0	2.2828	2016	9	1	2016
6984	0.1	71.1620	2016	9	1	2016
6985	0.0	14.4648	2016	9	1	2016
6986	0.0	2.9900	2016	9	1	2016
6987	0.0	12.6720	2016	9	1	2016
6988	0.2	130.2885	2017	11	11	2017
6990	0.0	18.6606	2015	11	8	2015

	Ship_month	Ship_date
6972	9	25
6973	9	25
6974	9	19
6975	7	2
6976	7	2
6977	11	28
6978	11	29
6979	1	19
6980	11	8
6981	11	8
6982	9	5
6983	9	5
6984	9	5
6985	9	5
6986	9	5
6987	9	5
6988	11	16
6990	11	13

[18 rows x 27 columns]

```
[50]: df = pd.concat([df, c], ignore_index=True) #concatinating the new dataframe
      ↪with existing dataframe df
```

```
[51]: df.shape #shape of data frame
```

```
[51]: (9872, 27)
```

```
[52]: df.drop(df[df["Postal Code"].isin(a)].index,inplace=True) #dropping rows that
      ↪have error in postal code after imputing with correct values
```

```
[53]: df[df["Row ID"]==6973]
```

```
[53]:      Row ID      Order ID Order Date  Ship Date      Ship Mode \
9854  6973.0  CA-2017-153822 2017-09-19 2017-09-25  Standard Class

      Customer ID Customer Name  Segment      Country      City      State \
9854    AB-10105  Adrian Barton  Consumer  United States  Phoenix  Arizona

      ...                                     Product Name  Sales Quantity \
9854  ...  Polycom VoiceStation 500 Conference phone  471.92      2.0

      Discount  Profit order_year  order_month  order_date  Ship_year \
9854         0.2  29.495         2017           9          19         2017

      Ship_month  Ship_date
9854           9         25

[1 rows x 27 columns]
```

```
[54]: df["Postal Code"].unique() #unique values in data frame
```

```
[54]: array([42420., 90036., 33311., 90032., 28027., 98103., 76106., 53711.,
      84084., 94109., 68025., 19140., 84057., 90049., 77095., 75080.,
      77041., 60540., 32935., 55122., 48185., 19901., 47150., 10024.,
      12180., 90004., 60610., 85234., 22153., 10009., 49201., 38109.,
      77070., 35601., 94122., 27707., 60623., 29203., 55901., 80013.,
      28205., 60462., 10035., 50322., 43229., 37620., 19805., 61701.,
      85023., 95661., 64055., 91104., 43055., 53132., 85254., 95123.,
      98105., 98115., 73034., 90045., 19134., 88220., 78207., 77036.,
      62521., 71203., 6824., 75051., 80219., 75220., 37064., 90604.,
      48601., 44256., 48227., 38401., 33614., 95051., 55044., 92037.,
      77506., 94513., 27514., 7960., 45231., 94110., 90301., 97206.,
      33319., 80906., 7109., 48180., 8701., 22204., 80004., 7601.,
      33710., 19143., 90805., 92345., 37130., 78745., 1852., 31907.,
      6040., 78550., 85705., 62301., 2038., 33024., 98198., 61604.,
      89115., 2886., 33180., 28403., 92646., 40475., 80027., 1841.,
      39212., 48187., 10801., 28052., 32216., 47201., 13021., 44312.,
      73071., 94521., 60068., 79109., 11757., 90008., 92024., 77340.,
      14609., 72701., 92627., 80134., 30318., 64118., 59405., 48234.,
      36116., 85204., 60653., 54302., 45503., 92804., 98270., 97301.,
      78041., 19120., 75217., 43123., 10011., 48126., 31088., 94591.,
      55407., 92691., 48307., 7060., 85635., 98661., 60505., 76017.,
      40214., 75081., 44105., 75701., 27217., 22980., 19013., 27511.,
      32137., 10550., 48205., 33012., 11572., 92105., 60201., 48183.,
      55016., 71111., 50315., 93534., 23223., 28806., 92530., 68104.,
      98026., 92704., 53209., 41042., 44052., 7036., 93905., 8901.,
      17602., 3301., 21044., 75043., 6360., 22304., 43615., 87401.,
```

92503., 90503., 78664., 92054., 33433., 23464., 92563., 28540.,
52601., 98502., 20016., 65109., 63376., 61107., 33142., 78521.,
10701., 94601., 28110., 20735., 30076., 72401., 47374., 94509.,
33030., 46350., 48911., 44221., 89502., 22801., 92025., 48073.,
20852., 33065., 14215., 33437., 39503., 93727., 27834., 11561.,
35630., 31204., 52402., 2908., 81001., 94533., 55106., 32725.,
42071., 6457., 11520., 90660., 84604., 84062., 30080., 24153.,
44134., 36608., 2740., 75061., 8360., 85301., 14304., 27360.,
92683., 38301., 75019., 91767., 89031., 18103., 19711., 85281.,
92677., 8302., 2149., 13601., 54915., 98006., 75002., 79907.,
76051., 75007., 37167., 98031., 70506., 97224., 60076., 75023.,
23434., 46203., 7002., 43017., 28314., 27405., 21215., 53142.,
66062., 98002., 74133., 97756., 27604., 74403., 6450., 42104.,
46614., 6010., 89015., 99207., 76248., 45014., 32127., 97504.,
22901., 59801., 33178., 29501., 97477., 32712., 19601., 80020.,
65807., 7501., 73120., 23320., 79424., 65203., 37604., 36830.,
92404., 1453., 59715., 85345., 44107., 8861., 91761., 91730.,
56560., 75150., 92374., 95207., 32174., 94086., 3820., 17403.,
77840., 63116., 2169., 95336., 44240., 76903., 84106., 35810.,
37918., 72209., 48146., 43302., 80122., 5408., 4401., 38671.,
47362., 48640., 57103., 80525., 47905., 37042., 95823., 91360.,
2148., 1040., 87105., 89431., 92236., 60126., 7055., 29406.,
23602., 14701., 46544., 43402., 7090., 92253., 32303., 37211.,
98226., 60098., 76117., 60090., 29483., 71901., 80112., 43130.,
88001., 35244., 75034., 95687., 84107., 53186., 93309., 33068.,
45373., 78415., 90278., 32839., 7050., 70601., 60035., 11550.,
46060., 55124., 29464., 48310., 54703., 78577., 59102., 97030.,
37421., 83642., 92307., 60440., 33801., 55369., 95695., 77489.,
77581., 94403., 49505., 93277., 66212., 92592., 92399., 2151.,
77301., 60477., 52001., 48127., 87505., 28601., 60188., 56301.,
33161., 46226., 33317., 34952., 29730., 79762., 53214., 91911.,
66502., 16602., 80229., 61821., 47401., 71854., 78539., 77520.,
46142., 90712., 2895., 54880., 76021., 98042., 74012., 33023.,
33021., 77536., 67212., 78501., 52240., 83704., 2920., 61032.,
77642., 95610., 75056., 98052., 32114., 86442., 46368., 58103.,
46514., 91776., 45011., 33063., 30328., 44060., 73505., 23666.,
13440., 54601., 83501., 39401., 94526., 48858., 84321., 6708.,
30605., 4240., 61832., 85323., 30062., 85364., 54401., 99301.,
60302., 32503., 77573., 20877., 84043., 35401., 92553., 40324.,
80538., 85224., 59601., 63122., 76706., 48066., 60423., 18018.,
55113., 68801., 55125., 48237., 72756., 88101., 33458., 93101.,
75104., 68701., 84020., 48104., 91941., 83201., 49423., 6460.,
60089., 92630., 96003., 95928., 13501., 72032., 82001., 42301.,
83605., 70065., 3060., 38134., 94061., 37087., 93454., 60016.,
98632., 37075., 50701., 2138., 60067., 1915., 97405., 93030.,
98059., 60025., 33445., 80022., 77590., 27893., 87124., 27534.,
98208., 90640., 92020., 33407., 61761., 60174., 93010., 97123.,

```
91505., 95351., 67846., 8401., 80501., 95616., 26003., 95037.,
7011., 53081., 30344., 57701., 1810., 34741., 6484., 6810.,
52302., 32771., 78666., 80634., 76063., 44035., 83301., 63301.,
33134., 60441., 1752., 20707., 77803., 71603., 57401., 21740.,
7017., 60004., 60543., 77705., 55433., 92672., 94568., 93405.,
72762., 95240., 77571., 45040., 30188.]
```

```
[55]: df[df["Postal Code"].isin(a)] #rechecking is there any error in postal code
```

```
[55]: Empty DataFrame
Columns: [Row ID, Order ID, Order Date, Ship Date, Ship Mode, Customer ID,
Customer Name, Segment, Country, City, State, Postal Code, Region, Product ID,
Category, Sub-Category, Product Name, Sales, Quantity, Discount, Profit,
order_year, order_month, order_date, Ship_year, Ship_month, Ship_date]
Index: []

[0 rows x 27 columns]
```

```
[56]: df.isnull().sum() #checking sum of null values
```

```
[56]: Row ID          0
Order ID          0
Order Date        0
Ship Date         0
Ship Mode         4
Customer ID       1
Customer Name     19
Segment          15
Country           18
City              0
State             0
Postal Code       0
Region            0
Product ID        2
Category          2
Sub-Category      4
Product Name      4
Sales             0
Quantity          0
Discount          0
Profit           13
order_year        0
order_month       0
order_date        0
Ship_year         0
Ship_month        0
Ship_date         0
```

dtype: int64

```
[57]: # filled with united states as it has only one country
df['Country'].fillna('United States', inplace=True)
```

```
[58]: #filling nan values in object dtype columns with mode
for i in df.select_dtypes(include="object"):
    df[i].fillna(df[i].mode()[0],inplace=True)
```

```
[59]: #used knnimputer to impute null values in numerical columns except sales(as it
      ↪is dependent variable)
impute=KNNImputer()
for i in df.select_dtypes(include="number").columns:
    if i != 'Sales':
        df[i] = impute.fit_transform(df[[i]])
```

```
[60]: df.isnull().sum() #checking sum of null values
```

```
[60]: Row ID          0
      Order ID       0
      Order Date     0
      Ship Date      0
      Ship Mode       0
      Customer ID    0
      Customer Name   0
      Segment        0
      Country        0
      City           0
      State          0
      Postal Code    0
      Region         0
      Product ID     0
      Category       0
      Sub-Category   0
      Product Name   0
      Sales          0
      Quantity       0
      Discount       0
      Profit         0
      order_year     0
      order_month    0
      order_date     0
      Ship_year      0
      Ship_month     0
      Ship_date      0
      dtype: int64
```

```
[61]: df.dtypes #checking dtypes
```

```
[61]: Row ID          float64
      Order ID       object
      Order Date     datetime64[ns]
      Ship Date      datetime64[ns]
      Ship Mode       object
      Customer ID     object
      Customer Name   object
      Segment        object
      Country         object
      City            object
      State           object
      Postal Code     float64
      Region          object
      Product ID      object
      Category        object
      Sub-Category    object
      Product Name     object
      Sales           float64
      Quantity        float64
      Discount        float64
      Profit          float64
      order_year      float64
      order_month     float64
      order_date      float64
      Ship_year       float64
      Ship_month      float64
      Ship_date       float64
      dtype: object
```

```
[62]: df.dropna(inplace=True) #dropping all nan value rows
```

```
[63]: df.isnull().sum() #checking null values in columns
```

```
[63]: Row ID          0
      Order ID        0
      Order Date      0
      Ship Date       0
      Ship Mode       0
      Customer ID     0
      Customer Name   0
      Segment         0
      Country         0
      City            0
      State           0
      Postal Code     0
```



```

Region          0
Product ID      0
Category        0
Sub-Category    0
Product Name    0
Sales           0
Quantity        0
Discount        0
Profit          0
order_year      0
order_month     0
order_date      0
Ship_year       0
Ship_month      0
Ship_date       0
dtype: int64

```

```
[64]: df.shape #checking shape
```

```
[64]: (9854, 27)
```

```
[65]: ((9994-9854)/9994)*100 #checking perecntage of rows i deleted
```

```
[65]: 1.4008405043025816
```

```
[66]: df.duplicated().sum() #checking number of duplicate rows
```

```
[66]: 10
```

```
[67]: df.drop_duplicates(inplace=True) #dropping duplicate rows
```

```
[68]: for i in df.columns:
        print(f"{i}:{df[i].unique()}\n") #checking all unique values for each
        ↪column
```

```
Row ID:[1.000e+00 2.000e+00 3.000e+00 ... 6.988e+03 6.989e+03 6.991e+03]
```

```
Order ID:['CA-2016-152156' 'CA-2016-138688' 'US-2015-108966' ...
'CA-2014-146990'
'CA-2016-116526' 'CA-2017-158561']
```

```
Order Date:['2016-11-08T00:00:00.000000000' '2016-06-12T00:00:00.000000000'
'2015-10-11T00:00:00.000000000' ... '2016-06-03T00:00:00.000000000'
'2015-04-12T00:00:00.000000000' '2014-01-21T00:00:00.000000000']
```

```
Ship Date:['2016-11-11T00:00:00.000000000' '2016-06-16T00:00:00.000000000'
'2015-10-18T00:00:00.000000000' ... '2015-05-23T00:00:00.000000000']
```

'2014-01-23T00:00:00.000000000' '2017-03-03T00:00:00.000000000']

Ship Mode:['Second Class' 'Standard Class' 'First Class' 'Same Day']

Customer ID:['CG-12520' 'DV-13045' 'SO-20335' 'BH-11710' 'AA-10480' 'IM-15070'
'HP-14815' 'PK-19075' 'AG-10270' 'ZD-21925' 'KB-16585' 'SF-20065'
'EB-13870' 'EH-13945' 'TB-21520' 'MA-17560' 'GH-14485' 'SN-20710'
'LC-16930' 'RA-19885' 'ES-14080' 'ON-18715' 'PO-18865' 'LH-16900'
'DP-13000' 'JM-15265' 'TB-21055' 'KM-16720' 'PS-18970' 'BS-11590'
'KD-16270' 'HM-14980' 'JE-15745' 'KB-16600' 'SC-20770' 'DN-13690'
'JC-16105' 'CS-12400' 'PG-18895' 'GM-14455' 'JS-15685' 'RB-19465'
'GZ-14470' 'LC-16870' 'JM-15250' 'PA-19060' 'CV-12805' 'CL-12565'
'RC-19960' 'DK-13090' 'GG-14650' 'SC-20725' 'AD-10180' 'PF-19165'
'TS-21610' 'LS-16975' 'DW-13585' 'LC-16885' 'JD-15895' 'SH-19975'
'SG-20080' 'HA-14920' 'MG-17680' 'JE-16165' 'TW-21025' 'SP-20650'
'NK-18490' 'DB-13060' 'NP-18670' 'TT-21070' 'EM-13960' 'RD-19900'
'MJ-17740' 'BM-11140' 'CS-12130' 'JB-15400' 'SJ-20500' 'JK-15640'
'DK-13150' 'RM-19675' 'SK-19990' 'FM-14290' 'AM-10360' 'MP-17470'
'BS-11755' 'LC-17140' 'HK-14890' 'LE-16810' 'JH-15985' 'MS-17980'
'VW-21775' 'JH-15910' 'DS-13180' 'VD-21670' 'EA-14035' 'DB-13120'
'KL-16645' 'DW-13480' 'LH-17155' 'KC-16540' 'DL-13315' 'DR-12880'
'CC-12670' 'DL-13600' 'SB-20290' 'RC-19825' 'AH-10210' 'CB-12535'
'CA-12310' 'KH-16690' 'BB-10990' 'AG-10495' 'JO-15280' 'AH-10195'
'NZ-18565' 'KL-16555' 'AS-10225' 'CR-12625' 'SH-20395' 'BP-11185'
'TS-21205' 'AG-10525' 'SP-20860' 'NM-18445' 'FA-14230' 'GK-14620'
'DJ-13510' 'PO-18850' 'JL-15850' 'DB-13615' 'CC-12550' 'TD-20995'
'AB-10060' 'JL-15505' 'VB-21745' 'KW-16435' 'JD-16060' 'MK-17905'
'GT-14755' 'AG-10900' 'MM-18280' 'AR-10405' 'RA-19915' 'AS-10285'
'LA-16780' 'DO-13435' 'DK-13225' 'NG-18430' 'MV-18190' 'JG-15115'
'BP-11095' 'VP-21730' 'SS-20140' 'AG-10675' 'LF-17185' 'RF-19840'
'KH-16510' 'KC-16675' 'CJ-12010' 'PB-19150' 'MP-17965' 'NF-18385'
'SD-20485' 'KH-16630' 'RB-19795' 'MK-18160' 'PO-19180' 'BB-11545'
'TB-21595' 'RB-19360' 'EB-13705' 'SC-20095' 'TN-21040' 'JS-15940'
'MH-17785' 'JP-15520' 'JE-15475' 'JG-15805' 'XP-21865' 'EM-14065'
'MT-18070' 'SA-20830' 'CW-11905' 'AJ-10960' 'SS-20590' 'RO-19780'
'MD-17350' 'MY-17380' 'CM-12385' 'LS-17245' 'BN-11515' 'DB-13210'
'MC-17605' 'BD-11605' 'MG-18145' 'KB-16240' 'JC-15340' 'RL-19615'
'AA-10375' 'EP-13915' 'DK-12985' 'BD-11500' 'LM-17065' 'AS-10135'
'BD-11320' 'GT-14710' 'AJ-10945' 'OT-18730' 'LP-17080' 'CA-12775'
'JF-15490' 'FP-14320' 'EB-13840' 'JF-15415' 'SF-20200' 'TG-21640'
'WB-21850' 'CC-12145' 'DV-13465' 'BD-11725' 'ZC-21910' 'MS-17830'
'LR-16915' 'TP-21130' 'CK-12205' 'AS-10240' 'AR-10510' 'NB-18655'
'GD-14590' 'CK-12595' 'NG-18355' 'CA-12265' 'SF-20965' 'MO-17800'
'AT-10735' 'FM-14380' 'DJ-13420' 'ME-17725' 'JD-16150' 'JL-15835'
'SC-20305' 'CC-12430' 'AR-10825' 'SR-20740' 'CR-12730' 'EH-14125'
'CB-12025' 'SP-20545' 'TH-21235' 'RP-19390' 'RB-19570' 'CD-11980'
'DJ-13630' 'GT-14635' 'MC-17845' 'RA-19285' 'NP-18325' 'AB-10165'
'JO-15550' 'JK-15370' 'BN-11470' 'DP-13165' 'TH-21550' 'AP-10915']

'RS-19765'	'SV-20365'	'CK-12325'	'RD-19810'	'MR-17545'	'SC-20695'
'JF-15355'	'EG-13900'	'DS-13030'	'PO-19195'	'SS-20875'	'PB-19105'
'RF-19735'	'YC-21895'	'DC-13285'	'CP-12340'	'BF-11020'	'LH-17020'
'CS-12250'	'AJ-10795'	'BV-11245'	'DL-12865'	'BM-11785'	'LT-17110'
'JK-15730'	'ES-14020'	'RH-19495'	'CD-11920'	'HW-14935'	'MC-18130'
'GM-14440'	'PJ-19015'	'BW-11110'	'TR-21325'	'PG-18820'	'JL-15175'
'BM-11650'	'EM-14095'	'AF-10885'	'GA-14725'	'CK-12760'	'DP-13105'
'BK-11260'	'SJ-20125'	'CM-12445'	'AJ-10780'	'LS-16945'	'GP-14740'
'PK-18910'	'SM-20005'	'AG-10765'	'PM-19135'	'LL-16840'	'JS-15595'
'EL-13735'	'PC-18745'	'HL-15040'	'MS-17365'	'GB-14530'	'JR-16210'
'BE-11335'	'SC-20050'	'RW-19630'	'SE-20110'	'AH-10075'	'JM-15535'
'JJ-15760'	'RK-19300'	'CG-12040'	'RP-19270'	'KC-16255'	'KH-16360'
'GH-14665'	'SW-20275'	'JA-15970'	'DL-12925'	'LW-16990'	'TB-21190'
'BS-11800'	'RW-19690'	'TZ-21580'	'AS-10630'	'TS-21340'	'SL-20155'
'MW-18235'	'RD-19585'	'RA-19945'	'MT-17815'	'VG-21790'	'JS-15880'
'KM-16225'	'HR-14770'	'DE-13255'	'AG-10390'	'JJ-15445'	'JH-15430'
'RD-19660'	'MO-17500'	'NS-18640'	'DG-13300'	'NF-18595'	'MG-17650'
'TS-21160'	'BD-11620'	'CM-12160'	'SN-20560'	'EH-14005'	'FO-14305'
'MS-17710'	'CC-12100'	'DW-13540'	'BT-11395'	'CY-12745'	'BT-11485'
'PS-19045'	'PV-18985'	'NM-18520'	'DL-13495'	'CS-12355'	'FH-14275'
'NC-18340'	'AA-10315'	'LT-16765'	'AP-10720'	'PM-18940'	'AT-10435'
'CA-12055'	'HR-14830'	'BT-11530'	'LH-16750'	'SW-20755'	'SP-20620'
'BF-11170'	'KT-16480'	'BG-11695'	'GM-14680'	'EJ-14155'	'NP-18700'
'MH-18115'	'JR-15700'	'SM-20950'	'CC-12220'	'PF-19225'	'DC-12850'
'BD-11770'	'GM-14500'	'TB-21355'	'JH-16180'	'EB-13975'	'QJ-19255'
'TC-21535'	'CS-12460'	'HG-14965'	'LW-16825'	'MC-17575'	'LP-17095'
'EB-14170'	'GZ-14545'	'CP-12085'	'FG-14260'	'LD-17005'	'AB-10255'
'MN-17935'	'JR-15670'	'JF-15190'	'CM-12115'	'AS-10045'	'KB-16315'
'BP-11290'	'ND-18370'	'LB-16735'	'KT-16465'	'HM-14860'	'AB-10600'
'SZ-20035'	'MG-17890'	'JK-16120'	'PP-18955'	'YS-21880'	'KM-16375'
'AB-10105'	'HA-14905'	'BT-11305'	'SV-20815'	'RW-19540'	'DK-12835'
'ST-20530'	'MM-17920'	'PW-19030'	'SC-20440'	'TS-21085'	'MC-17425'
'ME-17320'	'NH-18610'	'MB-18085'	'KD-16495'	'MB-17305'	'KN-16390'
'NP-18685'	'CS-12505'	'KD-16345'	'MS-17770'	'CM-12655'	'Co-12640'
'TS-21370'	'JW-15220'	'JD-15790'	'PC-19000'	'AR-10540'	'AI-10855'
'TB-21400'	'PL-18925'	'GH-14425'	'MP-18175'	'JM-15655'	'CL-11890'
'DB-13270'	'IG-15085'	'BO-11425'	'AB-10150'	'JW-16075'	'EB-13750'
'SG-20470'	'CM-12190'	'AW-10840'	'MC-18100'	'TT-21460'	'VG-21805'
'MY-18295'	'RD-19480'	'DP-13390'	'ML-17395'	'PN-18775'	'JC-15385'
'JG-15160'	'MC-17275'	'NW-18400'	'TB-21280'	'BS-11380'	'FH-14365'
'VM-21685'	'HH-15010'	'CD-12280'	'TH-21100'	'MM-18055'	'NS-18505'
'RB-19645'	'SW-20455'	'EB-13930'	'PS-18760'	'HF-14995'	'HZ-14950'
'CD-12790'	'JK-15205'	'FM-14215'	'ED-13885'	'DA-13450'	'JW-15955'
'RM-19375'	'ML-17755'	'CC-12685'	'JE-15610'	'RP-19855'	'TB-21175'
'BE-11455'	'JF-15565'	'PB-19210'	'BT-11680'	'JL-15235'	'CH-12070'
'ND-18460'	'BF-10975'	'KH-16330'	'GW-14605'	'AC-10420'	'NC-18625'
'ME-18010'	'BP-11230'	'JC-15775'	'AS-10090'	'AC-10450'	'MD-17860'
'DB-13660'	'EH-13990'	'EH-13765'	'MZ-17515'	'SC-20230'	'JE-15715'

'AC-10615'	'JD-16015'	'CB-12415'	'JS-16030'	'LW-17215'	'SC-20800'
'AM-10705'	'RH-19510'	'CT-11995'	'MC-17590'	'CC-12610'	'KA-16525'
'TC-20980'	'BF-11080'	'MM-17260'	'AH-10120'	'BW-11200'	'SW-20245'
'BS-11665'	'RF-19345'	'TB-21625'	'AF-10870'	'RB-19435'	'CS-11950'
'KF-16285'	'JH-15820'	'IL-15100'	'PB-18805'	'RH-19600'	'AW-10930'
'RB-19705'	'ML-17410'	'DB-13555'	'MH-17620'	'DK-13375'	'BT-11440'
'DB-13405'	'TG-21310'	'BF-11005'	'JM-16195'	'MZ-17335'	'MW-18220'
'MV-17485'	'SM-20320'	'TP-21415'	'JK-15625'	'PJ-18835'	'RS-19420'
'SV-20935'	'BC-11125'	'EM-13825'	'BM-11575'	'KN-16705'	'KW-16570'
'SC-20260'	'CV-12295'	'SG-20605'	'TM-21010'	'EM-13810'	'ML-18040'
'CR-12580'	'AZ-10750'	'PW-19240'	'SC-20380'	'CM-11935'	'GM-14695'
'TB-21250'	'JM-15865'	'SC-20575'	'LS-17200'	'RR-19315'	'DB-12910'
'TT-21220'	'LO-17170'	'KD-16615'	'NB-18580'	'BD-11635'	'CM-12235'
'EN-13780'	'KN-16450'	'BO-11350'	'AG-10300'	'MC-17635'	'TA-21385'
'JF-15295'	'TT-21265'	'SB-20170'	'CL-12700'	'HG-15025'	'NL-18310'
'RR-19525'	'TC-21295'	'SV-20785'	'BE-11410'	'SC-20680'	'DF-13135'
'FH-14350'	'MS-17530'	'RH-19555'	'GA-14515'	'JP-16135'	'Dp-13240'
'MO-17950'	'ER-13855'	'DL-13330'	'MH-18025'	'DR-12940'	'DM-13015'
'CA-11965'	'AC-10660'	'DM-13345'	'VF-21715'	'CC-12370'	'BF-11275'
'HG-14845'	'BP-11155'	'EM-14140'	'MA-17995'	'AY-10555'	'GB-14575'
'JB-16045'	'MG-17875'	'SR-20425'	'JB-16000'	'DM-12955'	'TC-21475'
'SW-20350'	'RE-19450'	'BF-11215'	'KB-16405'	'JG-15310'	'EC-14050'
'EB-14110'	'JP-15460'	'CS-11845'	'GH-14410'	'PT-19090'	'JL-15130'
'AH-10030'	'CC-12475'	'DW-13195'	'SJ-20215'	'BG-11740'	'LB-16795'
'CM-11815'	'EH-14185'	'TS-21505'	'PR-18880'	'LC-17050'	'CS-12490'
'DH-13075'	'JO-15145'	'AH-10690'	'HJ-14875'	'MH-17455'	'RD-19930'
'SC-20020'	'SU-20665'	'FC-14335'	'RB-19330'	'NC-18535'	'DB-12970'
'MF-17665'	'RM-19750'	'VM-21835'	'EJ-13720'	'NC-18415'	'LS-17230'
'KE-16420'	'DH-13675'	'PF-19120'	'VT-21700'	'MG-17695'	'SP-20920'
'CS-12175'	'DK-12895'	'KM-16660'	'AA-10645'	'DD-13570'	'AH-10465'
'TP-21565'	'EK-13795'	'CR-12820'	'SG-20890'	'AH-10585'	'NF-18475'
'BS-11365'	'SH-20635'	'RD-19720'	'AG-10330'	'GR-14560'	'VS-21820'
'TM-21490'	'TS-21430'	'DB-13360'	'NR-18550'	'JB-15925'	'CS-11860'
'MF-18250'	'LW-17125'	'AR-10345'	'KS-16300'	'AB-10015'	'LR-17035'
'SS-20410'	'JM-15580'	'JK-15325'	'DM-13525'	'ML-18265'	'MH-17290'
'FC-14245'	'TH-21115'	'JK-16090'	'SB-20185'	'BG-11035'	'BW-11065'
'MG-18205'	'DO-13645'	'BP-11050'	'TS-21655'	'EM-14200'	'AO-10810'
'MH-17440'	'SS-20515'	'LD-16855'	'VP-21760'	'TC-21145'	'IM-15055'
'AR-10570'	'CM-12715'	'FW-14395'	'LC-16960'	'HE-14800'	'BD-11560'
'HD-14785'	'CJ-11875'	'RS-19870'	'SC-20845'	'RE-19405'	'SM-20905']

Customer Name:['Claire Gute' 'Darrin Van Huff' "Sean O'Donnell" 'Brosina Hoffman'

'Andrew Allen' 'Irene Maddox' 'Harold Pawlan' 'Pete Kriz'
'Alejandro Grove' 'Zuschuss Donatelli' 'Ken Black' 'Sandra Flanagan'
'Emily Burns' 'Eric Hoffmann' 'Tracy Blumstein' 'Matt Abelman'
'Gene Hale' 'Steve Nguyen' 'Linda Cazamias' 'Ruben Ausman' 'Erin Smith'
'Odella Nelson' "Patrick O'Donnell" 'Lena Hernandez' 'Darren Powers'

'Janet Molinari' 'Ted Butterfield' 'Kunst Miller' 'Paul Stevenson'
 'Brendan Sweed' 'Karen Daniels' 'Henry MacAllister' 'Joel Eaton'
 'Ken Brennan' 'Stewart Carmichael' 'Duane Noonan' 'Julie Creighton'
 'Christopher Schild' 'Paul Gonzalez' 'Gary Mitchum' 'Jim Sink'
 'Rick Bensley' 'Gary Zandusky' 'Lena Cacioppo' 'Janet Martin'
 'Pete Armstrong' 'Cynthia Voltz' 'Clay Ludtke' 'Ryan Crowe' 'Dave Kipp'
 'Greg Guthrie' 'Steven Cartwright' 'Alan Dominguez' 'Philip Fox'
 'Troy Staebel' 'Lindsay Shagiari' 'Dorothy Wardle' 'Lena Creighton'
 'Jonathan Doherty' 'Sally Hughsby' 'Sandra Glassco' 'Helen Andreada'
 'Maureen Gastineau' 'Justin Ellison' 'Tamara Willingham'
 'Stephanie Phelps' 'Neil Knudson' 'Dave Brooks' 'Nora Paige'
 'Ted Trevino' 'Eric Murdock' 'Ruben Dartt' 'Max Jones' 'Becky Martin'
 'Chad Sievert' 'Jennifer Braxton' 'Shirley Jackson' 'Jim Kriz'
 'David Kendrick' 'Robert Marley' 'Sally Knutson' 'Frank Merwin'
 'Alice McCarthy' 'Mark Packer' 'Bruce Stewart' 'Logan Currie'
 'Heather Kirkland' 'Laurel Elliston' 'Joseph Holt' 'Michael Stewart'
 'Victoria Wilson' 'Jonathan Howell' 'David Smith' 'Valerie Dominguez'
 'Erin Ashbrook' 'David Bremer' 'Ken Lonsdale' 'Dianna Wilson'
 'Logan Haushalter' 'Kelly Collister' 'Delfina Latchford'
 'Dan Reichenbach' 'Craig Carreira' 'Dorris liebe' 'Sean Braxton'
 'Roy Collins' 'Alan Hwang' 'Claudia Bergmann' 'Christine Abelman'
 'Kristen Hastings' 'Barry Blumstein' 'Andrew Gjertsen' 'Jas O'Carroll'
 'Alan Haines' 'Nick Zandusky' 'Kelly Lampkin' 'Alan Schoenberger'
 'Corey Roper' 'Shahid Hopkins' 'Ben Peterman' 'Thomas Seio'
 'Andy Gerbode' 'Sung Pak' 'Nathan Mautz' 'Frank Atkinson' 'Grace Kelly'
 'Don Jones' 'Patrick O'Brill' 'John Lucas' 'Doug Bickford'
 'Clay Cheatham' 'Tamara Dahlen' 'Adam Bellavance' 'Jeremy Lonsdale'
 'Victoria Brennan' 'Katrina Willman' 'Julia Dunbar' 'Michael Kennedy'
 'Guy Thornton' 'Arthur Gainer' 'Muhammed MacIntyre' 'Allen Rosenblatt'
 'Russell Applegate' 'Alejandro Savely' 'Laura Armstrong' 'Denny Ordway'
 'Dean Katz' 'Nathan Gelder' 'Mike Vittorini' 'Jack Garza' 'Bart Pistole'
 'Victor Preis' 'Saphhira Shifley' 'Anna Gayman' 'Luke Foster'
 'Roy Französisch' 'Keith Herrera' 'Kimberly Carter' 'Caroline Jumper'
 'Philip Brown' 'William Brown' 'Michael Paige' 'Natalie Fritzler'
 'Shirley Daniels' 'Ken Heidel' 'Ross Baird' 'Mike Kennedy'
 'Philisse Overcash' 'Brenda Bowman' 'Troy Blackwell' 'Raymond Buch'
 'Ed Braxton' 'Sanjit Chand' 'Tanja Norvell' 'Joni Sundaresam'
 'Maya Herman' 'Jeremy Pistek' 'Jeremy Ellison' 'John Grady'
 'Xylona Preis' 'Erin Mull' 'Michelle Tran' 'Sue Ann Reed' 'Carl Weiss'
 'Astrea Jones' 'Sonia Sunley' 'Rose O'Brian' 'Maribeth Dona'
 'Maribeth Yedwab' 'Christopher Martinez' 'Lynn Smith' 'Bradley Nguyen'
 'Dean Braden' 'Matt Connell' 'Brian Dahlen' 'Mike Gockenbach'
 'Karen Bern' 'Jasper Cacioppo' 'Rob Lucas' 'Allen Arnold' 'Emily Phan'
 'Darren Koutras' 'Bradley Drucker' 'Liz MacKendrick' 'Adrian Shami'
 'Bill Donatelli' 'Greg Tran' 'Ashley Jarboe' 'Olvera Toch'
 'Liz Pelletier' 'Cynthia Arntzen' 'Jeremy Farry' 'Frank Preis'
 'Ellis Ballard' 'Jennifer Ferguson' 'Sarah Foster' 'Trudy Glocke'
 'Carlos Soltero' 'Charles Crestani' 'Dianna Vittorini' 'Bruce Degenhardt'

'Zuschuss Carroll' 'Melanie Seite' 'Lena Radford' 'Theone Pippenger'
 'Chloris Kastensmidt' 'Alan Shonely' 'Andrew Roberts' 'Nona Balk'
 'Giulietta Dortch' 'Clytie Kelty' 'Nat Gilpin' 'Christina Anderson'
 'Sylvia Foulston' 'Meg O'Connel' 'Annie Thurman' 'Fred McMath'
 'Denny Joy' 'Max Engle' 'Justin Deggeller' 'John Lee' 'Sean Christensen'
 'Chuck Clark' 'Anthony Rawles' 'Steven Roelle' 'Craig Reiter'
 'Eugene Hildebrand' 'Cassandra Brandow' 'Sibella Parks' 'Tiffany House'
 'Resi Pölking' 'Rob Beeghly' 'Carol Darley' 'Doug Jacobs'
 'Grant Thornton' 'Michael Chen' 'Ralph Arnett' 'Naresj Patel'
 'Alan Barnes' 'Jesus Ocampo' 'Jay Kimmel' 'Brad Norvell' 'David Philippe'
 'Tracy Hopkins' 'Arthur Prichep' 'Roland Schwarz' 'Seth Vernon'
 'Christine Kargatis' 'Ross DeVincents' 'Mathew Reese' 'Steve Chapman'
 'Jay Fein' 'Emily Grady' 'Darrin Sayre' 'Phillina Ober' 'Sung Shariari'
 'Peter Bühler' 'Roland Fjeld' 'Yoseph Carroll' 'Debra Catini'
 'Christine Phan' 'Barry Französisch' 'Lisa Hazard' 'Chris Selesnick'
 'Anthony Johnson' 'Benjamin Venier' 'Dan Lawera' 'Bryan Mills'
 'Liz Thompson' 'Joe Kamberova' 'Erica Smith' 'Rick Hansen' 'Carlos Daly'
 'Helen Wasserman' 'Mike Caudle' 'Gary McGarr' 'Pauline Johnson'
 'Bart Watters' 'Toby Ritter' 'Patrick Gardner' 'James Lanier'
 'Brian Moss' 'Eudokia Martin' 'Art Foster' 'Guy Armstrong' 'Cyma Kinney'
 'Dave Poirier' 'Berenike Kampe' 'Sanjit Jacobs' 'Chuck Magee'
 'Anthony Jacobs' 'Linda Southworth' 'Guy Phonely' 'Paul Knutson'
 'Sally Matthias' 'Anthony Garverick' 'Peter McVee' 'Lauren Leatherbury'
 'Jill Stevenson' 'Ed Ludwig' 'Pamela Coakley' 'Hunter Lopez'
 'Maribeth Schnelling' 'George Bell' 'Justin Ritter' 'Bill Eplett'
 'Sample Company A' 'Rob Williams' 'Sanjit Engle' 'Adam Hart'
 'Jessica Myrick' 'Joel Jenkins' 'Ralph Kennedy' 'Catherine Glotzbach'
 'Rachel Payne' 'Karen Carlisle' 'Katherine Hughes' 'Greg Hansen'
 'Scott Williamson' 'Joseph Airdo' 'Daniel Lacy' 'Lindsay Williams'
 'Thomas Brumley' 'Bryan Spruell' 'Robert Waldorf' 'Tracy Zic'
 'Ann Steele' 'Toby Swindell' 'Sara Luxemburg' 'Mitch Willingham'
 'Rob Dowd' 'Ryan Akin' 'Meg Tillman' 'Vivek Gonzalez' 'John Stevenson'
 'Kalyca Meade' 'Hallie Redmond' 'Deanra Eno' 'Allen Goldenen'
 'Jennifer Jackson' 'Jennifer Halladay' 'Robert Dilbeck' 'Mary O'Rourke'
 'Noel Staavos' 'Deirdre Greer' 'Nicole Fjeld' 'Matthew Grinstein'
 'Theresa Swint' 'Brian DeCherney' 'Charles McCrossin' 'Skye Norling'
 'Erica Hernandez' 'Frank Olsen' 'Maurice Satty' 'Chad Cunningham'
 'Don Weiss' 'Bill Tyler' 'Craig Yedwab' 'Brad Thomas' 'Penelope Sewall'
 'Paul Van Hugh' 'Neoma Murray' 'Dionis Lloyd' 'Christine Sundaresam'
 'Frank Hawley' 'Nat Carroll' 'Alex Avila' 'Larry Tron' 'Anne Pryor'
 'Paul MacIntyre' 'Alyssa Tate' 'Cathy Armstrong' 'Harold Ryan'
 'Bradley Talbott' 'Larry Hughes' 'Steven Ward' 'Stefania Perrino'
 'Ben Ferrer' 'Kean Thornton' 'Brooke Gillingham' 'Greg Matthias'
 'Eva Jacobs' 'Nora Preis' 'Mick Hernandez' 'Jocasta Rupert'
 'Suzanne McNair' 'Chris Cortes' 'Phillip Flathmann' 'Dan Campbell'
 'Bryan Davis' 'Gene McClure' 'Todd Boyes' 'Justin Hirsh' 'Erica Bern'
 'Quincy Jones' 'Tracy Collins' 'Chuck Sachs' 'Henry Goldwyn'
 'Laurel Workman' 'Matt Collins' 'Liz Preis' 'Evan Bailliet'

'George Zrebassa' 'Cathy Prescott' 'Frank Gastineau' 'Lisa DeCherney'
 'Alejandro Ballentine' 'Michael Nguyen' 'Jim Radford' 'Jamie Frazer'
 'Chad McGuire' 'Aaron Smayling' 'Karl Braun' 'Beth Paige'
 'Natalie DeCherney' 'Larry Blacks' 'Kean Takahito' 'Harry Marie'
 'Ann Blume' 'Sam Zeldin' 'Michael Granlund' 'Julie Kriz' 'Paul Prost'
 'Yana Sorensen' 'Katherine Murray' 'Adrian Barton' 'Helen Abelman'
 'Beth Thompson' 'Stuart Van' 'Rick Wilson' 'Damala Kotsonis' 'Shui Tom'
 'Michael Moore' 'Pauline Webber' 'Shaun Chance' 'Thais Sissman'
 'Mark Cousins' 'Maria Etezadi' 'Nicole Hansen' 'Mick Brown'
 'Keith Dawkins' 'Maria Bertelson' 'Katherine Nockton' 'Nora Pelletier'
 'Cindy Stewart' 'Katherine Ducich' 'Maxwell Schwartz' 'Corinna Mitchell'
 'Corey-Lock' 'Todd Sumrall' 'Jane Waco' 'John Dryer' 'Pauline Chand'
 'Andy Reiter' 'Arianne Irving' 'Tom Boeckenhauer' 'Paul Lucas'
 'Gary Hwang' 'Mike Pelletier' 'Jim Mitchum' 'Carl Ludwig'
 'Deborah Brumfield' 'Ivan Gibson' 'Bobby Odegard' 'Aimee Bixby'
 'Julia West' 'Edward Becker' 'Sheri Gordon' 'Charlotte Melton'
 'Anthony Witt' 'Mick Crebagga' 'Tonja Turnell' 'Vivek Grady'
 'Muhammed Yedwab' 'Rick Duston' 'Dennis Pardue' 'Marina Lichtenstein'
 'Parhena Norris' 'Jenna Caffey' 'James Galang' 'Marc Crier'
 'Natalie Webber' 'Toby Braunhardt' 'Bill Stewart' 'Fred Hopkins'
 'Valerie Mitchum' 'Hilary Holden' 'Christina DeMoss' 'Thea Hendricks'
 'Michelle Moray' 'Neola Schneider' 'Robert Barroso' 'Shaun Weien'
 'Eric Barreto' 'Pamela Stobb' 'Herbert Flentye' 'Henia Zydlo'
 'Cynthia Delaney' 'Jamie Kunitz' 'Filia McAdams' 'Emily Ducich'
 'Dianna Arnett' 'Joni Wasserman' 'Raymond Messe' 'Max Ludwig'
 'Craig Carroll' 'Jim Epp' 'Roy Phan' 'Thomas Boland' 'Brad Eason'
 'Jill Fjeld' 'Phillip Breyer' 'Brian Thompson' 'Janet Lee' 'Cathy Hwang'
 'Neil Ducich' 'Barbara Fisher' 'Katharine Harms' 'Giulietta Weimer'
 'Alyssa Crouse' 'Noah Childs' 'Michelle Ellison' 'Benjamin Patterson'
 'John Castell' 'Adam Shillingsburg' 'Amy Cox' 'Michael Dominguez'
 'Duane Benoit' 'Erica Hackney' 'Edward Hooks' 'Mary Zewe' 'Scot Coram'
 'Joe Elijah' 'Ann Chong' 'Joy Daniels' 'Christy Brittain' 'Joy Smith'
 'Luke Weiss' 'Stuart Calhoun' 'Anne McFarland' 'Rick Huthwaite'
 'Carol Triggs' 'Matt Collister' 'Corey Catlett' 'Kelly Andreada'
 'Tamara Chand' 'Bart Folk' 'Magdelene Morse' 'Adrian Hane' 'Ben Wallace'
 'Scot Wooten' 'Brian Stugart' 'Randy Ferguson' 'Trudy Brown'
 'Art Ferguson' 'Richard Bierner' 'Karen Ferguson' 'John Huston'
 'Ivan Liston' 'Patrick Bzostek' 'Rob Haberlin' 'Arthur Wiediger'
 'Roger Barcio' 'Maris LaWare' 'Dorothy Badders' 'Matt Hagelstein'
 'Dennis Kane' 'Bobby Trafton' 'Denny Blanton' 'Toby Gnade' 'Barry Franz'
 'Justin MacKendrick' 'Maria Zettner' 'Mitch Webber' 'Mark Van Huff'
 'Sean Miller' 'Tom Prescott' 'Jim Karlsson' 'Patrick Jones'
 'Ricardo Sperren' 'Susan Vittorini' 'Becky Castell' 'Elizabeth Moffitt'
 'Brendan Murry' 'Kristina Nunn' 'Kelly Williams' 'Scott Cohen'
 'Christina VanderZanden' 'Speros Goranitis' 'Tamara Manning'
 'Eleni McCrary' 'Michelle Lonsdale' 'Clay Rozendal' 'Annie Zypern'
 'Pierre Wener' 'Shahid Collister' 'Carlos Meador' 'Greg Maxwell'
 'Tim Brockman' 'John Murray' 'Sonia Cooley' 'Luke Schmidt' 'Ralph Ritter'

'Daniel Byrd' 'Thomas Thornton' 'Lori Olson' 'Ken Dana' 'Nicole Brennan'
 'Brian Derr' 'Chris McAfee' 'Edward Nazzal' 'Kean Nguyen' 'Bill Overfelt'
 'Aleksandra Gannaway' 'Matthew Clasen' 'Tom Ashbrook' 'Jason Fortune-'
 'Tim Taslimi' 'Sarah Bern' 'Craig Leslie' 'Hunter Glantz'
 'Nancy Lomonaco' 'Rick Reed' 'Toby Carlisle' 'Stewart Visinsky'
 'Bobby Elias' 'Steve Carroll' 'David Flashing' 'Fred Harton'
 'MaryBeth Skach' 'Ritsa Hightower' 'George Ashbrook' 'Julie Prescott'
 'Dean percer' 'Michael Oakman' 'Elpida Rittenbach' 'Denise Leinenbach'
 'Michelle Huthwaite' 'Daniel Raglin' 'Darrin Martin' 'Carol Adams'
 'Anna Chung' 'Denise Monton' 'Vicky Freymann' 'Christopher Conant'
 'Beth Fritzler' 'Harry Greene' 'Becky Pak' 'Eugene Moren'
 'Michelle Arnett' 'Andy Yotov' 'Giulietta Baptist' 'Julia Barnett'
 'Michael Grace' 'Sharelle Roach' 'Joy Bell-' 'Dario Medina'
 'Tony Chapman' 'Sean Wendt' 'Richard Eichhorn' 'Benjamin Farhat'
 'Katrina Bavinger' 'Jason Gross' 'Erin Creighton' 'Eugene Barchas'
 'Jennifer Patt' 'Cari Sayre' 'Gary Hansen' 'Pete Takahito' 'Jack Lebron'
 'Aaron Hawkins' 'Cindy Chapman' 'David Wiener' 'Sarah Jordon'
 'Bruce Geld' 'Laurel Beltran' 'Candace McMahon' 'Evan Henry' 'Tony Sayre'
 'Patrick Ryan' 'Liz Carlisle' 'Cindy Schnelling' 'Dave Hallsten'
 "Jack O'Briant" 'Anna Häberlin' 'Heather Jas' 'Mark Hamilton'
 "Russell D'Ascenzo" 'Sam Craven' 'Stephanie Ulpright' 'Fred Chung'
 'Randy Bradley' 'Nick Crebassa' 'Darren Budd' 'Maureen Fritzler'
 'Roland Murray' 'Vivian Mathis' 'Ed Jacobs' 'Nathan Cano'
 'Lycoris Saunders' 'Katrina Edelman' 'Duane Huffman' 'Peter Fuller'
 'Valerie Takahito' 'Maureen Gnade' 'Susan Pistek' 'Charles Sheldon'
 'Dana Kaydos' 'Khloe Miller' 'Anna Andreadi' 'Dorothy Dickinson'
 'Amy Hunt' 'Tracy Poddar' 'Eileen Kiefer' 'Cyra Reiten' 'Susan Gilcrest'
 'Angele Hood' 'Neil Französisch' 'Bill Shonely' 'Stefanie Holloman'
 'Roger Demir' 'Alex Grayson' 'Georgia Rosenberg' 'Vivek Sundaresam'
 'Tony Molinari' 'Tom Stivers' 'Dennis Bolton' 'Nick Radford'
 'Joni Blumstein' 'Cari Schnelling' 'Monica Federle' 'Liz Willingham'
 'Alex Russell' 'Karen Seio' 'Aaron Bergman' 'Lisa Ryan' 'Shahid Shariari'
 'Jill Matthias' 'Jason Klamczynski' 'Don Miller' 'Muhammed Lee'
 'Marc Harrigan' 'Frank Carlisle' 'Thea Hudgings' 'Juliana Krohn'
 'Sarah Brown' 'Barry Gonzalez' 'Barry Weirich' 'Mitch Gastineau'
 "Doug O'Connell" 'Barry Pond' 'Trudy Schmidt' 'Evan Minnotte'
 "Anthony O'Donnell" 'Mark Haberlin' 'Shirley Schmidt' 'Lela Donovan'
 'Victoria Pisteka' 'Theresa Coyne' 'Ionia McGrath' 'Anemone Ratner'
 'Craig Molinari' 'Fred Wasserman' 'Lindsay Castell' 'Harold Engle'
 'Brendan Dodson' 'Harold Dahlen' 'Carl Jackson' 'Roy Skaria' 'Sung Chung'
 'Ricardo Emerson' 'Susan MacKendrick']

Segment:['Consumer' 'Corporate' 'Home Office']

Country:['United States']

City:['Henderson' 'Los Angeles' 'Fort Lauderdale' 'Concord' 'Seattle'
 'Fort Worth' 'Madison' 'West Jordan' 'San Francisco' 'Fremont'

'Philadelphia' 'Orem' 'Houston' 'Richardson' 'Naperville' 'Melbourne'
'Eagan' 'Westland' 'Dover' 'New Albany' 'New York City' 'Troy' 'Chicago'
'Gilbert' 'Springfield' 'Jackson' 'Memphis' 'Decatur' 'Durham' 'Columbia'
'Rochester' 'Aurora' 'Charlotte' 'Orland Park' 'Urbandale' 'Columbus'
'Bristol' 'Wilmington' 'Bloomington' 'Phoenix' 'Roseville' 'Independence'
'Pasadena' 'Newark' 'Franklin' 'Scottsdale' 'San Jose' 'Edmond'
'Carlsbad' 'San Antonio' 'Monroe' 'Fairfield' 'Grand Prairie' 'Denver'
'Dallas' 'Whittier' 'Saginaw' 'Medina' 'Detroit' 'Tampa' 'Santa Clara'
'Lakeville' 'San Diego' 'Brentwood' 'Chapel Hill' 'Morristown'
'Cincinnati' 'Inglewood' 'Portland' 'Tamarac' 'Colorado Springs'
'Belleville' 'Taylor' 'Lakewood' 'Arlington' 'Arvada' 'Hackensack'
'Saint Petersburg' 'Long Beach' 'Hesperia' 'Murfreesboro' 'Austin'
'Lowell' 'Manchester' 'Harlingen' 'Tucson' 'Quincy' 'Pembroke Pines'
'Des Moines' 'Peoria' 'Las Vegas' 'Warwick' 'Miami' 'Huntington Beach'
'Richmond' 'Louisville' 'Lawrence' 'Canton' 'New Rochelle' 'Gastonia'
'Jacksonville' 'Auburn' 'Akron' 'Norman' 'Park Ridge' 'Amarillo'
'Lindenhurst' 'Huntsville' 'Fayetteville' 'Costa Mesa' 'Parker' 'Atlanta'
'Gladstone' 'Great Falls' 'Montgomery' 'Mesa' 'Green Bay' 'Anaheim'
'Marysville' 'Salem' 'Laredo' 'Grove City' 'Dearborn' 'Warner Robins'
'Vallejo' 'Minneapolis' 'Mission Viejo' 'Rochester Hills' 'Plainfield'
'Sierra Vista' 'Vancouver' 'Cleveland' 'Tyler' 'Burlington' 'Waynesboro'
'Chester' 'Cary' 'Palm Coast' 'Mount Vernon' 'Hialeah' 'Oceanside'
'Evanston' 'Trenton' 'Cottage Grove' 'Bossier City' 'Lancaster'
'Asheville' 'Lake Elsinore' 'Omaha' 'Edmonds' 'Santa Ana' 'Milwaukee'
'Florence' 'Lorain' 'Linden' 'Salinas' 'New Brunswick' 'Garland'
'Norwich' 'Alexandria' 'Toledo' 'Farmington' 'Riverside' 'Torrance'
'Round Rock' 'Boca Raton' 'Virginia Beach' 'Murrieta' 'Olympia'
'Washington' 'Jefferson City' 'Saint Peters' 'Rockford' 'Brownsville'
'Yonkers' 'Oakland' 'Clinton' 'Encinitas' 'Roswell' 'Jonesboro' 'Antioch'
'Homestead' 'La Porte' 'Lansing' 'Cuyahoga Falls' 'Reno' 'Harrisonburg'
'Escondido' 'Royal Oak' 'Rockville' 'Coral Springs' 'Buffalo'
'Boynton Beach' 'Gulfport' 'Fresno' 'Greenville' 'Macon' 'Cedar Rapids'
'Providence' 'Pueblo' 'Saint Paul' 'Deltona' 'Murray' 'Middletown'
'Freeport' 'Pico Rivera' 'Provo' 'Pleasant Grove' 'Smyrna' 'Parma'
'Mobile' 'New Bedford' 'Irving' 'Vineland' 'Glendale' 'Niagara Falls'
'Thomasville' 'Westminster' 'Coppell' 'Pomona' 'North Las Vegas'
'Allentown' 'Tempe' 'Laguna Niguel' 'Bridgeton' 'Everett' 'Watertown'
'Appleton' 'Bellevue' 'Allen' 'El Paso' 'Grapevine' 'Carrollton' 'Kent'
'Lafayette' 'Tigard' 'Skokie' 'Plano' 'Suffolk' 'Indianapolis' 'Bayonne'
'Dublin' 'Greensboro' 'Baltimore' 'Kenosha' 'Olathe' 'Tulsa' 'Redmond'
'Raleigh' 'Muskogee' 'Meriden' 'Bowling Green' 'South Bend' 'Spokane'
'Keller' 'Port Orange' 'Medford' 'Charlottesville' 'Missoula' 'Apopka'
'Reading' 'Broomfield' 'Paterson' 'Oklahoma City' 'Chesapeake' 'Lubbock'
'Johnson City' 'San Bernardino' 'Leominster' 'Bozeman' 'Perth Amboy'
'Ontario' 'Rancho Cucamonga' 'Moorhead' 'Mesquite' 'Redlands' 'Stockton'
'Ormond Beach' 'Sunnyvale' 'York' 'College Station' 'Saint Louis'
'Manteca' 'San Angelo' 'Salt Lake City' 'Knoxville' 'Little Rock'
'Lincoln Park' 'Marion' 'Littleton' 'Bangor' 'Southaven' 'New Castle'

'Midland' 'Sioux Falls' 'Fort Collins' 'Clarksville' 'Sacramento'
 'Thousand Oaks' 'Malden' 'Holyoke' 'Albuquerque' 'Sparks' 'Coachella'
 'Elmhurst' 'Passaic' 'North Charleston' 'Newport News' 'Jamestown'
 'Mishawaka' 'Westfield' 'La Quinta' 'Tallahassee' 'Nashville'
 'Bellingham' 'Woodstock' 'Haltom City' 'Wheeling' 'Summerville'
 'Hot Springs' 'Englewood' 'Las Cruces' 'Hoover' 'Frisco' 'Vacaville'
 'Waukesha' 'Bakersfield' 'Pompano Beach' 'Corpus Christi' 'Redondo Beach'
 'Orlando' 'Orange' 'Lake Charles' 'Highland Park' 'Hempstead'
 'Noblesville' 'Apple Valley' 'Mount Pleasant' 'Sterling Heights'
 'Eau Claire' 'Pharr' 'Billings' 'Gresham' 'Chattanooga' 'Meridian'
 'Bolingbrook' 'Lakeland' 'Maple Grove' 'Woodland' 'Missouri City'
 'Pearland' 'San Mateo' 'Grand Rapids' 'Visalia' 'Overland Park'
 'Temecula' 'Yucaipa' 'Revere' 'Conroe' 'Tinley Park' 'Dubuque'
 'Dearborn Heights' 'Santa Fe' 'Hickory' 'Carol Stream' 'Saint Cloud'
 'North Miami' 'Plantation' 'Port Saint Lucie' 'Rock Hill' 'Odessa'
 'West Allis' 'Chula Vista' 'Manhattan' 'Altoona' 'Thornton' 'Champaign'
 'Texarkana' 'Edinburg' 'Baytown' 'Greenwood' 'Woonsocket' 'Superior'
 'Bedford' 'Covington' 'Broken Arrow' 'Miramar' 'Hollywood' 'Deer Park'
 'Wichita' 'McAllen' 'Iowa City' 'Boise' 'Cranston' 'Port Arthur'
 'Citrus Heights' 'The Colony' 'Daytona Beach' 'Bullhead City' 'Portage'
 'Fargo' 'Elkhart' 'San Gabriel' 'Hamilton' 'Margate' 'Sandy Springs'
 'Mentor' 'Lawton' 'Hampton' 'Rome' 'La Crosse' 'Lewiston' 'Hattiesburg'
 'Danville' 'Logan' 'Waterbury' 'Athens' 'Avondale' 'Marietta' 'Yuma'
 'Wausau' 'Pasco' 'Oak Park' 'Pensacola' 'League City' 'Gaithersburg'
 'Lehi' 'Tuscaloosa' 'Moreno Valley' 'Georgetown' 'Loveland' 'Chandler'
 'Helena' 'Kirkwood' 'Waco' 'Frankfort' 'Bethlehem' 'Grand Island'
 'Woodbury' 'Rogers' 'Clovis' 'Jupiter' 'Santa Barbara' 'Cedar Hill'
 'Norfolk' 'Draper' 'Ann Arbor' 'La Mesa' 'Pocatello' 'Holland' 'Milford'
 'Buffalo Grove' 'Lake Forest' 'Redding' 'Chico' 'Utica' 'Conway'
 'Cheyenne' 'Owensboro' 'Caldwell' 'Kenner' 'Nashua' 'Bartlett'
 'Redwood City' 'Lebanon' 'Santa Maria' 'Des Plaines' 'Longview'
 'Hendersonville' 'Waterloo' 'Cambridge' 'Palatine' 'Beverly' 'Eugene'
 'Oxnard' 'Renton' 'Glenview' 'Delray Beach' 'Commerce City' 'Texas City'
 'Wilson' 'Rio Rancho' 'Goldsboro' 'Montebello' 'El Cajon'
 'West Palm Beach' 'Normal' 'Saint Charles' 'Camarillo' 'Hillsboro'
 'Burbank' 'Modesto' 'Garden City' 'Atlantic City' 'Longmont' 'Davis'
 'Morgan Hill' 'Clifton' 'Sheboygan' 'East Point' 'Rapid City' 'Andover'
 'Kissimmee' 'Shelton' 'Danbury' 'Sanford' 'San Marcos' 'Greeley'
 'Mansfield' 'Elyria' 'Twin Falls' 'Coral Gables' 'Romeoville'
 'Marlborough' 'Laurel' 'Bryan' 'Pine Bluff' 'Aberdeen' 'Hagerstown'
 'East Orange' 'Arlington Heights' 'Oswego' 'Beaumont' 'Coon Rapids'
 'San Clemente' 'San Luis Obispo' 'Springdale' 'Lodi' 'Mason']

State:['Kentucky' 'California' 'Florida' 'North Carolina' 'Washington' 'Texas'
 'Wisconsin' 'Utah' 'Nebraska' 'Pennsylvania' 'Illinois' 'Minnesota'
 'Michigan' 'Delaware' 'Indiana' 'New York' 'Arizona' 'Virginia'
 'Tennessee' 'Alabama' 'South Carolina' 'Colorado' 'Iowa' 'Ohio'
 'Missouri' 'Oklahoma' 'New Mexico' 'Louisiana' 'Connecticut' 'New Jersey']

'Oregon' 'Massachusetts' 'Georgia' 'Nevada' 'Rhode Island' 'Mississippi'
'Arkansas' 'Montana' 'New Hampshire' 'Maryland' 'District of Columbia'
'Kansas' 'Vermont' 'Maine' 'South Dakota' 'Idaho' 'North Dakota'
'Wyoming' 'West Virginia']

Postal Code:[42420. 90036. 33311. 90032. 28027. 98103. 76106. 53711. 84084.
94109.

68025. 19140. 84057. 90049. 77095. 75080. 77041. 60540. 32935. 55122.
48185. 19901. 47150. 10024. 12180. 90004. 60610. 85234. 22153. 10009.
49201. 38109. 77070. 35601. 94122. 27707. 60623. 29203. 55901. 80013.
28205. 60462. 10035. 50322. 43229. 37620. 19805. 61701. 85023. 95661.
64055. 91104. 43055. 53132. 85254. 95123. 98105. 98115. 73034. 90045.
19134. 88220. 78207. 77036. 62521. 71203. 6824. 75051. 80219. 75220.
37064. 90604. 48601. 44256. 48227. 38401. 33614. 95051. 55044. 92037.
77506. 94513. 27514. 7960. 45231. 94110. 90301. 97206. 33319. 80906.
7109. 48180. 8701. 22204. 80004. 7601. 33710. 19143. 90805. 92345.
37130. 78745. 1852. 31907. 6040. 78550. 85705. 62301. 2038. 33024.
98198. 61604. 89115. 2886. 33180. 28403. 92646. 40475. 80027. 1841.
39212. 48187. 10801. 28052. 32216. 47201. 13021. 44312. 73071. 94521.
60068. 79109. 11757. 90008. 92024. 77340. 14609. 72701. 92627. 80134.
30318. 64118. 59405. 48234. 36116. 85204. 60653. 54302. 45503. 92804.
98270. 97301. 78041. 19120. 75217. 43123. 10011. 48126. 31088. 94591.
55407. 92691. 48307. 7060. 85635. 98661. 60505. 76017. 40214. 75081.
44105. 75701. 27217. 22980. 19013. 27511. 32137. 10550. 48205. 33012.
11572. 92105. 60201. 48183. 55016. 71111. 50315. 93534. 23223. 28806.
92530. 68104. 98026. 92704. 53209. 41042. 44052. 7036. 93905. 8901.
17602. 3301. 21044. 75043. 6360. 22304. 43615. 87401. 92503. 90503.
78664. 92054. 33433. 23464. 92563. 28540. 52601. 98502. 20016. 65109.
63376. 61107. 33142. 78521. 10701. 94601. 28110. 20735. 30076. 72401.
47374. 94509. 33030. 46350. 48911. 44221. 89502. 22801. 92025. 48073.
20852. 33065. 14215. 33437. 39503. 93727. 27834. 11561. 35630. 31204.
52402. 2908. 81001. 94533. 55106. 32725. 42071. 6457. 11520. 90660.
84604. 84062. 30080. 24153. 44134. 36608. 2740. 75061. 8360. 85301.
14304. 27360. 92683. 38301. 75019. 91767. 89031. 18103. 19711. 85281.
92677. 8302. 2149. 13601. 54915. 98006. 75002. 79907. 76051. 75007.
37167. 98031. 70506. 97224. 60076. 75023. 23434. 46203. 7002. 43017.
28314. 27405. 21215. 53142. 66062. 98002. 74133. 97756. 27604. 74403.
6450. 42104. 46614. 6010. 89015. 99207. 76248. 45014. 32127. 97504.
22901. 59801. 33178. 29501. 97477. 32712. 19601. 80020. 65807. 7501.
73120. 23320. 79424. 65203. 37604. 36830. 92404. 1453. 59715. 85345.
44107. 8861. 91761. 91730. 56560. 75150. 92374. 95207. 32174. 94086.
3820. 17403. 77840. 63116. 2169. 95336. 44240. 76903. 84106. 35810.
37918. 72209. 48146. 43302. 80122. 5408. 4401. 38671. 47362. 48640.
57103. 80525. 47905. 37042. 95823. 91360. 2148. 1040. 87105. 89431.
92236. 60126. 7055. 29406. 23602. 14701. 46544. 43402. 7090. 92253.
32303. 37211. 98226. 60098. 76117. 60090. 29483. 71901. 80112. 43130.
88001. 35244. 75034. 95687. 84107. 53186. 93309. 33068. 45373. 78415.
90278. 32839. 7050. 70601. 60035. 11550. 46060. 55124. 29464. 48310.

54703. 78577. 59102. 97030. 37421. 83642. 92307. 60440. 33801. 55369.
 95695. 77489. 77581. 94403. 49505. 93277. 66212. 92592. 92399. 2151.
 77301. 60477. 52001. 48127. 87505. 28601. 60188. 56301. 33161. 46226.
 33317. 34952. 29730. 79762. 53214. 91911. 66502. 16602. 80229. 61821.
 47401. 71854. 78539. 77520. 46142. 90712. 2895. 54880. 76021. 98042.
 74012. 33023. 33021. 77536. 67212. 78501. 52240. 83704. 2920. 61032.
 77642. 95610. 75056. 98052. 32114. 86442. 46368. 58103. 46514. 91776.
 45011. 33063. 30328. 44060. 73505. 23666. 13440. 54601. 83501. 39401.
 94526. 48858. 84321. 6708. 30605. 4240. 61832. 85323. 30062. 85364.
 54401. 99301. 60302. 32503. 77573. 20877. 84043. 35401. 92553. 40324.
 80538. 85224. 59601. 63122. 76706. 48066. 60423. 18018. 55113. 68801.
 55125. 48237. 72756. 88101. 33458. 93101. 75104. 68701. 84020. 48104.
 91941. 83201. 49423. 6460. 60089. 92630. 96003. 95928. 13501. 72032.
 82001. 42301. 83605. 70065. 3060. 38134. 94061. 37087. 93454. 60016.
 98632. 37075. 50701. 2138. 60067. 1915. 97405. 93030. 98059. 60025.
 33445. 80022. 77590. 27893. 87124. 27534. 98208. 90640. 92020. 33407.
 61761. 60174. 93010. 97123. 91505. 95351. 67846. 8401. 80501. 95616.
 26003. 95037. 7011. 53081. 30344. 57701. 1810. 34741. 6484. 6810.
 52302. 32771. 78666. 80634. 76063. 44035. 83301. 63301. 33134. 60441.
 1752. 20707. 77803. 71603. 57401. 21740. 7017. 60004. 60543. 77705.
 55433. 92672. 94568. 93405. 72762. 95240. 77571. 45040. 30188.]

Region:['South' 'West' 'Central' 'East']

Product ID:['FUR-B0-10001798' 'FUR-CH-10000454' 'OFF-LA-10000240' ...
 71.16199999999998 14.4648 12.672]

Category:['Furniture' 'Office Supplies' 'Technology']

Sub-Category:['Bookcases' 'Chairs' 'Labels' 'Tables' 'Storage' 'Furnishings'
 'Art'
 'Phones' 'Binders' 'Appliances' 'Paper' 'Accessories' 'Envelopes'
 'Fasteners' 'Supplies' 'Machines' 'Copiers']

Product Name:['Bush Somerset Collection Bookcase'
 'Hon Deluxe Fabric Upholstered Stacking Chairs, Rounded Back'
 'Self-Adhesive Address Labels for Typewriters by Universal' ...
 'Hoover Commercial Lightweight Upright Vacuum' 'LG G2'
 'Eldon Jumbo ProFile Portable File Boxes Graphite/Black']

Sales:[261.96 731.94 14.62 ... 8.78 376.74 44.43]

Quantity:[2. 3. 5. 7. 4. 6. 9. 1. 8. 14. 11. 13. 10. 12.]

Discount:[0. 0.45 0.2 0.8 0.3 0.5 0.7 0.6 0.32 0.1 0.4 0.15]

Profit:[41.9136 219.582 6.8714 ... 2.99 130.2885 18.6606]

```

order_year:[2016. 2015. 2014. 2017.]

order_month:[11.  6. 10.  4. 12.  5.  8.  7.  9.  1.  3.  2.]

order_date:[ 8. 12. 11.  9. 15.  5. 22. 13. 27. 16. 25. 17. 19. 10. 20. 18. 24.
30.
 4. 14. 26.  3. 28. 29.  1. 23.  2.  6.  7. 31. 21.]

Ship_year:[2016. 2015. 2014. 2017. 2018.]

Ship_month:[11.  6. 10.  4. 12.  5.  9.  7.  1.  3.  8.  2.]

Ship_date:[11. 16. 18. 14. 20. 10. 26. 15.  1. 13. 30. 21. 23. 31. 22. 25. 17.
5.
 6.  2.  8.  4. 28. 27. 12. 19.  7.  9.  3. 24. 29.]

```

```
[69]: df.columns #prints all columns
```

```
[69]: Index(['Row ID', 'Order ID', 'Order Date', 'Ship Date', 'Ship Mode',
'Customer ID', 'Customer Name', 'Segment', 'Country', 'City', 'State',
'Postal Code', 'Region', 'Product ID', 'Category', 'Sub-Category',
'Product Name', 'Sales', 'Quantity', 'Discount', 'Profit', 'order_year',
'order_month', 'order_date', 'Ship_year', 'Ship_month', 'Ship_date'],
dtype='object')
```

```
[70]: #dropping columns that are not required for analysis
df.drop(columns=['Row ID', 'Order ID',
'Customer ID', 'Customer Name',
'Postal Code', 'Product ID', 'Order Date', 'Ship_
Date'],axis=1,inplace=True)
```

```
[71]: df.sample(5) #display random 5 samples
```

```
[71]:
```

	Ship Mode	Segment	Country	City \
1938	Standard Class	Consumer	United States	Oklahoma City
4419	First Class	Home Office	United States	Aurora
7769	Standard Class	Consumer	United States	Morgan Hill
6763	Standard Class	Corporate	United States	Springfield
1296	First Class	Corporate	United States	Jacksonville

	State	Region	Category	Sub-Category \
1938	Oklahoma	Central	Technology	Phones
4419	Colorado	West	Office Supplies	Paper
7769	California	West	Office Supplies	Art
6763	Virginia	South	Technology	Accessories
1296	North Carolina	South	Technology	Machines

	Product Name	Sales	Quantity \
1938	Jabra Supreme Plus Driver Edition Headset	479.960	4.0
4419	Xerox 1992	14.352	3.0
7769	Rogers Handheld Barrel Pencil Sharpener	21.920	8.0
6763	Imation 30456 USB Flash Drive 8GB	20.700	3.0
1296	Cisco CP-7937G Unified IP Conference Station P...	695.700	2.0

	Discount	Profit	order_year	order_month	order_date	Ship_year \
1938	0.0	134.3888	2014.0	12.0	3.0	2014.0
4419	0.2	5.2026	2017.0	9.0	28.0	2017.0
7769	0.0	5.9184	2016.0	12.0	3.0	2016.0
6763	0.0	1.6560	2016.0	9.0	12.0	2016.0
1296	0.5	-27.8280	2017.0	1.0	2.0	2017.0

	Ship_month	Ship_date
1938	12.0	9.0
4419	10.0	1.0
7769	12.0	7.0
6763	9.0	17.0
1296	1.0	4.0

```
[72]: #data descriptive
df.describe()
```

```
[72]:
```

	Sales	Quantity	Discount	Profit	order_year \
count	9844.000000	9844.000000	9844.000000	9844.000000	9844.000000
mean	230.386939	3.790431	0.155312	28.970794	2015.719829
std	626.024933	2.224033	0.205817	235.713664	1.123452
min	0.444000	1.000000	0.000000	-6599.978000	2014.000000
25%	17.310000	2.000000	0.000000	1.757325	2015.000000
50%	54.804000	3.000000	0.200000	8.709500	2016.000000
75%	209.970000	5.000000	0.200000	29.460650	2017.000000
max	22638.480000	14.000000	0.800000	8399.976000	2017.000000

	order_month	order_date	Ship_year	Ship_month	Ship_date
count	9844.000000	9844.000000	9844.000000	9844.000000	9844.000000
mean	7.811560	15.460585	2015.735270	7.742686	15.854429
std	3.284929	8.756378	1.126109	3.342078	8.807584
min	1.000000	1.000000	2014.000000	1.000000	1.000000
25%	5.000000	8.000000	2015.000000	5.000000	8.000000
50%	9.000000	15.000000	2016.000000	9.000000	16.000000
75%	11.000000	23.000000	2017.000000	11.000000	24.000000
max	12.000000	31.000000	2018.000000	12.000000	31.000000

```
[73]: # df.to_excel('C:/Users/ganesh/Downloads/cleanedata1.xlsx', index=False)
```

3 Outlier Analysis

4 Quantity Outlier Analysis

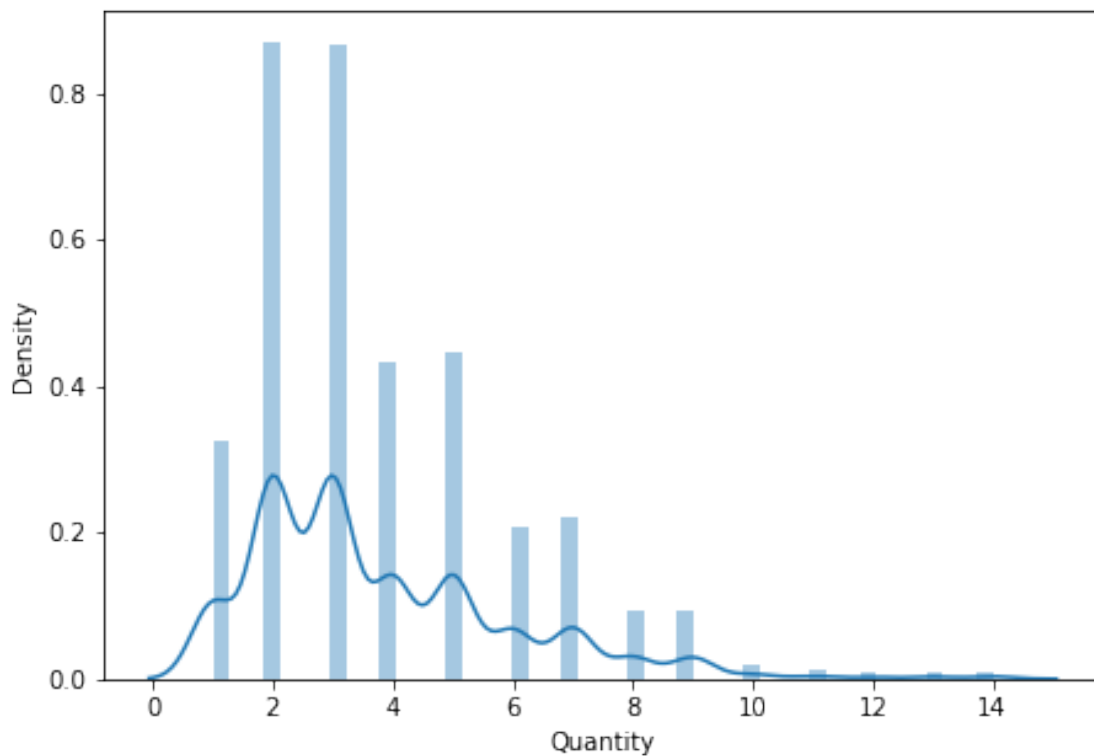
```
[74]: #unique values  
df["Quantity"].unique()
```

```
[74]: array([ 2.,  3.,  5.,  7.,  4.,  6.,  9.,  1.,  8., 14., 11., 13., 10.,  
        12.])
```

```
[75]: #distribution graph for quantity  
plt.figure(figsize=(16,5))  
plt.subplot(1,2,2)  
sns.distplot(df["Quantity"])  
plt.show()
```

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn\distributions.py:2619:
FutureWarning: `distplot` is a deprecated function and will be removed in a
future version. Please adapt your code to use either `displot` (a figure-level
function with similar flexibility) or `histplot` (an axes-level function for
histograms).

```
warnings.warn(msg, FutureWarning)
```



```
[76]: #right skew  
df["Quantity"].skew()
```

```
[76]: 1.272389553355905
```

```
[77]: df["Quantity"].describe()
```

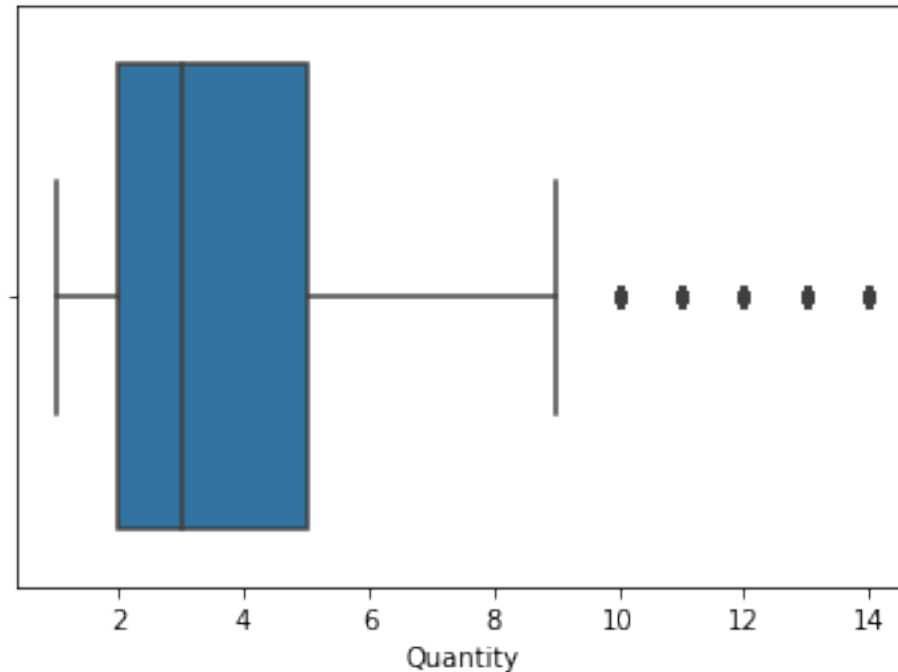
```
[77]: count      9844.000000  
      mean         3.790431  
      std         2.224033  
      min         1.000000  
      25%         2.000000  
      50%         3.000000  
      75%         5.000000  
      max        14.000000  
      Name: Quantity, dtype: float64
```

```
[78]: #box plot  
sns.boxplot(df["Quantity"])
```

```
C:\Users\ganesh\anaconda3\lib\site-packages\seaborn\_decorators.py:36:  
FutureWarning: Pass the following variable as a keyword arg: x. From version  
0.12, the only valid positional argument will be `data`, and passing other  
arguments without an explicit keyword will result in an error or  
misinterpretation.  
    warnings.warn(  

```

```
[78]: <AxesSubplot:xlabel='Quantity'>
```

```
[79]: #Percentile Method
low=df["Quantity"].quantile(0.05)
high=df["Quantity"].quantile(0.95)
low,high
```

```
[79]: (1.0, 8.0)
```

```
[80]: df[df["Quantity"]>high]
```

```
[80]:
```

	Ship Mode	Segment	Country	City	State \
10	Standard Class	Consumer	United States	Los Angeles	California
37	Standard Class	Home Office	United States	Houston	Texas
102	Second Class	Consumer	United States	Columbus	Ohio
111	First Class	Consumer	United States	Wilmington	Delaware
124	Standard Class	Consumer	United States	Roseville	California
...
9801	Standard Class	Consumer	United States	San Francisco	California
9802	Standard Class	Consumer	United States	Anaheim	California
9839	Standard Class	Home Office	United States	Los Angeles	California
9844	Standard Class	Consumer	United States	Long Beach	New York
9860	Second Class	Consumer	United States	Atlanta	Georgia

	Region	Category	Sub-Category \
10	West	Furniture	Tables
37	Central	Office Supplies	Envelopes

102	East	Office Supplies	Fasteners
111	East	Office Supplies	Envelopes
124	West	Furniture	Furnishings
...
9801	West	Technology	Accessories
9802	West	Office Supplies	Storage
9839	West	Office Supplies	Binders
9844	East	Office Supplies	Labels
9860	South	Office Supplies	Paper

		Product Name	Sales	Quantity \
10		Chromcraft Rectangular Conference Tables	1706.184	9.0
37	#10-4 1/8" x 9 1/2"	Premium Diagonal Seam Enve...	113.328	9.0
102		OIC Colored Binder Clips, Assorted Sizes	40.096	14.0
111	Globe Weis Peel & Seel	First Class Envelopes	115.020	9.0
124		Longer-Life Soft White Bulbs	43.120	14.0
...
9801	Memorex Mini Travel Drive	16 GB USB 2.0 Flash ...	223.580	14.0
9802	Carina Mini System Audio Rack, Model AR050B		998.820	9.0
9839	Ibico Recycled Linen-Style Covers		437.472	14.0
9844	Self-Adhesive Removable Labels		31.500	10.0
9860	Wirebound Message Book, 4 per Page		48.870	9.0

	Discount	Profit	order_year	order_month	order_date	Ship_year \
10	0.2	85.3092	2014.0	6.0	9.0	2014.0
37	0.2	35.4150	2015.0	12.0	27.0	2015.0
102	0.2	14.5348	2014.0	8.0	25.0	2014.0
111	0.0	51.7590	2016.0	6.0	12.0	2016.0
124	0.0	20.6976	2016.0	10.0	13.0	2016.0
...
9801	0.0	87.1962	2017.0	11.0	24.0	2017.0
9802	0.0	29.9646	2014.0	12.0	28.0	2015.0
9839	0.2	153.1152	2016.0	12.0	6.0	2016.0
9844	0.0	15.1200	2015.0	5.0	17.0	2015.0
9860	0.0	23.9463	2017.0	11.0	25.0	2017.0

	Ship_month	Ship_date
10	6.0	14.0
37	12.0	31.0
102	8.0	27.0
111	6.0	15.0
124	10.0	19.0
...
9801	11.0	30.0
9802	1.0	3.0
9839	12.0	10.0
9844	5.0	23.0

9860 11.0 29.0

[418 rows x 19 columns]

```
[81]: #did because felt outliers are valid (so used capping)
```

```
[82]: #capping outliers
df["Quantity"]=np.where(df["Quantity"]>high,high,np.
↪where(df["Quantity"]<low,low,df["Quantity"]))
```

```
[83]: df.shape
```

```
[83]: (9844, 19)
```

```
[84]: df[df["Quantity"]>=5]
```

```
[84]:
```

	Ship Mode	Segment	Country	City	State \
3	Standard Class	Consumer	United States	Fort Lauderdale	Florida
5	Standard Class	Consumer	United States	Los Angeles	California
7	Standard Class	Consumer	United States	Los Angeles	California
9	Standard Class	Consumer	United States	Los Angeles	California
10	Standard Class	Consumer	United States	Los Angeles	California
...
9860	Second Class	Consumer	United States	Atlanta	Georgia
9861	Standard Class	Consumer	United States	Los Angeles	California
9864	Standard Class	Consumer	United States	Detroit	Michigan
9869	Standard Class	Consumer	United States	Detroit	Michigan
9870	Second Class	Corporate	United States	Fort Lauderdale	Florida

	Region	Category	Sub-Category \
3	South	Furniture	Tables
5	West	Furniture	Furnishings
7	West	Technology	Phones
9	West	Office Supplies	Appliances
10	West	Furniture	Tables
...
9860	South	Office Supplies	Paper
9861	West	Office Supplies	Paper
9864	Central	Office Supplies	Binders
9869	Central	Office Supplies	Binders
9870	South	Office Supplies	Appliances

	Product Name	Sales	Quantity \
3	Bretford CR4500 Series Slim Rectangular Table	957.5775	5.0
5	Eldon Expressions Wood and Plastic Desk Access...	48.8600	7.0
7	Mitel 5320 IP Phone VoIP phone	907.1520	6.0
9	Belkin F5C206VTEL 6 Outlet Surge	114.9000	5.0

10	Chromcraft Rectangular Conference Tables	1706.1840	8.0
...
9860	Wirebound Message Book, 4 per Page	48.8700	8.0
9861	Xerox 19	154.9000	5.0
9864	Wilson Jones Turn Tabs Binder Tool for Ring Bi...	24.1000	5.0
9869	Wilson Jones 1" Hanging DublLock Ring Binders	26.4000	5.0
9870	Hoover Upright Vacuum With Dirt Cup	1158.1200	5.0

	Discount	Profit	order_year	order_month	order_date	Ship_year	\
3	0.45	-383.0310	2015.0	10.0	11.0	2015.0	
5	0.00	14.1694	2014.0	6.0	9.0	2014.0	
7	0.20	90.7152	2014.0	6.0	9.0	2014.0	
9	0.00	34.4700	2014.0	6.0	9.0	2014.0	
10	0.20	85.3092	2014.0	6.0	9.0	2014.0	
...	
9860	0.00	23.9463	2017.0	11.0	25.0	2017.0	
9861	0.00	69.7050	2017.0	1.0	14.0	2017.0	
9864	0.00	11.0860	2016.0	9.0	1.0	2016.0	
9869	0.00	12.6720	2016.0	9.0	1.0	2016.0	
9870	0.20	130.2885	2017.0	11.0	11.0	2017.0	

	Ship_month	Ship_date
3	10.0	18.0
5	6.0	14.0
7	6.0	14.0
9	6.0	14.0
10	6.0	14.0
...
9860	11.0	29.0
9861	1.0	19.0
9864	9.0	5.0
9869	9.0	5.0
9870	11.0	16.0

[3053 rows x 19 columns]

```
[85]: #distribution and box plot for quantity
plt.figure(figsize=(16,8))
plt.subplot(2,2,1)
sns.distplot(df["Quantity"])

plt.subplot(2,2,2)
sns.boxplot(df["Quantity"])

plt.show()
```

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn\distributions.py:2619:

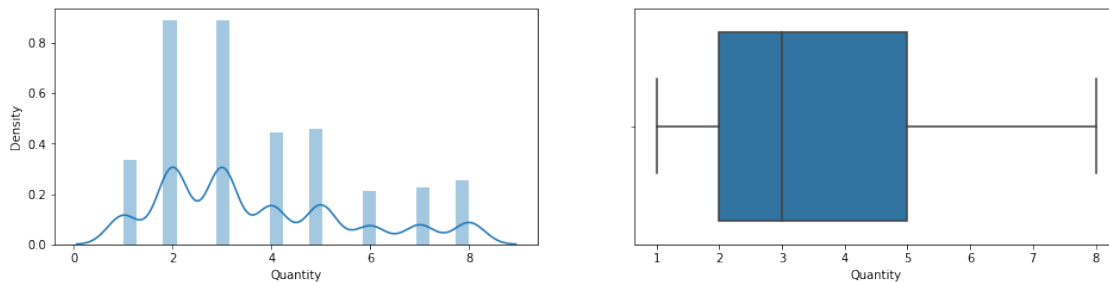
FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

```
warnings.warn(msg, FutureWarning)
```

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn_decorators.py:36:

FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
warnings.warn(
```



```
[86]: #right skew
df["Quantity"].skew()
```

```
[86]: 0.7107622703947052
```

5 Disount Outlier Analysis

```
[87]: #unique values in discount
df["Discount"].unique()
```

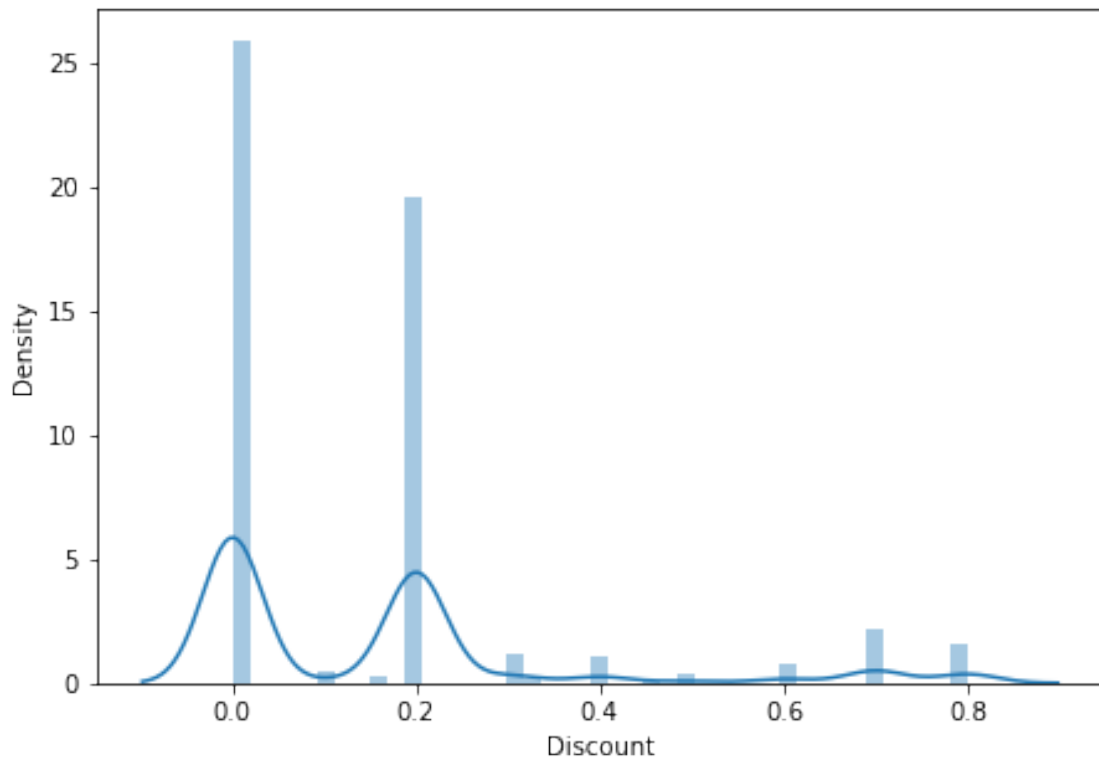
```
[87]: array([0. , 0.45, 0.2 , 0.8 , 0.3 , 0.5 , 0.7 , 0.6 , 0.32, 0.1 , 0.4 ,
        0.15])
```

```
[88]: #distribution graph for discount
plt.figure(figsize=(16,5))
plt.subplot(1,2,2)
sns.distplot(df["Discount"])
plt.show()
```

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn\distributions.py:2619:

FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for

```
histograms).  
warnings.warn(msg, FutureWarning)
```



```
[89]: #right skew  
df["Discount"].skew()
```

```
[89]: 1.6922053394888577
```

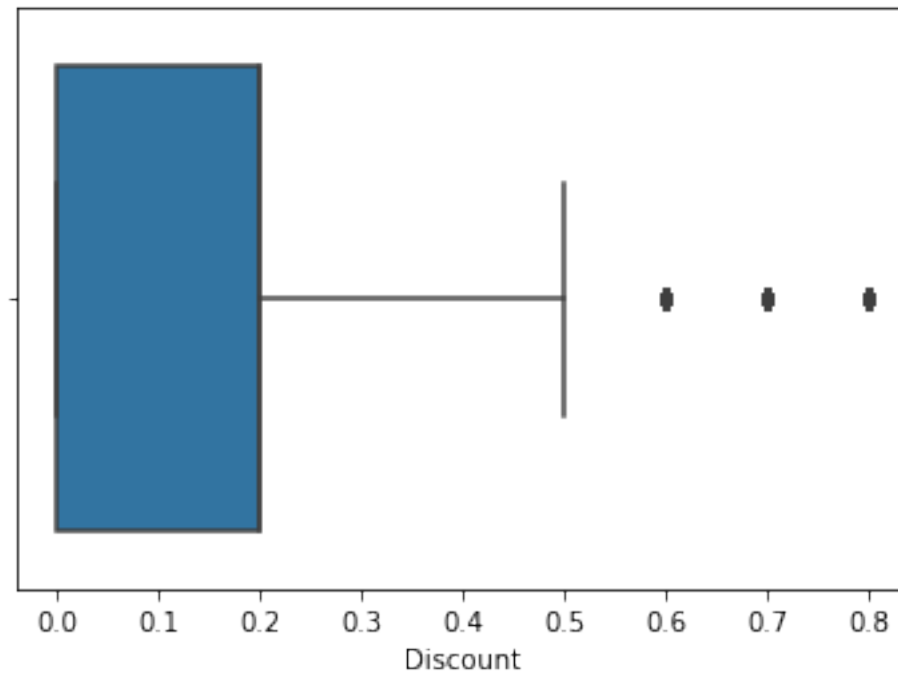
```
[90]: df["Discount"].describe()
```

```
[90]: count      9844.000000  
mean         0.155312  
std          0.205817  
min          0.000000  
25%          0.000000  
50%          0.200000  
75%          0.200000  
max          0.800000  
Name: Discount, dtype: float64
```

```
[91]: #boxplot  
sns.boxplot(df["Discount"])
```

```
C:\Users\ganesh\anaconda3\lib\site-packages\seaborn\_decorators.py:36:
FutureWarning: Pass the following variable as a keyword arg: x. From version
0.12, the only valid positional argument will be `data`, and passing other
arguments without an explicit keyword will result in an error or
misinterpretation.
  warnings.warn(
```

```
[91]: <AxesSubplot:xlabel='Discount'>
```



```
[92]: per25=df["Discount"].quantile(0.25)
      per75=df["Discount"].quantile(0.75)
      per25,per75
```

```
[92]: (0.0, 0.2)
```

```
[93]: #IQR(method)
      iqr=per75-per25
      low=per25-1.5*iqr
      high=per75+1.5*iqr
      low,high
```

```
[93]: (-0.30000000000000004, 0.5)
```

```
[94]: df[df["Discount"]>high]
```

[94]:

	Ship Mode	Segment	Country	City	State \
14	Standard Class	Home Office	United States	Fort Worth	Texas
15	Standard Class	Home Office	United States	Fort Worth	Texas
28	Standard Class	Consumer	United States	Philadelphia	Pennsylvania
32	Standard Class	Consumer	United States	Philadelphia	Pennsylvania
36	First Class	Corporate	United States	Richardson	Texas
...
9737	First Class	Home Office	United States	Cleveland	Ohio
9763	Standard Class	Consumer	United States	Carrollton	Texas
9780	Standard Class	Corporate	United States	Bryan	Texas
9781	Standard Class	Home Office	United States	Akron	Ohio
9855	Standard Class	Consumer	United States	Phoenix	Arizona

	Region	Category	Sub-Category \
14	Central	Office Supplies	Appliances
15	Central	Office Supplies	Binders
28	East	Office Supplies	Binders
32	East	Office Supplies	Binders
36	Central	Furniture	Furnishings
...
9737	East	Office Supplies	Binders
9763	Central	Furniture	Furnishings
9780	Central	Office Supplies	Binders
9781	East	Office Supplies	Binders
9855	West	Office Supplies	Binders

	Product Name	Sales	Quantity \
14	Holmes Replacement Filter for HEPA Air Cleaner...	68.810	5.0
15	Storex DuraTech Recycled Plastic Frosted Binders	2.544	3.0
28	Avery Recycled Flexi-View Covers for Binding S...	9.618	2.0
32	Acco Pressboard Covers with Storage Hooks, 14 ...	6.858	6.0
36	Electrix Architect's Clamp-On Swing Arm Lamp, ...	190.920	5.0
...
9737	Wilson Jones Clip & Carry Folder Binder Tool f...	8.700	5.0
9763	GE General Use Halogen Bulbs, 100 Watts, 1 Bul...	25.128	3.0
9780	GBC Pre-Punched Binding Paper, Plastic, White,...	22.386	7.0
9781	Acco Expandable Hanging Binders	5.742	3.0
9855	Plastic Binding Combs	18.180	4.0

	Discount	Profit	order_year	order_month	order_date	Ship_year \
14	0.8	-123.8580	2015.0	11.0	22.0	2015.0
15	0.8	-3.8160	2015.0	11.0	22.0	2015.0
28	0.7	-7.0532	2015.0	9.0	17.0	2015.0
32	0.7	-5.7150	2015.0	9.0	17.0	2015.0
36	0.6	-147.9630	2016.0	12.0	8.0	2016.0
...
9737	0.7	-6.3800	2017.0	4.0	20.0	2017.0

9763	0.6	-6.9102	2014.0	11.0	12.0	2014.0
9780	0.8	-35.8176	2016.0	3.0	15.0	2016.0
9781	0.7	-4.5936	2014.0	11.0	24.0	2014.0
9855	0.7	-13.9380	2017.0	9.0	19.0	2017.0

	Ship_month	Ship_date
14	11.0	26.0
15	11.0	26.0
28	9.0	21.0
32	9.0	21.0
36	12.0	10.0
...
9737	4.0	21.0
9763	11.0	18.0
9780	3.0	19.0
9781	11.0	30.0
9855	9.0	25.0

[834 rows x 19 columns]

```
[95]: #capping outliers with high and low values
df["Discount"] = np.where(df["Discount"] > high, high, np.
    ↳ where(df["Discount"] < low, low, df["Discount"]))
```

```
[96]: df.shape
```

```
[96]: (9844, 19)
```

```
[97]: #distribution and box plot for discount
plt.figure(figsize=(16,8))
plt.subplot(2,2,1)
sns.distplot(df["Discount"])
plt.subplot(2,2,2)
sns.boxplot(df["Discount"])

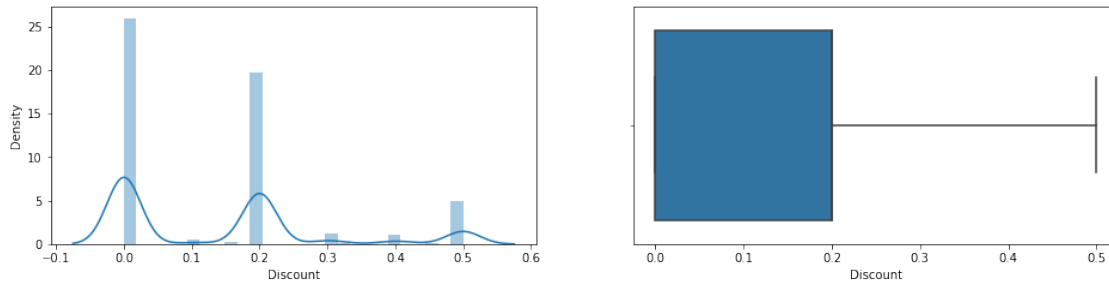
plt.show()
```

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn\distributions.py:2619:
FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn_decorators.py:36:
FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or

```
misinterpretation.
warnings.warn(
```



```
[98]: df.shape
```

```
[98]: (9844, 19)
```

```
[99]: #right skew
df["Discount"].skew()
```

```
[99]: 0.9644235759205968
```

6 Profit Outlier Analysis

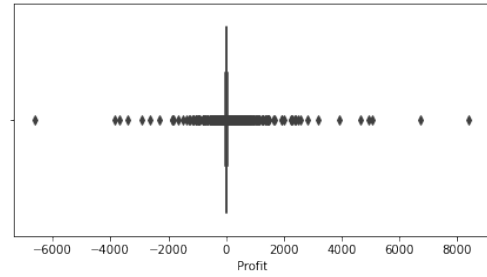
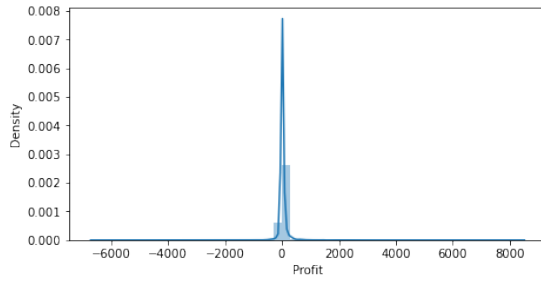
```
[100]: #distribution and boxplot graph for profit
plt.figure(figsize=(16,8))
plt.subplot(2,2,1)
sns.distplot(df["Profit"])
plt.subplot(2,2,2)
sns.boxplot(df["Profit"])
plt.show()
```

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn\distributions.py:2619:
FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

```
warnings.warn(msg, FutureWarning)
```

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn_decorators.py:36:
FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
warnings.warn(
```



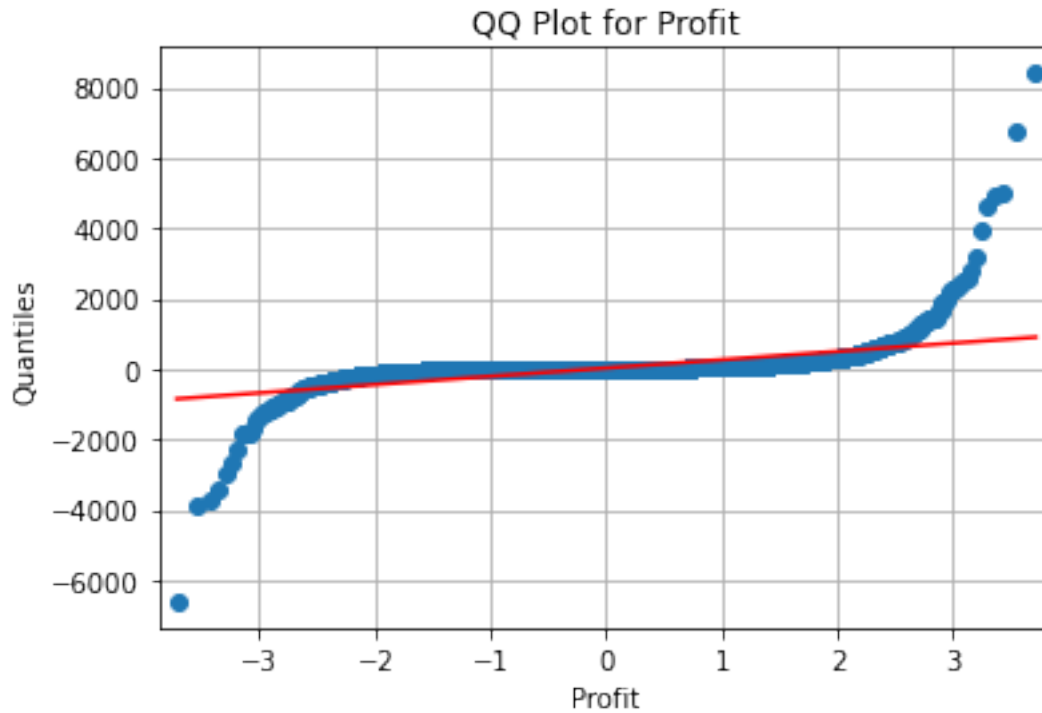
```
[101]: #right skew
df["Profit"].skew()
```

```
[101]: 7.532178434791842
```

```
[102]: data = df

# Select the column you want to analyze
column_name = 'Profit'
column_data = data[column_name]

# Create the QQ plot using statsmodels.api
sm.qqplot(column_data, line='s') # 's' for straight reference line
plt.xlabel(column_name)
plt.ylabel('Quantiles')
plt.title('QQ Plot for ' + column_name)
plt.grid(True)
plt.show()
```



```
[103]: df[df["Profit"]<=0]
```

```
[103]:
```

	Ship Mode	Segment	Country	City \
3	Standard Class	Consumer	United States	Fort Lauderdale
14	Standard Class	Home Office	United States	Fort Worth
15	Standard Class	Home Office	United States	Fort Worth
23	Second Class	Consumer	United States	Philadelphia
27	Standard Class	Consumer	United States	Philadelphia
...
9791	Standard Class	Consumer	United States	San Bernardino
9797	Second Class	Corporate	United States	Los Angeles
9822	First Class	Home Office	United States	Houston
9837	Standard Class	Home Office	United States	Los Angeles
9855	Standard Class	Consumer	United States	Phoenix

	State	Region	Category	Sub-Category \
3	Florida	South	Furniture	Tables
14	Texas	Central	Office Supplies	Appliances
15	Texas	Central	Office Supplies	Binders
23	Pennsylvania	East	Furniture	Chairs
27	Pennsylvania	East	Furniture	Bookcases
...
9791	California	West	Furniture	Bookcases

9797	California	West	Furniture	Tables
9822	Texas	Central	Furniture	Bookcases
9837	California	West	Office Supplies	Fasteners
9855	Arizona	West	Office Supplies	Binders

	Product Name	Sales	Quantity \
3	Bretford CR4500 Series Slim Rectangular Table	957.5775	5.0
14	Holmes Replacement Filter for HEPA Air Cleaner...	68.8100	5.0
15	Storex DuraTech Recycled Plastic Frosted Binders	2.5440	3.0
23	Global Deluxe Stacking Chair, Gray	71.3720	2.0
27	Riverside Palais Royal Lawyers Bookcase, Royal...	3083.4300	7.0
...
9791	O'Sullivan Living Dimensions 3-Shelf Bookcases	683.3320	4.0
9797	Hon 61000 Series Interactive Training Tables	71.0880	2.0
9822	Bush Heritage Pine Collection 5-Shelf Bookcase...	383.4656	4.0
9837	Alliance Big Bands Rubber Bands, 12/Pack	13.8600	7.0
9855	Plastic Binding Combs	18.1800	4.0

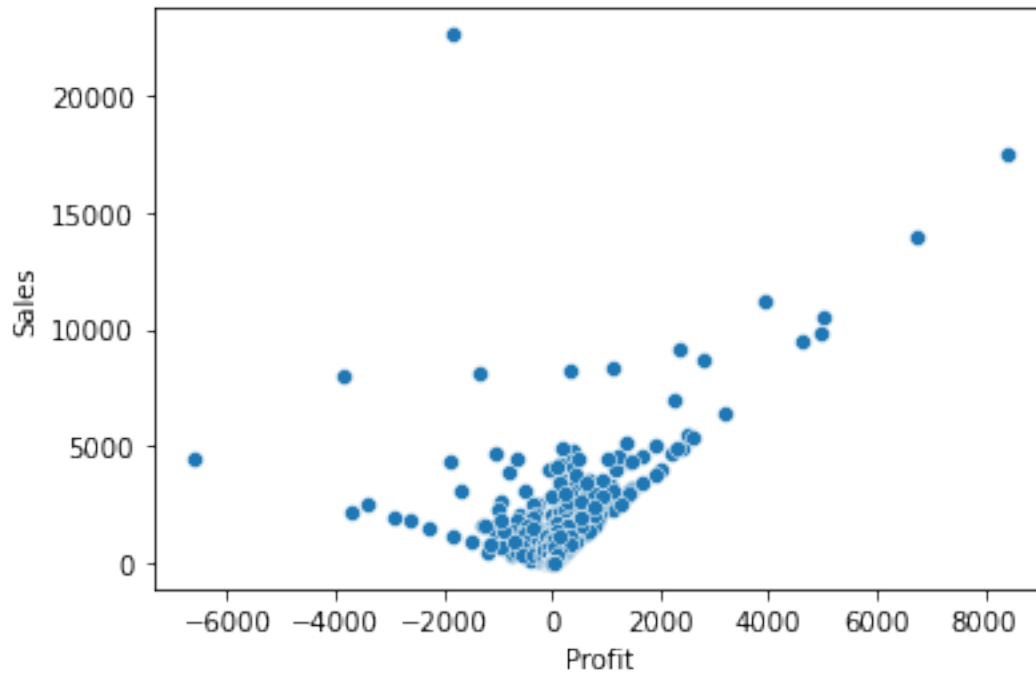
	Discount	Profit	order_year	order_month	order_date	Ship_year \
3	0.45	-383.0310	2015.0	10.0	11.0	2015.0
14	0.50	-123.8580	2015.0	11.0	22.0	2015.0
15	0.50	-3.8160	2015.0	11.0	22.0	2015.0
23	0.30	-1.0196	2017.0	7.0	16.0	2017.0
27	0.50	-1665.0522	2015.0	9.0	17.0	2015.0
...
9791	0.15	-40.1960	2015.0	11.0	13.0	2015.0
9797	0.20	-1.7772	2016.0	6.0	3.0	2016.0
9822	0.32	-67.6704	2015.0	3.0	19.0	2015.0
9837	0.00	0.0000	2016.0	12.0	6.0	2016.0
9855	0.50	-13.9380	2017.0	9.0	19.0	2017.0

	Ship_month	Ship_date
3	10.0	18.0
14	11.0	26.0
15	11.0	26.0
23	7.0	18.0
27	9.0	21.0
...
9791	11.0	17.0
9797	6.0	6.0
9822	3.0	22.0
9837	12.0	10.0
9855	9.0	25.0

[1896 rows x 19 columns]

```
[104]: #scatterplot prfit vs sales
sns.scatterplot(x="Profit",y="Sales",data=df)
```

```
[104]: <AxesSubplot:xlabel='Profit', ylabel='Sales'>
```



```
[105]: df[df["Profit"]<=0]
```

```
[105]:
```

	Ship Mode	Segment	Country	City \
3	Standard Class	Consumer	United States	Fort Lauderdale
14	Standard Class	Home Office	United States	Fort Worth
15	Standard Class	Home Office	United States	Fort Worth
23	Second Class	Consumer	United States	Philadelphia
27	Standard Class	Consumer	United States	Philadelphia
...
9791	Standard Class	Consumer	United States	San Bernardino
9797	Second Class	Corporate	United States	Los Angeles
9822	First Class	Home Office	United States	Houston
9837	Standard Class	Home Office	United States	Los Angeles
9855	Standard Class	Consumer	United States	Phoenix

	State	Region	Category	Sub-Category \
3	Florida	South	Furniture	Tables
14	Texas	Central	Office Supplies	Appliances
15	Texas	Central	Office Supplies	Binders
23	Pennsylvania	East	Furniture	Chairs

27	Pennsylvania	East	Furniture	Bookcases
...
9791	California	West	Furniture	Bookcases
9797	California	West	Furniture	Tables
9822	Texas	Central	Furniture	Bookcases
9837	California	West	Office Supplies	Fasteners
9855	Arizona	West	Office Supplies	Binders

		Product Name	Sales	Quantity \
3		Bretford CR4500 Series Slim Rectangular Table	957.5775	5.0
14		Holmes Replacement Filter for HEPA Air Cleaner...	68.8100	5.0
15		Storex DuraTech Recycled Plastic Frosted Binders	2.5440	3.0
23		Global Deluxe Stacking Chair, Gray	71.3720	2.0
27		Riverside Palais Royal Lawyers Bookcase, Royal...	3083.4300	7.0
...	
9791		O'Sullivan Living Dimensions 3-Shelf Bookcases	683.3320	4.0
9797		Hon 61000 Series Interactive Training Tables	71.0880	2.0
9822		Bush Heritage Pine Collection 5-Shelf Bookcase...	383.4656	4.0
9837		Alliance Big Bands Rubber Bands, 12/Pack	13.8600	7.0
9855		Plastic Binding Combs	18.1800	4.0

	Discount	Profit	order_year	order_month	order_date	Ship_year \
3	0.45	-383.0310	2015.0	10.0	11.0	2015.0
14	0.50	-123.8580	2015.0	11.0	22.0	2015.0
15	0.50	-3.8160	2015.0	11.0	22.0	2015.0
23	0.30	-1.0196	2017.0	7.0	16.0	2017.0
27	0.50	-1665.0522	2015.0	9.0	17.0	2015.0
...
9791	0.15	-40.1960	2015.0	11.0	13.0	2015.0
9797	0.20	-1.7772	2016.0	6.0	3.0	2016.0
9822	0.32	-67.6704	2015.0	3.0	19.0	2015.0
9837	0.00	0.0000	2016.0	12.0	6.0	2016.0
9855	0.50	-13.9380	2017.0	9.0	19.0	2017.0

	Ship_month	Ship_date
3	10.0	18.0
14	11.0	26.0
15	11.0	26.0
23	7.0	18.0
27	9.0	21.0
...
9791	11.0	17.0
9797	6.0	6.0
9822	3.0	22.0
9837	12.0	10.0
9855	9.0	25.0

[1896 rows x 19 columns]

```
[106]: EPSILON = 1e-8
positive_profit = df['Profit'] + EPSILON

# Take the absolute value
abs_profit = np.abs(positive_profit)

# Apply the logarithm
log_profit = np.log(abs_profit)

# Restore the sign
transformed_profit = np.sign(df['Profit']) * log_profit

# Replace -0 with 0
transformed_profit = np.where(transformed_profit == -0, 0, transformed_profit)

# Add the transformed profit column to the DataFrame
df['Transformed_Profit'] = transformed_profit
```

```
[107]: df[df["Profit"]==0]
```

```
[107]:
```

	Ship Mode	Segment	Country	City	State \
201	Standard Class	Home Office	United States	Tampa	Florida
509	Second Class	Consumer	United States	San Francisco	California
520	First Class	Consumer	United States	Seattle	Washington
526	Standard Class	Corporate	United States	Seattle	Washington
775	Standard Class	Consumer	United States	Philadelphia	Pennsylvania
...
9272	First Class	Consumer	United States	Los Angeles	California
9501	Standard Class	Corporate	United States	Seattle	Washington
9746	Standard Class	Consumer	United States	Lafayette	Indiana
9758	Standard Class	Consumer	United States	Fairfield	Ohio
9837	Standard Class	Home Office	United States	Los Angeles	California

	Region	Category	Sub-Category \
201	South	Furniture	Furnishings
509	West	Furniture	Chairs
520	West	Office Supplies	Fasteners
526	West	Furniture	Chairs
775	East	Furniture	Chairs
...
9272	West	Furniture	Chairs
9501	West	Office Supplies	Storage
9746	Central	Office Supplies	Fasteners
9758	East	Furniture	Furnishings
9837	West	Office Supplies	Fasteners

	Product Name	Sales	Quantity \
201	Tenex Contemporary Contur Chairmats for Low an...	258.072	3.0
509	HON 5400 Series Task Chairs for Big and Tall	1121.568	2.0
520	Alliance Big Bands Rubber Bands, 12/Pack	3.960	2.0
526	Hon Every-Day Series Multi-Task Chairs	451.152	3.0
775	Global Leather Executive Chair	1228.465	5.0
...
9272	HON 5400 Series Task Chairs for Big and Tall	2803.920	5.0
9501	Contico 72"H Heavy-Duty Storage System	204.900	5.0
9746	Alliance Big Bands Rubber Bands, 12/Pack	5.940	3.0
9758	Deflect-o EconoMat Nonstudded, No Bevel Mat	82.640	2.0
9837	Alliance Big Bands Rubber Bands, 12/Pack	13.860	7.0

	Discount	Profit	order_year	order_month	order_date	Ship_year \
201	0.2	0.0	2017.0	4.0	7.0	2017.0
509	0.2	0.0	2016.0	4.0	15.0	2016.0
520	0.0	0.0	2015.0	12.0	7.0	2015.0
526	0.2	0.0	2017.0	10.0	1.0	2017.0
775	0.3	0.0	2014.0	6.0	28.0	2014.0
...
9272	0.2	0.0	2015.0	1.0	27.0	2015.0
9501	0.0	0.0	2014.0	3.0	7.0	2014.0
9746	0.0	0.0	2014.0	1.0	23.0	2014.0
9758	0.2	0.0	2016.0	6.0	6.0	2016.0
9837	0.0	0.0	2016.0	12.0	6.0	2016.0

	Ship_month	Ship_date	Transformed_Profit
201	4.0	12.0	0.0
509	4.0	17.0	0.0
520	12.0	9.0	0.0
526	10.0	8.0	0.0
775	7.0	2.0	0.0
...
9272	1.0	29.0	0.0
9501	3.0	12.0	0.0
9746	1.0	27.0	0.0
9758	6.0	10.0	0.0
9837	12.0	10.0	0.0

[65 rows x 20 columns]

```
[108]: # sum of null values
df["Transformed_Profit"].isnull().sum()
```

[108]: 0

```
[109]: df[df["Transformed_Profit"]<=0]
```

```
[109]:
```

	Ship Mode	Segment	Country	City \
3	Standard Class	Consumer	United States	Fort Lauderdale
14	Standard Class	Home Office	United States	Fort Worth
15	Standard Class	Home Office	United States	Fort Worth
23	Second Class	Consumer	United States	Philadelphia
27	Standard Class	Consumer	United States	Philadelphia
...
9797	Second Class	Corporate	United States	Los Angeles
9804	Second Class	Home Office	United States	Seattle
9822	First Class	Home Office	United States	Houston
9837	Standard Class	Home Office	United States	Los Angeles
9855	Standard Class	Consumer	United States	Phoenix

	State	Region	Category	Sub-Category \
3	Florida	South	Furniture	Tables
14	Texas	Central	Office Supplies	Appliances
15	Texas	Central	Office Supplies	Binders
23	Pennsylvania	East	Furniture	Chairs
27	Pennsylvania	East	Furniture	Bookcases
...
9797	California	West	Furniture	Tables
9804	Washington	West	Office Supplies	Storage
9822	Texas	Central	Furniture	Bookcases
9837	California	West	Office Supplies	Fasteners
9855	Arizona	West	Office Supplies	Binders

	Product Name	Sales	Quantity \
3	Bretford CR4500 Series Slim Rectangular Table	957.5775	5.0
14	Holmes Replacement Filter for HEPA Air Cleaner...	68.8100	5.0
15	Storex DuraTech Recycled Plastic Frosted Binders	2.5440	3.0
23	Global Deluxe Stacking Chair, Gray	71.3720	2.0
27	Riverside Palais Royal Lawyers Bookcase, Royal...	3083.4300	7.0
...
9797	Hon 61000 Series Interactive Training Tables	71.0880	2.0
9804	Rogers Jumbo File, Granite	40.7400	3.0
9822	Bush Heritage Pine Collection 5-Shelf Bookcase...	383.4656	4.0
9837	Alliance Big Bands Rubber Bands, 12/Pack	13.8600	7.0
9855	Plastic Binding Combs	18.1800	4.0

	Discount	Profit	order_year	order_month	order_date	Ship_year \
3	0.45	-383.0310	2015.0	10.0	11.0	2015.0
14	0.50	-123.8580	2015.0	11.0	22.0	2015.0
15	0.50	-3.8160	2015.0	11.0	22.0	2015.0
23	0.30	-1.0196	2017.0	7.0	16.0	2017.0
27	0.50	-1665.0522	2015.0	9.0	17.0	2015.0

...
9797	0.20	-1.7772	2016.0	6.0	3.0	2016.0
9804	0.00	0.4074	2015.0	4.0	12.0	2015.0
9822	0.32	-67.6704	2015.0	3.0	19.0	2015.0
9837	0.00	0.0000	2016.0	12.0	6.0	2016.0
9855	0.50	-13.9380	2017.0	9.0	19.0	2017.0

	Ship_month	Ship_date	Transformed_Profit
3	10.0	18.0	-5.948116
14	11.0	26.0	-4.819136
15	11.0	26.0	-1.339203
23	7.0	18.0	-0.019410
27	9.0	21.0	-7.417612
...
9797	6.0	6.0	-0.575039
9804	4.0	17.0	-0.897960
9822	3.0	22.0	-4.214649
9837	12.0	10.0	0.000000
9855	9.0	25.0	-2.634619

[2123 rows x 20 columns]

```
[110]: df.head() #top 5 rows
```

```
[110]:
```

	Ship Mode	Segment	Country	City	State \
0	Second Class	Consumer	United States	Henderson	Kentucky
1	Second Class	Consumer	United States	Henderson	Kentucky
2	Second Class	Corporate	United States	Los Angeles	California
3	Standard Class	Consumer	United States	Fort Lauderdale	Florida
4	Standard Class	Consumer	United States	Fort Lauderdale	Florida

	Region	Category	Sub-Category \
0	South	Furniture	Bookcases
1	South	Furniture	Chairs
2	West	Office Supplies	Labels
3	South	Furniture	Tables
4	South	Office Supplies	Storage

	Product Name	Sales	Quantity \
0	Bush Somerset Collection Bookcase	261.9600	2.0
1	Hon Deluxe Fabric Upholstered Stacking Chairs,...	731.9400	3.0
2	Self-Adhesive Address Labels for Typewriters b...	14.6200	2.0
3	Bretford CR4500 Series Slim Rectangular Table	957.5775	5.0
4	Eldon Fold 'N Roll Cart System	22.3680	2.0

	Discount	Profit	order_year	order_month	order_date	Ship_year \
0	0.00	41.9136	2016.0	11.0	8.0	2016.0

1	0.00	219.5820	2016.0	11.0	8.0	2016.0
2	0.00	6.8714	2016.0	6.0	12.0	2016.0
3	0.45	-383.0310	2015.0	10.0	11.0	2015.0
4	0.20	2.5164	2015.0	10.0	11.0	2015.0

	Ship_month	Ship_date	Transformed_Profit
0	11.0	11.0	3.735610
1	11.0	11.0	5.391726
2	6.0	16.0	1.927368
3	10.0	18.0	-5.948116
4	10.0	18.0	0.922829

```
[111]: #distribution and boxplot for transformed profit
plt.figure(figsize=(16,8))
plt.subplot(2,2,1)
sns.distplot(df["Transformed_Profit"])
plt.subplot(2,2,2)
sns.boxplot(df["Transformed_Profit"])

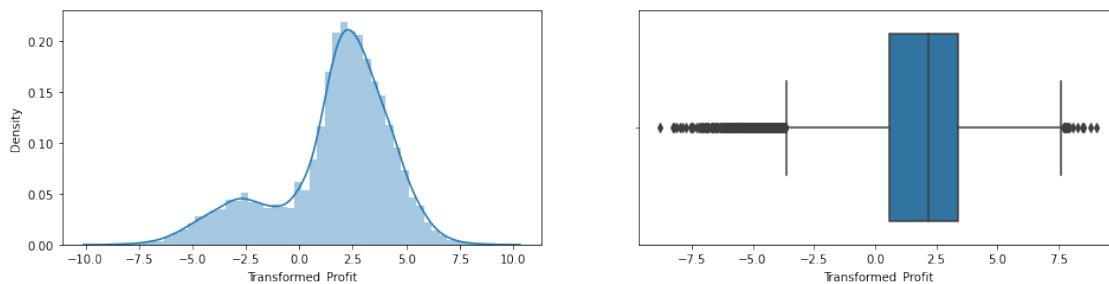
plt.show()
```

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn\distributions.py:2619:
FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

```
warnings.warn(msg, FutureWarning)
```

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn_decorators.py:36:
FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
warnings.warn(
```



```
[112]: #left skew
df["Transformed_Profit"].skew()
```

```
[112]: -0.8874539960016707
```

```
[113]: df.shape
```

```
[113]: (9844, 20)
```

```
[114]: Q1 = df['Transformed_Profit'].quantile(0.25)
Q3 = df['Transformed_Profit'].quantile(0.75)

# Calculate IQR(Inter Quantile Range)
IQR = Q3 - Q1

# Define lower and upper bounds
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Remove outliers
df_no_outliers = df[(df['Transformed_Profit'] >= lower_bound) &
                    (df['Transformed_Profit'] <= upper_bound)]
```

```
[115]: df.shape
```

```
[115]: (9844, 20)
```

```
[116]: #copying data in df
df=df_no_outliers.copy()
```

```
[117]: df.shape
```

```
[117]: (9233, 20)
```

```
[118]: #sum of null values in transformed profit
df['Transformed_Profit'].isnull().sum()
```

```
[118]: 0
```

```
[119]: df.shape
```

```
[119]: (9233, 20)
```

```
[120]: #distribution and box plot for transformed profit
plt.figure(figsize=(16,8))
plt.subplot(2,2,1)
sns.distplot(df["Transformed_Profit"])
plt.subplot(2,2,2)
sns.boxplot(df["Transformed_Profit"])
```

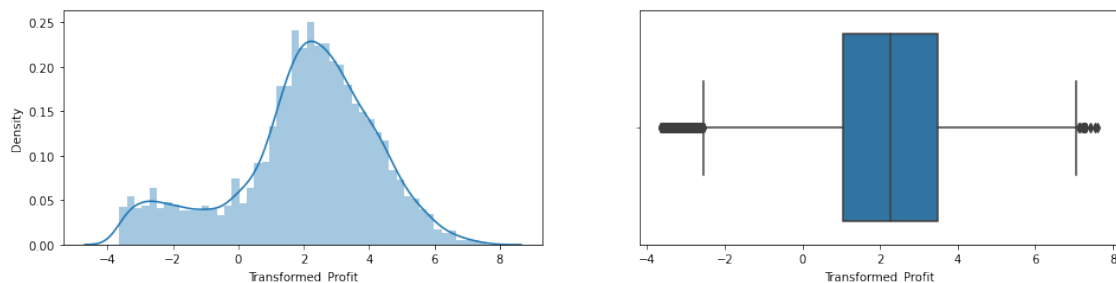
```
plt.show()
```

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn\distributions.py:2619:
FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

```
warnings.warn(msg, FutureWarning)
```

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn_decorators.py:36:
FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
warnings.warn(
```



```
[121]: df["Transformed_Profit"].skew()
```

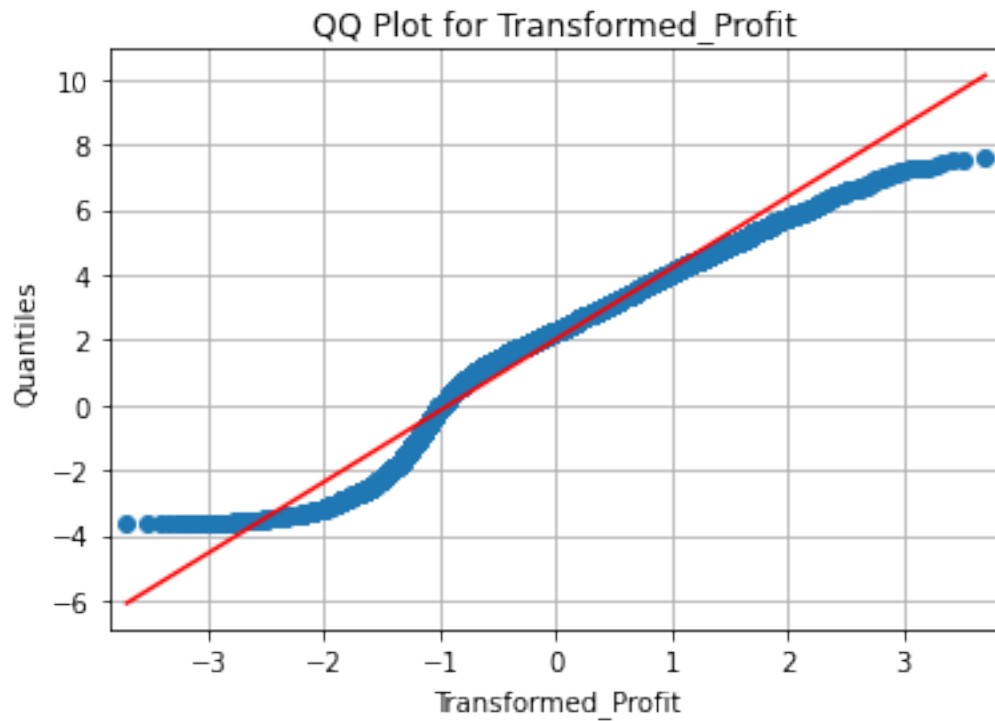
```
[121]: -0.6187364044245487
```

```
[122]: #display all columns  
df.columns
```

```
[122]: Index(['Ship Mode', 'Segment', 'Country', 'City', 'State', 'Region',  
        'Category', 'Sub-Category', 'Product Name', 'Sales', 'Quantity',  
        'Discount', 'Profit', 'order_year', 'order_month', 'order_date',  
        'Ship_year', 'Ship_month', 'Ship_date', 'Transformed_Profit'],  
        dtype='object')
```

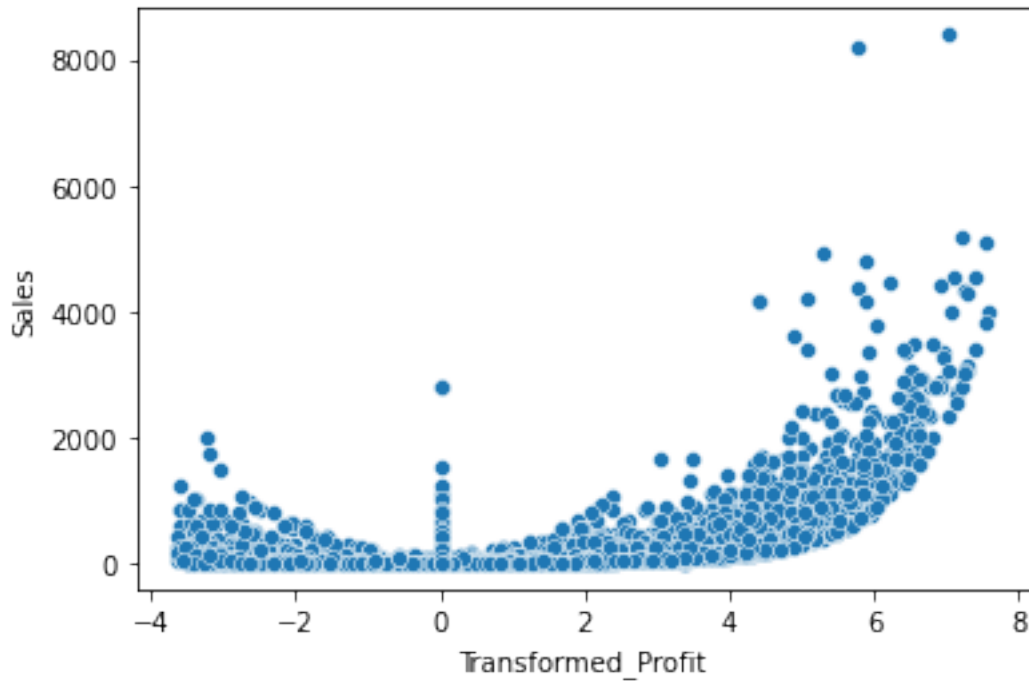
```
[123]: data = df  
  
# Select the column you want to analyze  
column_name = 'Transformed_Profit'  
column_data = data[column_name]  
  
# Create the QQ plot using statsmodels.api  
sm.qqplot(column_data, line='s') # 's' for straight reference line
```

```
plt.xlabel(column_name)
plt.ylabel('Quantiles')
plt.title('QQ Plot for ' + column_name)
plt.grid(True)
plt.show()
```



```
[124]: #Scatter plot sales vs transformed_profit
sns.scatterplot(x="Transformed_Profit",y="Sales",data=df)
```

```
[124]: <AxesSubplot:xlabel='Transformed_Profit', ylabel='Sales'>
```

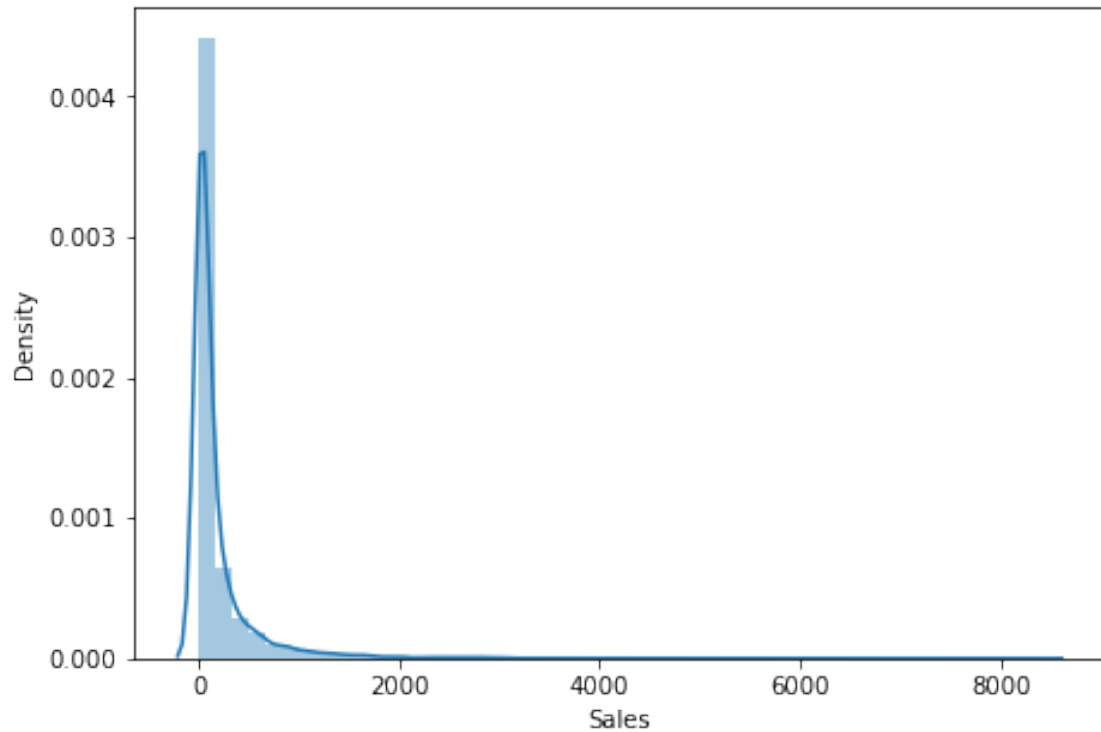


7 Sales Outlier Analysis

[125]: *# sales data distribution graph*

```
plt.figure(figsize=(16,5))
plt.subplot(1,2,2)
sns.distplot(df["Sales"])
plt.show()
```

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn\distributions.py:2619:
FutureWarning: `distplot` is a deprecated function and will be removed in a
future version. Please adapt your code to use either `displot` (a figure-level
function with similar flexibility) or `histplot` (an axes-level function for
histograms).
warnings.warn(msg, FutureWarning)

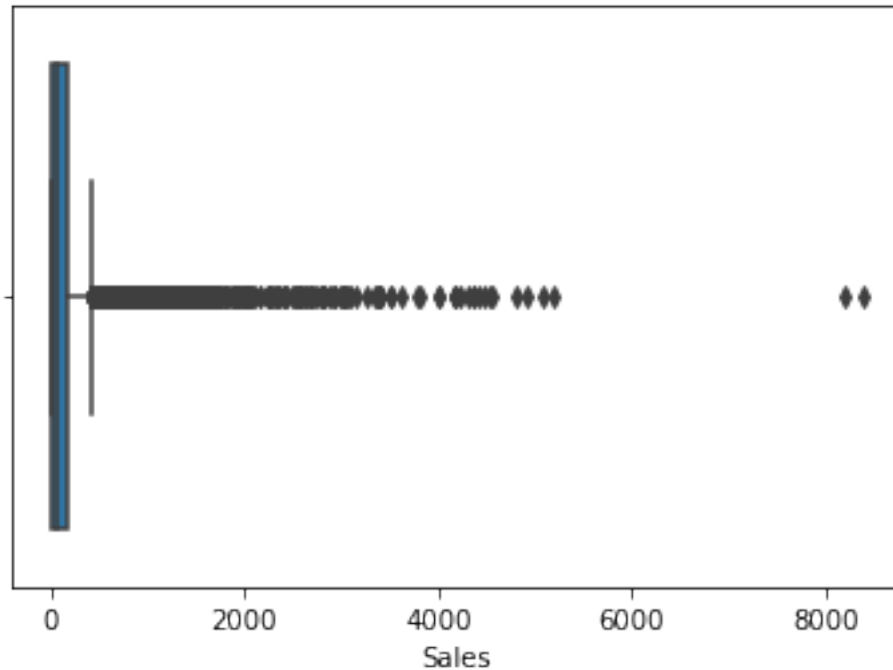


```
[126]: #box plot to check outliers  
sns.boxplot(df["Sales"])
```

```
C:\Users\ganesh\anaconda3\lib\site-packages\seaborn\_decorators.py:36:  
FutureWarning: Pass the following variable as a keyword arg: x. From version  
0.12, the only valid positional argument will be `data`, and passing other  
arguments without an explicit keyword will result in an error or  
misinterpretation.  
  warnings.warn(  

```

```
[126]: <AxesSubplot:xlabel='Sales'>
```



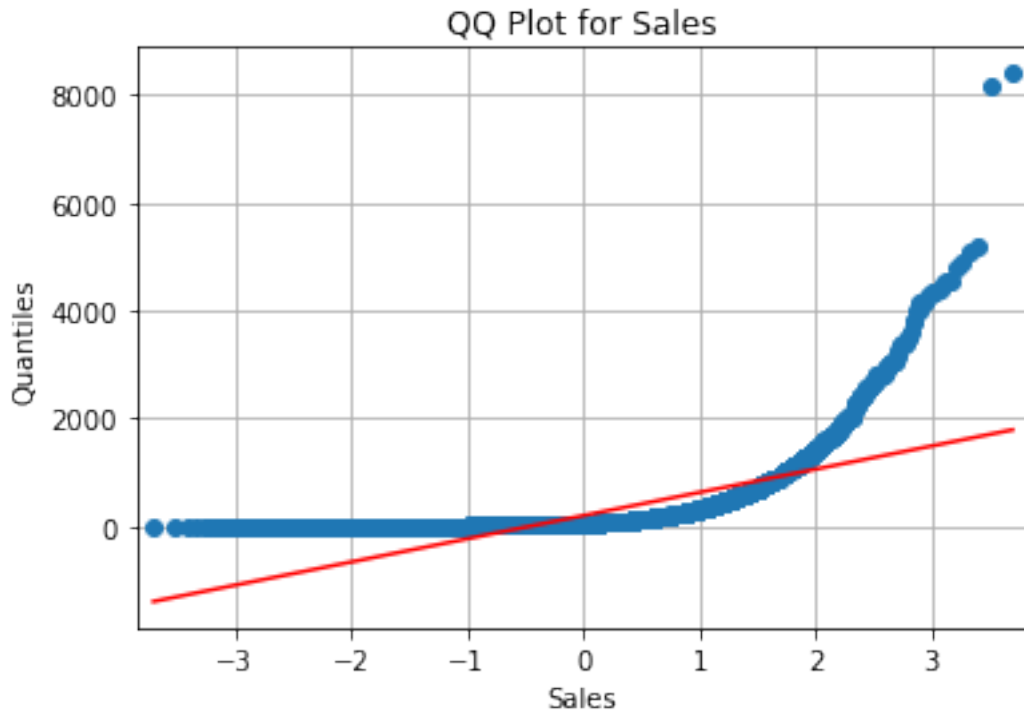
```
[127]: #checking skewness in data
df["Sales"].skew() #right skew
```

```
[127]: 5.945170113821082
```

```
[128]: data = df

# Select the column you want to analyze
column_name = 'Sales'
column_data = data[column_name]

# Create the QQ plot using statsmodels.api
sm.qqplot(column_data, line='s') # 's' for straight reference line
plt.xlabel(column_name)
plt.ylabel('Quantiles')
plt.title('QQ Plot for ' + column_name)
plt.grid(True)
plt.show()
```



```
[129]: #data transformation using log
df['Transformed_Sales'] = np.log(df['Sales'])
```

```
[130]: df.head() #top 5 rows display
```

```
[130]:
```

	Ship Mode	Segment	Country	City	State	\
0	Second Class	Consumer	United States	Henderson	Kentucky	
1	Second Class	Consumer	United States	Henderson	Kentucky	
2	Second Class	Corporate	United States	Los Angeles	California	
4	Standard Class	Consumer	United States	Fort Lauderdale	Florida	
5	Standard Class	Consumer	United States	Los Angeles	California	

	Region	Category	Sub-Category	\
0	South	Furniture	Bookcases	
1	South	Furniture	Chairs	
2	West	Office Supplies	Labels	
4	South	Office Supplies	Storage	
5	West	Furniture	Furnishings	

	Product Name	Sales	Quantity	\
0	Bush Somerset Collection Bookcase	261.960	2.0	
1	Hon Deluxe Fabric Upholstered Stacking Chairs,...	731.940	3.0	
2	Self-Adhesive Address Labels for Typewriters b...	14.620	2.0	

4	Eldon Fold 'N Roll Cart System	22.368	2.0
5	Eldon Expressions Wood and Plastic Desk Access...	48.860	7.0

	Discount	Profit	order_year	order_month	order_date	Ship_year \
0	0.0	41.9136	2016.0	11.0	8.0	2016.0
1	0.0	219.5820	2016.0	11.0	8.0	2016.0
2	0.0	6.8714	2016.0	6.0	12.0	2016.0
4	0.2	2.5164	2015.0	10.0	11.0	2015.0
5	0.0	14.1694	2014.0	6.0	9.0	2014.0

	Ship_month	Ship_date	Transformed_Profit	Transformed_Sales
0	11.0	11.0	3.735610	5.568192
1	11.0	11.0	5.391726	6.595699
2	6.0	16.0	1.927368	2.682390
4	10.0	18.0	0.922829	3.107631
5	6.0	14.0	2.651085	3.888959

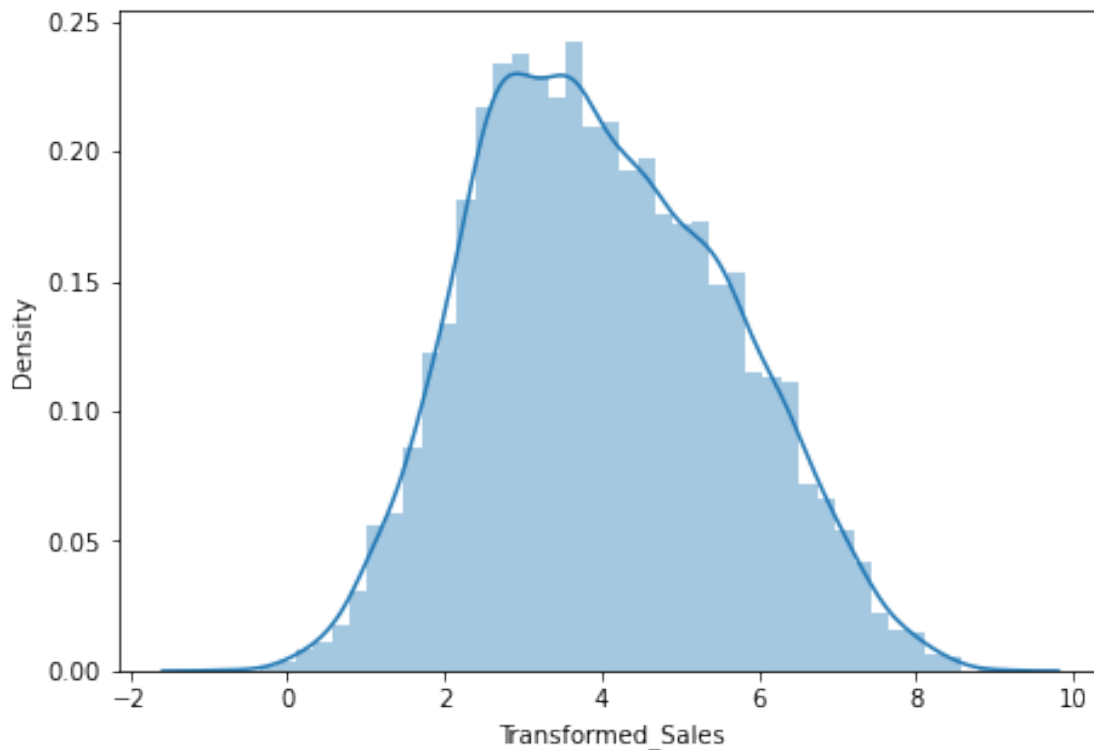
```
[131]: df.shape
```

```
[131]: (9233, 21)
```

```
[132]: #distribution graph after transformation
plt.figure(figsize=(16,5))
plt.subplot(1,2,2)
sns.distplot(df["Transformed_Sales"])
plt.show()
```

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn\distributions.py:2619:
FutureWarning: `distplot` is a deprecated function and will be removed in a
future version. Please adapt your code to use either `displot` (a figure-level
function with similar flexibility) or `histplot` (an axes-level function for
histograms).

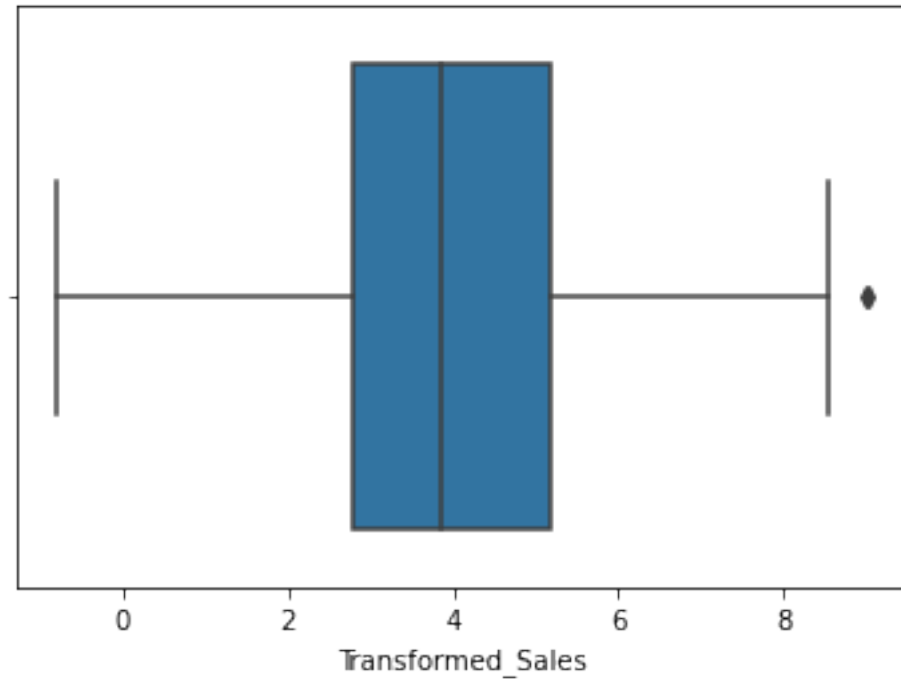
```
warnings.warn(msg, FutureWarning)
```



```
[133]: #boxplot
sns.boxplot(df["Transformed_Sales"])
```

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn_decorators.py:36:
FutureWarning: Pass the following variable as a keyword arg: x. From version
0.12, the only valid positional argument will be `data`, and passing other
arguments without an explicit keyword will result in an error or
misinterpretation.
warnings.warn(

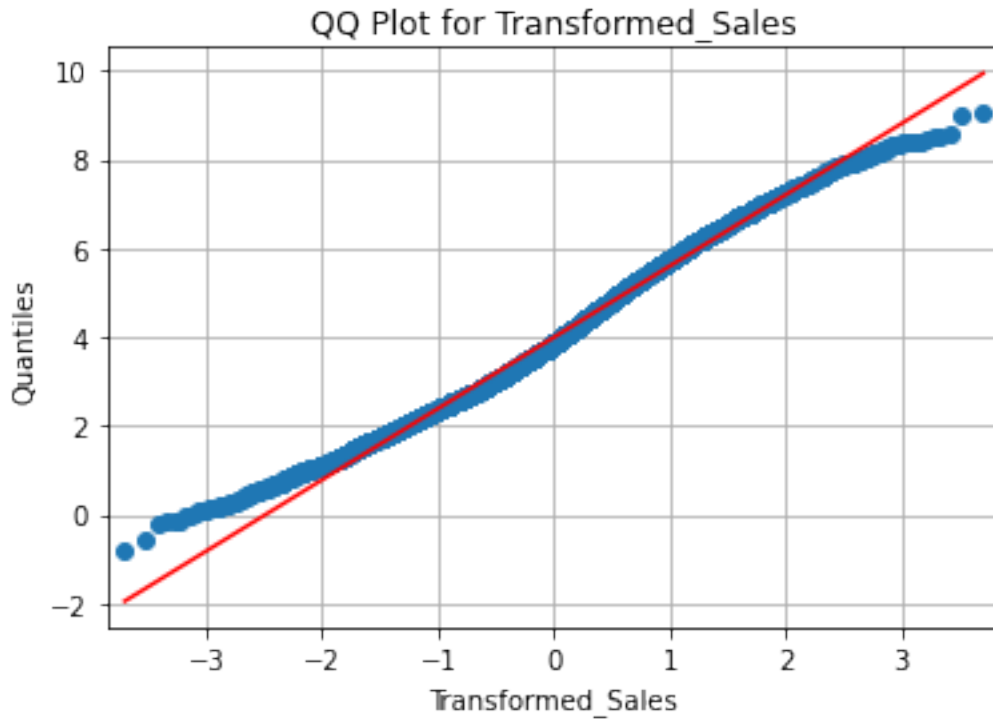
```
[133]: <AxesSubplot:xlabel='Transformed_Sales'>
```



```
[134]: #right skewness  
df["Transformed_Sales"].skew()
```

```
[134]: 0.2193355558713904
```

```
[135]: import statsmodels.api as sm  
data = df  
  
# Select the column you want to analyze  
column_name = 'Transformed_Sales'  
column_data = data[column_name]  
  
# Create the QQ plot using statsmodels.api  
sm.qqplot(column_data, line='s') # 's' for straight reference line  
plt.xlabel(column_name)  
plt.ylabel('Quantiles')  
plt.title('QQ Plot for ' + column_name)  
plt.grid(True)  
plt.show()
```



```
[136]: df.shape
```

```
[136]: (9233, 21)
```

```
[137]: #z_score technique to delete outliers
z_scores = (df['Transformed_Sales'] - df['Transformed_Sales'].mean()) /
           df['Transformed_Sales'].std()
outliers = df[np.abs(z_scores) > 3]
```

```
[138]: #drop outliers
df = df.drop(outliers.index)
```

```
[139]: df.shape
```

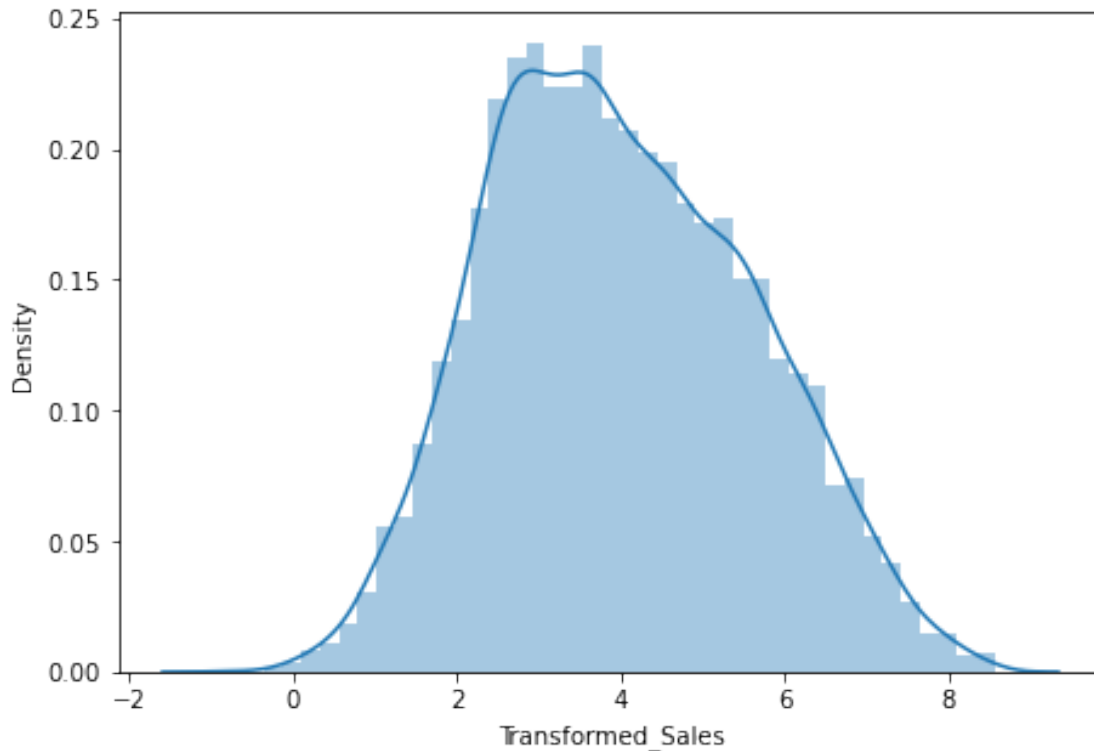
```
[139]: (9231, 21)
```

```
[140]: #distribution graph after outlier deletion
plt.figure(figsize=(16,5))
plt.subplot(1,2,2)
sns.distplot(df["Transformed_Sales"])
plt.show()
```

C:\Users\ganesh\anaconda3\lib\site-packages\seaborn\distributions.py:2619:
FutureWarning: `distplot` is a deprecated function and will be removed in a

future version. Please adapt your code to use either ``displot`` (a figure-level function with similar flexibility) or ``histplot`` (an axes-level function for histograms).

```
warnings.warn(msg, FutureWarning)
```

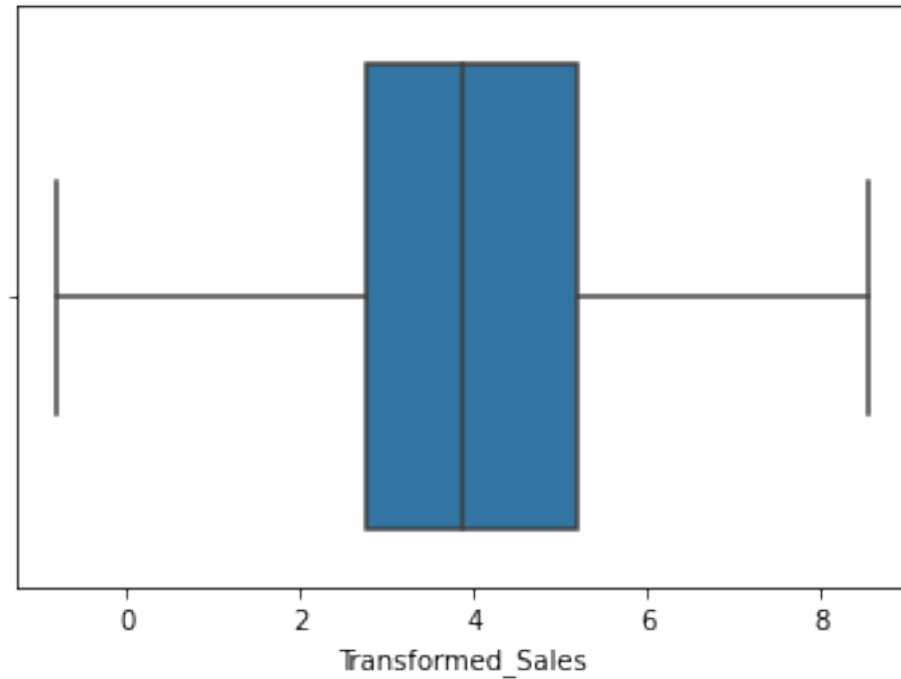


```
[141]: #box plot
sns.boxplot(df["Transformed_Sales"])
```

```
C:\Users\ganesh\anaconda3\lib\site-packages\seaborn\_decorators.py:36:
FutureWarning: Pass the following variable as a keyword arg: x. From version
0.12, the only valid positional argument will be `data`, and passing other
arguments without an explicit keyword will result in an error or
misinterpretation.
```

```
warnings.warn(
```

```
[141]: <AxesSubplot:xlabel='Transformed_Sales'>
```

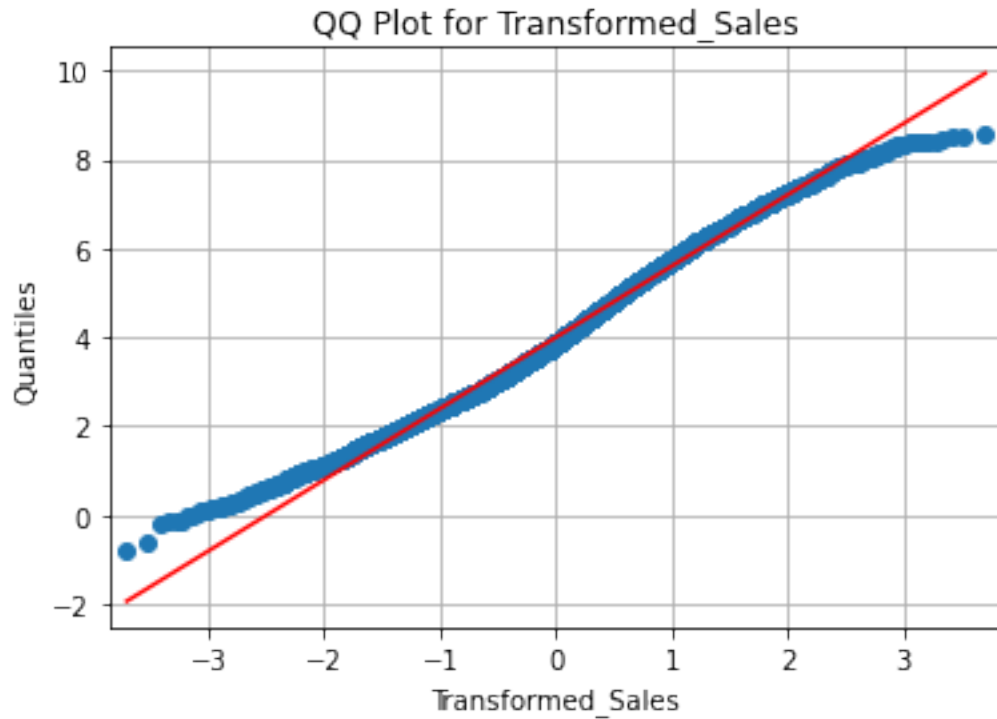
```
[142]: #skew
df["Transformed_Sales"].skew()
```

```
[142]: 0.21540098687838372
```

```
[143]: data = df

# Select the column you want to analyze
column_name = 'Transformed_Sales'
column_data = data[column_name]

# Create the QQ plot using statsmodels.api
sm.qqplot(column_data, line='s') # 's' for straight reference line
plt.xlabel(column_name)
plt.ylabel('Quantiles')
plt.title('QQ Plot for ' + column_name)
plt.grid(True)
plt.show()
```



```
[144]: df
```

```
[144]:
```

	Ship Mode	Segment	Country	City	State \
0	Second Class	Consumer	United States	Henderson	Kentucky
1	Second Class	Consumer	United States	Henderson	Kentucky
2	Second Class	Corporate	United States	Los Angeles	California
4	Standard Class	Consumer	United States	Fort Lauderdale	Florida
5	Standard Class	Consumer	United States	Los Angeles	California
...
9867	Standard Class	Consumer	United States	Detroit	Michigan
9868	Standard Class	Consumer	United States	Detroit	Michigan
9869	Standard Class	Consumer	United States	Detroit	Michigan
9870	Second Class	Corporate	United States	Fort Lauderdale	Florida
9871	Second Class	Consumer	United States	Hampton	Virginia

	Region	Category	Sub-Category \
0	South	Furniture	Bookcases
1	South	Furniture	Chairs
2	West	Office Supplies	Labels
4	South	Office Supplies	Storage
5	West	Furniture	Furnishings
...
9867	Central	Office Supplies	Binders

9868	Central	Office Supplies	Art
9869	Central	Office Supplies	Binders
9870	South	Office Supplies	Appliances
9871	South	Office Supplies	Appliances

		Product Name	Sales	Quantity \
0		Bush Somerset Collection Bookcase	261.960	2.0
1		Hon Deluxe Fabric Upholstered Stacking Chairs,...	731.940	3.0
2		Self-Adhesive Address Labels for Typewriters b...	14.620	2.0
4		Eldon Fold 'N Roll Cart System	22.368	2.0
5		Eldon Expressions Wood and Plastic Desk Access...	48.860	7.0
...	
9867		GBC Plastic Binding Combs	29.520	4.0
9868		Newell 315	11.960	2.0
9869		Wilson Jones 1" Hanging DublLock Ring Binders	26.400	5.0
9870		Hoover Upright Vacuum With Dirt Cup	1158.120	5.0
9871		Holmes Replacement Filter for HEPA Air Cleaner...	44.430	3.0

	Discount	Profit	order_year	order_month	order_date	Ship_year \
0	0.0	41.9136	2016.0	11.0	8.0	2016.0
1	0.0	219.5820	2016.0	11.0	8.0	2016.0
2	0.0	6.8714	2016.0	6.0	12.0	2016.0
4	0.2	2.5164	2015.0	10.0	11.0	2015.0
5	0.0	14.1694	2014.0	6.0	9.0	2014.0
...
9867	0.0	14.4648	2016.0	9.0	1.0	2016.0
9868	0.0	2.9900	2016.0	9.0	1.0	2016.0
9869	0.0	12.6720	2016.0	9.0	1.0	2016.0
9870	0.2	130.2885	2017.0	11.0	11.0	2017.0
9871	0.0	18.6606	2015.0	11.0	8.0	2015.0

	Ship_month	Ship_date	Transformed_Profit	Transformed_Sales
0	11.0	11.0	3.735610	5.568192
1	11.0	11.0	5.391726	6.595699
2	6.0	16.0	1.927368	2.682390
4	10.0	18.0	0.922829	3.107631
5	6.0	14.0	2.651085	3.888959
...
9867	9.0	5.0	2.671718	3.385068
9868	9.0	5.0	1.095273	2.481568
9869	9.0	5.0	2.539395	3.273364
9870	11.0	16.0	4.869751	7.054553
9871	11.0	13.0	2.926414	3.793915

[9231 rows x 21 columns]

8 Feature Engineering:

9 Identify relevant features that could impact sales, like trend or seasonality.

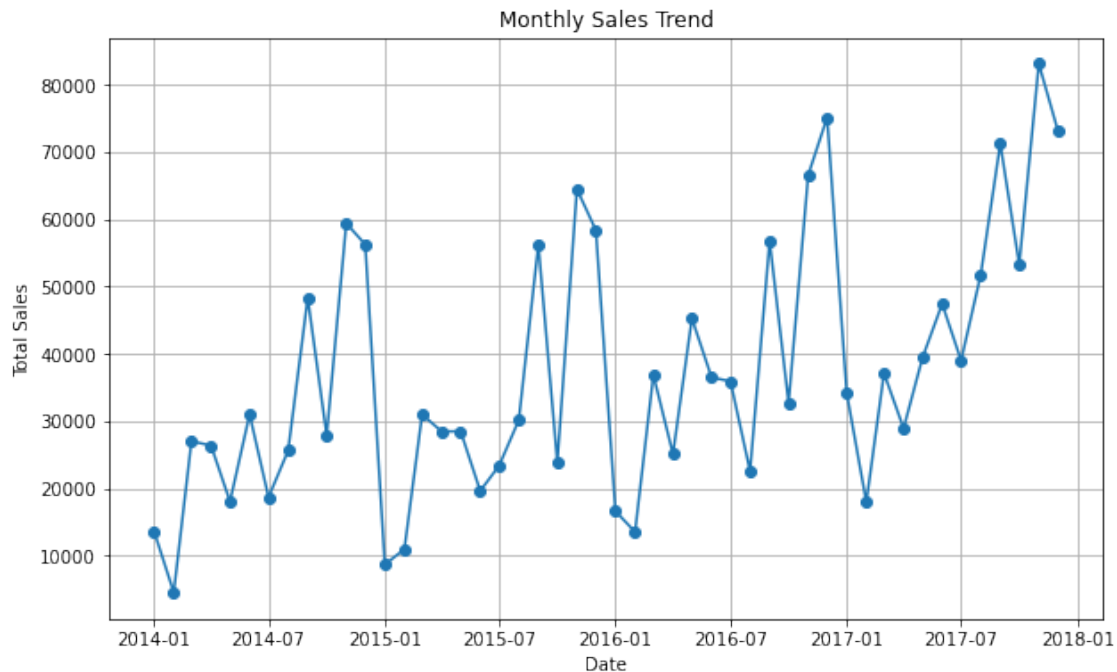
```
[145]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

# Assuming you have loaded your sales data into a pandas DataFrame named
# 'sales_data' with columns 'Order_year', 'Order_month', and 'sales'
dfs=df.copy()
# Convert 'Order_year' and 'Order_month' to integers
dfs['Order_year'] = dfs['order_year'].astype(int)
dfs['Order_month'] = dfs['order_month'].astype(int)

# Create a new column for 'order_date' combining year and month
dfs['order_date'] = pd.to_datetime(dfs['order_year'] * 100 +
    dfs['order_month'], format='%Y%m')

# Group by 'order_date' and calculate total sales
monthly_sales = dfs.groupby('order_date')['Sales'].sum().reset_index()

# Plot time series of sales
plt.figure(figsize=(10, 6))
plt.plot(monthly_sales['order_date'], monthly_sales['Sales'], marker='o',
    linestyle='-')
plt.title('Monthly Sales Trend')
plt.xlabel('Date')
plt.ylabel('Total Sales')
plt.grid(True)
plt.show()
```



```
[146]: # Analyzing sales by ship mode
ship_mode_sales = df.groupby('Ship Mode')['Sales'].sum().
    ↪sort_values(ascending=False)
plt.figure(figsize=(8, 5))
ship_mode_sales.plot(kind='bar')
plt.title('Sales by Ship Mode')
plt.xlabel('Ship Mode')
plt.ylabel('Sales')
plt.xticks(rotation=45)
plt.show()

# Analyzing sales by segment
segment_sales = df.groupby('Segment')['Sales'].sum().
    ↪sort_values(ascending=False)
plt.figure(figsize=(8, 5))
segment_sales.plot(kind='bar')
plt.title('Sales by Customer Segment')
plt.xlabel('Segment')
plt.ylabel('Sales')
plt.xticks(rotation=45)
plt.show()

# Analyzing sales by region
region_sales = df.groupby('Region')['Sales'].sum().sort_values(ascending=False)
plt.figure(figsize=(8, 5))
```

```

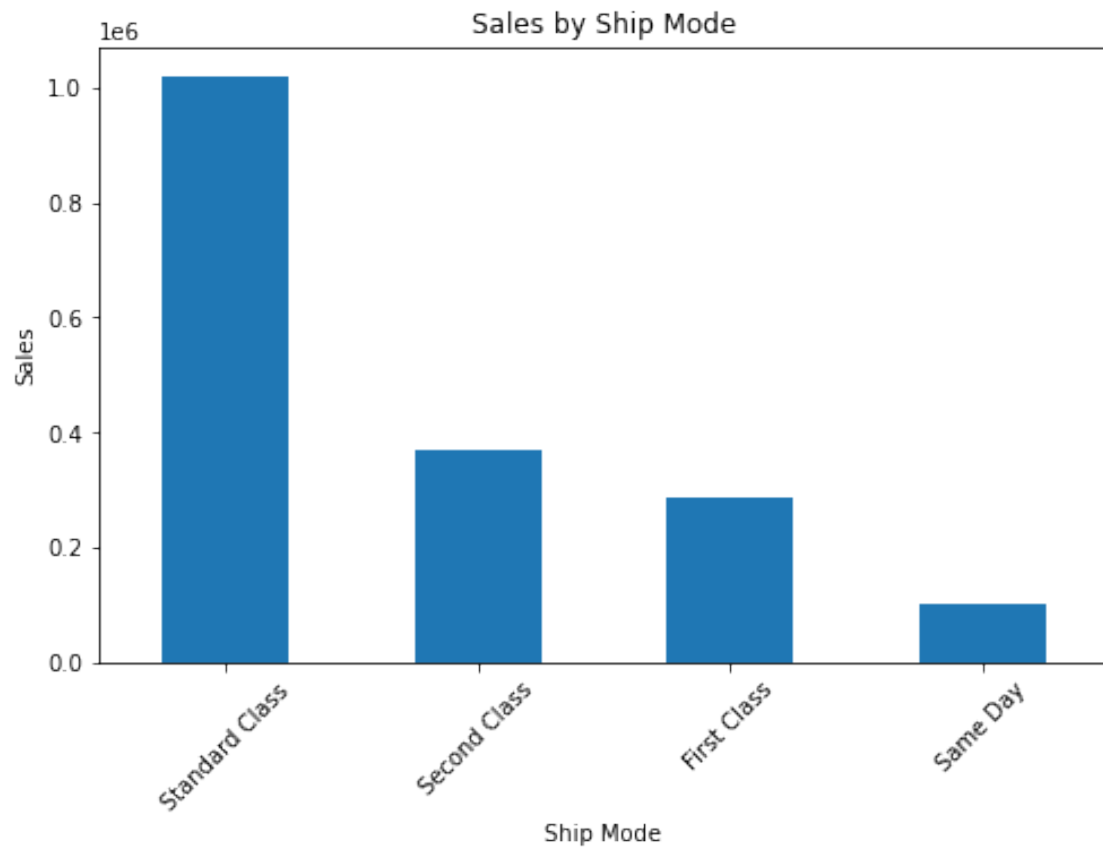
region_sales.plot(kind='bar')
plt.title('Sales by Region')
plt.xlabel('Region')
plt.ylabel('Sales')
plt.xticks(rotation=45)
plt.show()

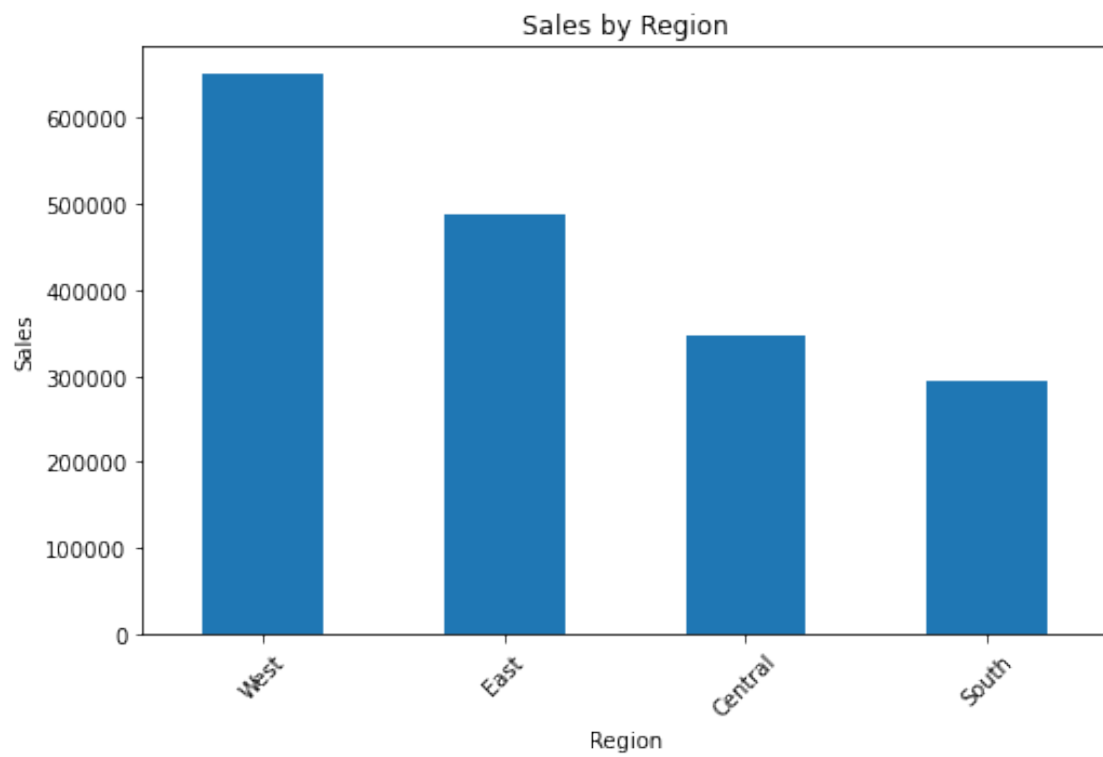
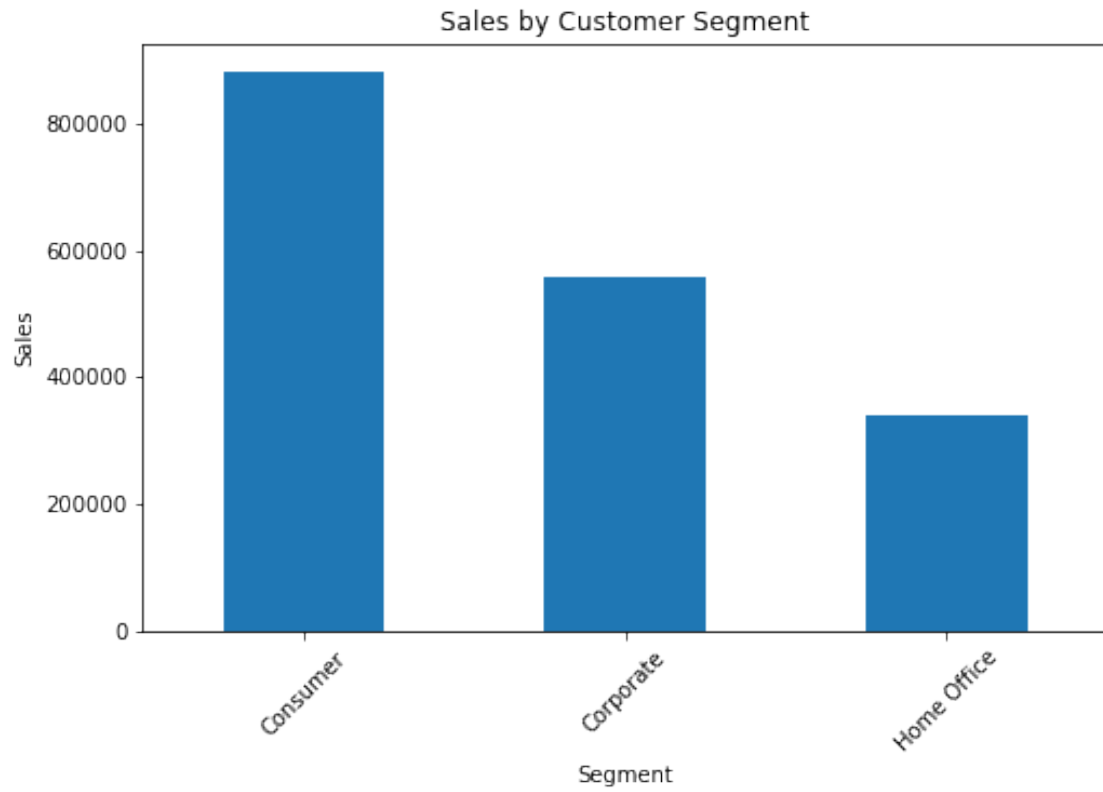
# Analyzing sales by category
category_sales = df.groupby('Category')['Sales'].sum().
    ↪sort_values(ascending=False)
plt.figure(figsize=(8, 5))
category_sales.plot(kind='bar')
plt.title('Sales by Category')
plt.xlabel('Category')
plt.ylabel('Sales')
plt.xticks(rotation=45)
plt.show()

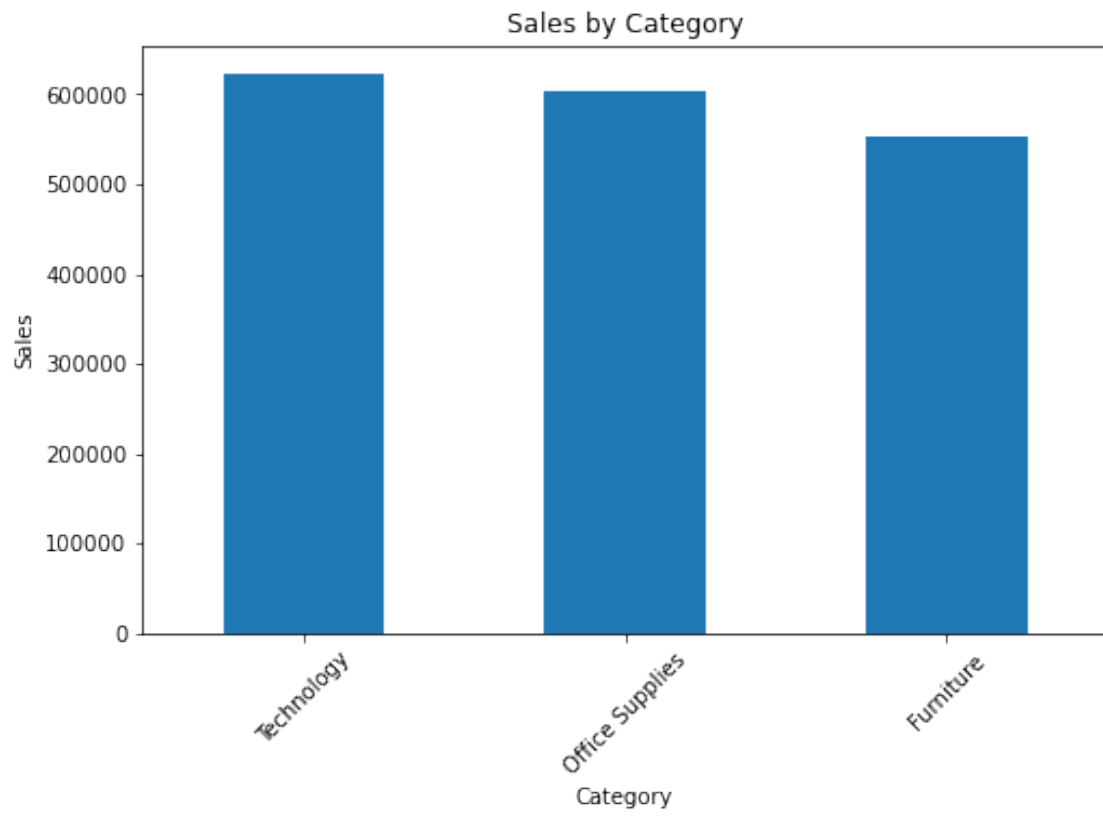
# Analyzing the impact of discount on sales
plt.figure(figsize=(8, 5))
sns.scatterplot(x='Discount', y='Sales', data=df)
plt.title('Impact of Discount on Sales')
plt.xlabel('Discount')
plt.ylabel('Sales')
plt.show()

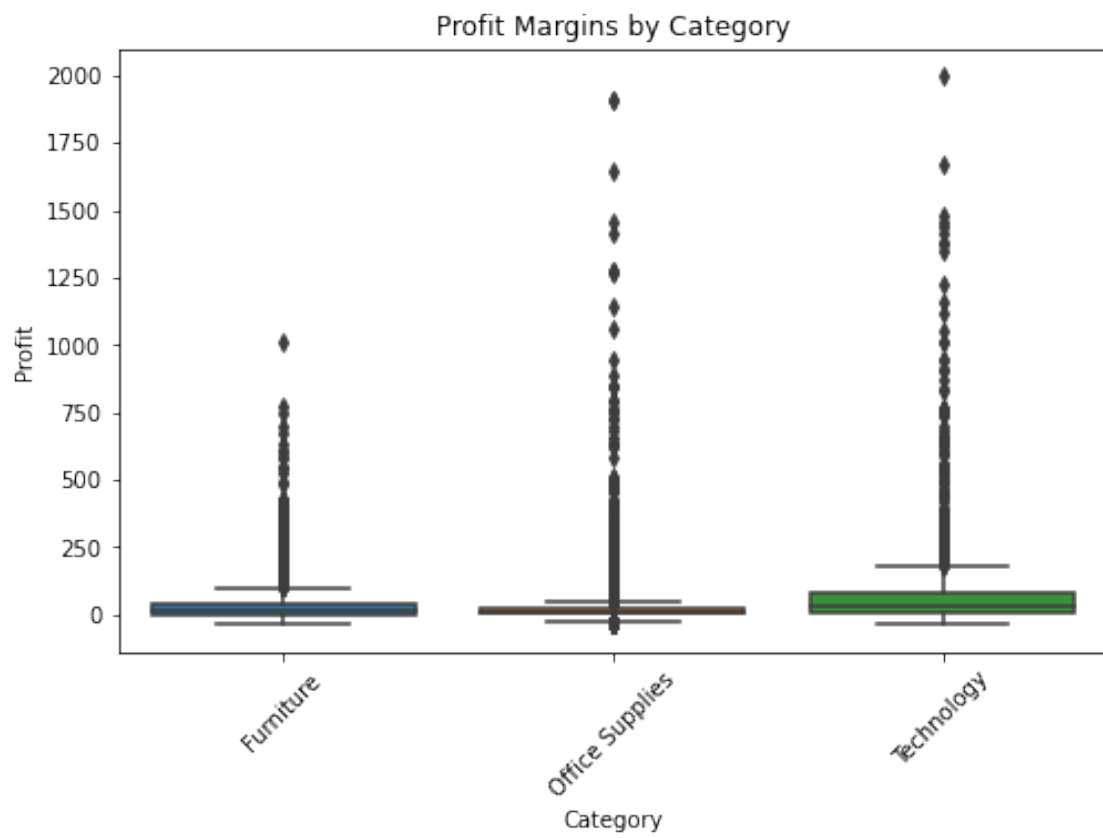
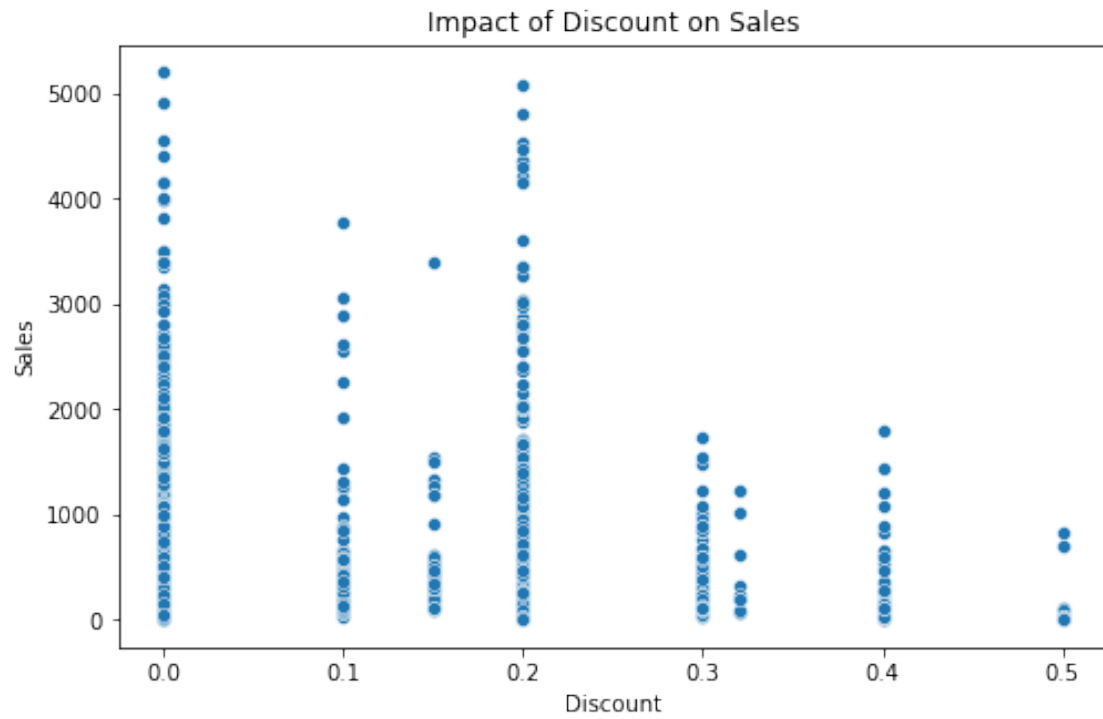
# Analyzing profit margins
plt.figure(figsize=(8, 5))
sns.boxplot(x='Category', y='Profit', data=df)
plt.title('Profit Margins by Category')
plt.xlabel('Category')
plt.ylabel('Profit')
plt.xticks(rotation=45)
plt.show()

```





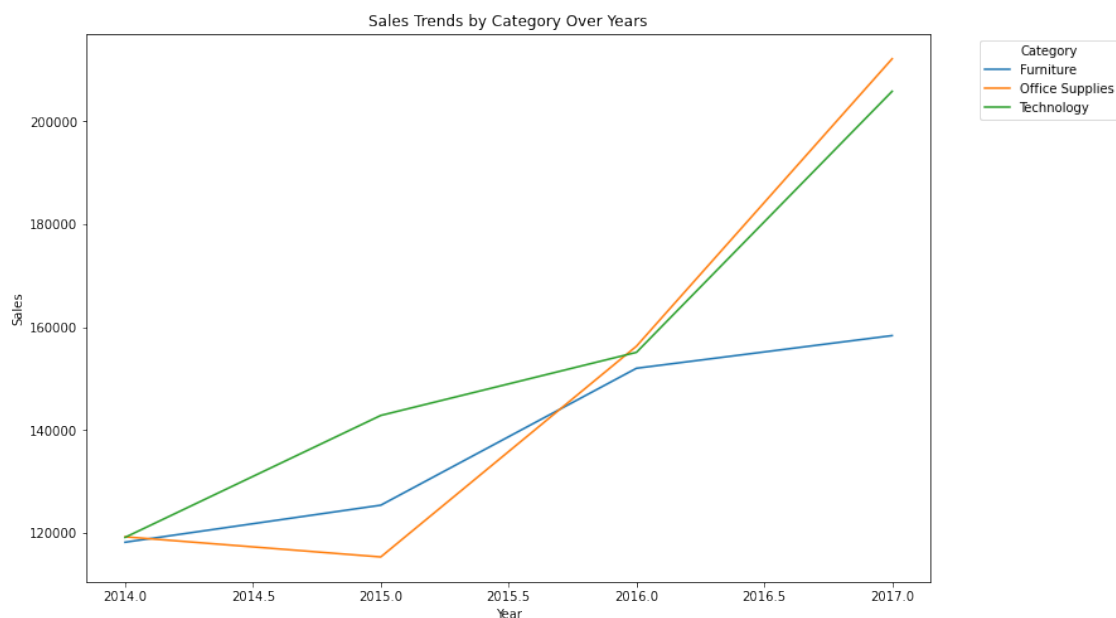




10 Investigate sales patterns for different product categories (Increasing or decreasing over years).

```
[147]: category_yearly_sales = df.groupby(['Category', 'order_year'])['Sales'].sum().  
        ↪reset_index()  
plt.figure(figsize=(12, 8))  
sns.lineplot(x='order_year', y='Sales', hue='Category',  
        ↪data=category_yearly_sales)  
plt.title('Sales Trends by Category Over Years')  
plt.xlabel('Year')  
plt.ylabel('Sales')  
plt.legend(title='Category', bbox_to_anchor=(1.05, 1), loc='upper left')
```

```
[147]: <matplotlib.legend.Legend at 0x21b0d020550>
```



11 Analyze the correlation between sales and other variables within each category.

```
[148]: category_correlations = {}  
for category in df['Category'].unique():  
    category_data = df[df['Category'] == category]
```

```

category_correlation_matrix = category_data.corr()
category_correlations[category] = category_correlation_matrix

for category, correlation_matrix in category_correlations.items():
    print(f"Correlation Matrix for {category} Category:")
    print(correlation_matrix)
    print()

```

Correlation Matrix for Furniture Category:

	Sales	Quantity	Discount	Profit	order_year	\
Sales	1.000000	0.414408	-0.085544	0.765578	-0.045909	
Quantity	0.414408	1.000000	-0.078189	0.362231	-0.070076	
Discount	-0.085544	-0.078189	1.000000	-0.333928	-0.034070	
Profit	0.765578	0.362231	-0.333928	1.000000	-0.030386	
order_year	-0.045909	-0.070076	-0.034070	-0.030386	1.000000	
order_month	0.034092	-0.001358	0.015186	0.026142	-0.012841	
order_date	0.004778	0.012220	0.007750	-0.003917	-0.003544	
Ship_year	-0.047723	-0.069455	-0.031428	-0.030472	0.993575	
Ship_month	0.041604	0.001825	0.003635	0.026501	-0.009276	
Ship_date	0.003479	-0.033666	0.010094	-0.000326	-0.036496	
Transformed_Profit	0.374131	0.247231	-0.675976	0.625102	-0.016565	
Transformed_Sales	0.741005	0.364254	-0.120222	0.534449	-0.030388	

	order_month	order_date	Ship_year	Ship_month	Ship_date	\
Sales	0.034092	0.004778	-0.047723	0.041604	0.003479	
Quantity	-0.001358	0.012220	-0.069455	0.001825	-0.033666	
Discount	0.015186	0.007750	-0.031428	0.003635	0.010094	
Profit	0.026142	-0.003917	-0.030472	0.026501	-0.000326	
order_year	-0.012841	-0.003544	0.993575	-0.009276	-0.036496	
order_month	1.000000	-0.054651	0.005182	0.906667	-0.028408	
order_date	-0.054651	1.000000	0.017885	-0.086814	0.385182	
Ship_year	0.005182	0.017885	1.000000	-0.039000	-0.058611	
Ship_month	0.906667	-0.086814	-0.039000	1.000000	0.008354	
Ship_date	-0.028408	0.385182	-0.058611	0.008354	1.000000	
Transformed_Profit	0.032409	-0.008701	-0.016583	0.031189	0.014099	
Transformed_Sales	0.012867	0.000152	-0.032783	0.023049	-0.003641	

	Transformed_Profit	Transformed_Sales
Sales	0.374131	0.741005
Quantity	0.247231	0.364254
Discount	-0.675976	-0.120222
Profit	0.625102	0.534449
order_year	-0.016565	-0.030388
order_month	0.032409	0.012867
order_date	-0.008701	0.000152
Ship_year	-0.016583	-0.032783
Ship_month	0.031189	0.023049

Ship_date	0.014099	-0.003641
Transformed_Profit	1.000000	0.351710
Transformed_Sales	0.351710	1.000000

Correlation Matrix for Office Supplies Category:

	Sales	Quantity	Discount	Profit	order_year	\
Sales	1.000000	0.167235	-0.120066	0.857103	0.014741	
Quantity	0.167235	1.000000	-0.016831	0.142620	0.022021	
Discount	-0.120066	-0.016831	1.000000	-0.153288	0.000224	
Profit	0.857103	0.142620	-0.153288	1.000000	0.017693	
order_year	0.014741	0.022021	0.000224	0.017693	1.000000	
order_month	-0.007454	0.024422	0.000188	0.004639	-0.011874	
order_date	0.005073	-0.003157	-0.003023	0.000185	-0.021896	
Ship_year	0.014565	0.020591	-0.000291	0.017404	0.993790	
Ship_month	-0.007163	0.029372	0.000444	0.005948	-0.001428	
Ship_date	0.012367	0.005077	0.020189	-0.001403	-0.013676	
Transformed_Profit	0.435989	0.255263	-0.618825	0.502668	0.017911	
Transformed_Sales	0.628206	0.373706	-0.340726	0.526248	0.015244	

	order_month	order_date	Ship_year	Ship_month	Ship_date	\
Sales	-0.007454	0.005073	0.014565	-0.007163	0.012367	
Quantity	0.024422	-0.003157	0.020591	0.029372	0.005077	
Discount	0.000188	-0.003023	-0.000291	0.000444	0.020189	
Profit	0.004639	0.000185	0.017404	0.005948	-0.001403	
order_year	-0.011874	-0.021896	0.993790	-0.001428	-0.013676	
order_month	1.000000	-0.031090	0.006544	0.906240	-0.020464	
order_date	-0.031090	1.000000	-0.000052	-0.066594	0.389597	
Ship_year	0.006544	-0.000052	1.000000	-0.030102	-0.034899	
Ship_month	0.906240	-0.066594	-0.030102	1.000000	0.013945	
Ship_date	-0.020464	0.389597	-0.034899	0.013945	1.000000	
Transformed_Profit	0.016034	-0.003491	0.016949	0.020760	-0.022837	
Transformed_Sales	0.007906	0.000924	0.015433	0.007314	-0.006283	

	Transformed_Profit	Transformed_Sales
Sales	0.435989	0.628206
Quantity	0.255263	0.373706
Discount	-0.618825	-0.340726
Profit	0.502668	0.526248
order_year	0.017911	0.015244
order_month	0.016034	0.007906
order_date	-0.003491	0.000924
Ship_year	0.016949	0.015433
Ship_month	0.020760	0.007314
Ship_date	-0.022837	-0.006283
Transformed_Profit	1.000000	0.700100
Transformed_Sales	0.700100	1.000000

Correlation Matrix for Technology Category:

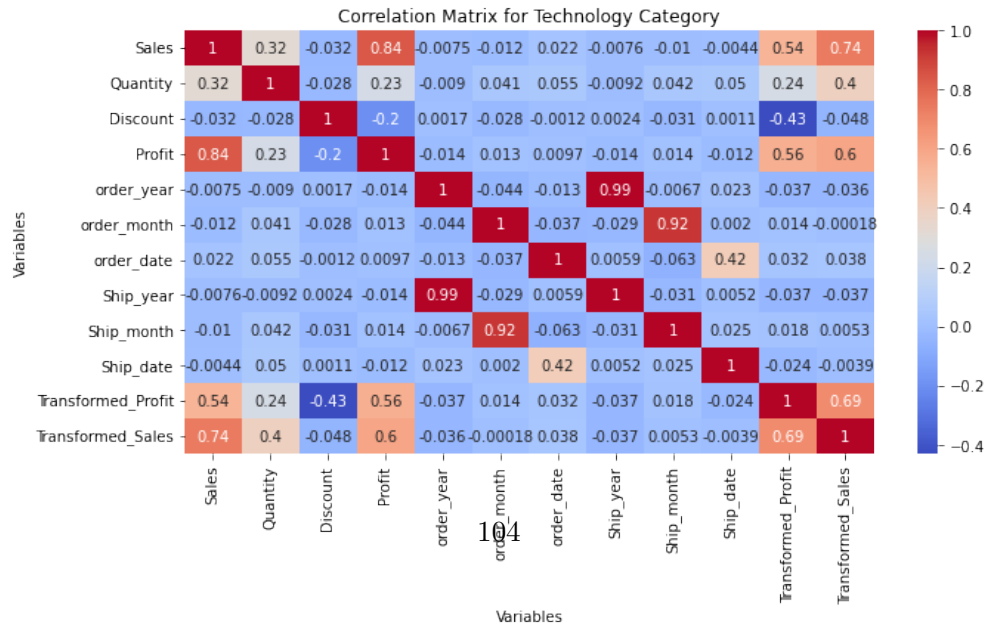
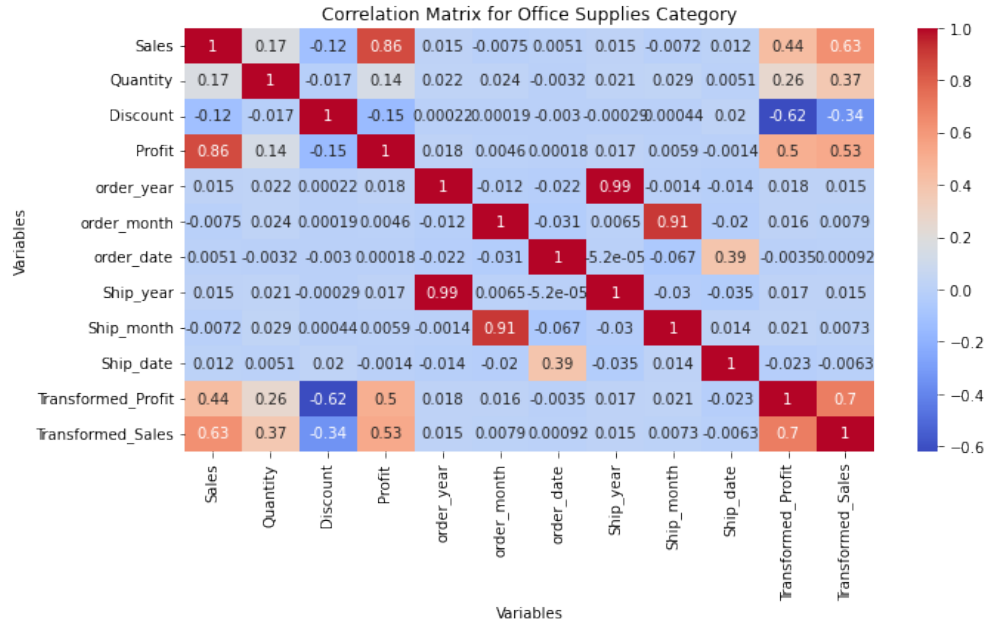
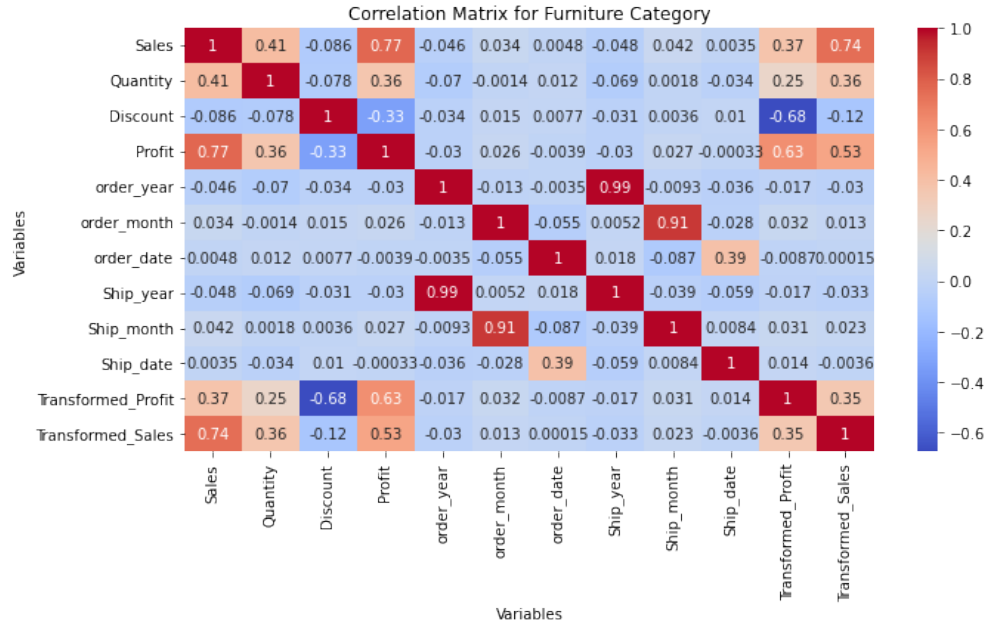
	Sales	Quantity	Discount	Profit	order_year	\
Sales	1.000000	0.319961	-0.032084	0.838539	-0.007476	
Quantity	0.319961	1.000000	-0.028492	0.234063	-0.008968	
Discount	-0.032084	-0.028492	1.000000	-0.204603	0.001747	
Profit	0.838539	0.234063	-0.204603	1.000000	-0.013748	
order_year	-0.007476	-0.008968	0.001747	-0.013748	1.000000	
order_month	-0.012163	0.041371	-0.028493	0.013373	-0.043843	
order_date	0.021743	0.054802	-0.001203	0.009719	-0.013018	
Ship_year	-0.007600	-0.009185	0.002401	-0.013821	0.994769	
Ship_month	-0.010463	0.042095	-0.031123	0.014313	-0.006737	
Ship_date	-0.004449	0.050469	0.001069	-0.012495	0.023132	
Transformed_Profit	0.539624	0.243737	-0.431459	0.558680	-0.036687	
Transformed_Sales	0.742175	0.395415	-0.048060	0.599116	-0.035851	

	order_month	order_date	Ship_year	Ship_month	Ship_date	\
Sales	-0.012163	0.021743	-0.007600	-0.010463	-0.004449	
Quantity	0.041371	0.054802	-0.009185	0.042095	0.050469	
Discount	-0.028493	-0.001203	0.002401	-0.031123	0.001069	
Profit	0.013373	0.009719	-0.013821	0.014313	-0.012495	
order_year	-0.043843	-0.013018	0.994769	-0.006737	0.023132	
order_month	1.000000	-0.037157	-0.028908	0.923858	0.001977	
order_date	-0.037157	1.000000	0.005905	-0.062725	0.416332	
Ship_year	-0.028908	0.005905	1.000000	-0.030615	0.005172	
Ship_month	0.923858	-0.062725	-0.030615	1.000000	0.025193	
Ship_date	0.001977	0.416332	0.005172	0.025193	1.000000	
Transformed_Profit	0.014378	0.031783	-0.036676	0.017842	-0.023510	
Transformed_Sales	-0.000178	0.037611	-0.036669	0.005336	-0.003888	

	Transformed_Profit	Transformed_Sales
Sales	0.539624	0.742175
Quantity	0.243737	0.395415
Discount	-0.431459	-0.048060
Profit	0.558680	0.599116
order_year	-0.036687	-0.035851
order_month	0.014378	-0.000178
order_date	0.031783	0.037611
Ship_year	-0.036676	-0.036669
Ship_month	0.017842	0.005336
Ship_date	-0.023510	-0.003888
Transformed_Profit	1.000000	0.691990
Transformed_Sales	0.691990	1.000000

```
[149]: fig, axes = plt.subplots(nrows=len(category_correlations), figsize=(10, 6 * len(category_correlations)))
```

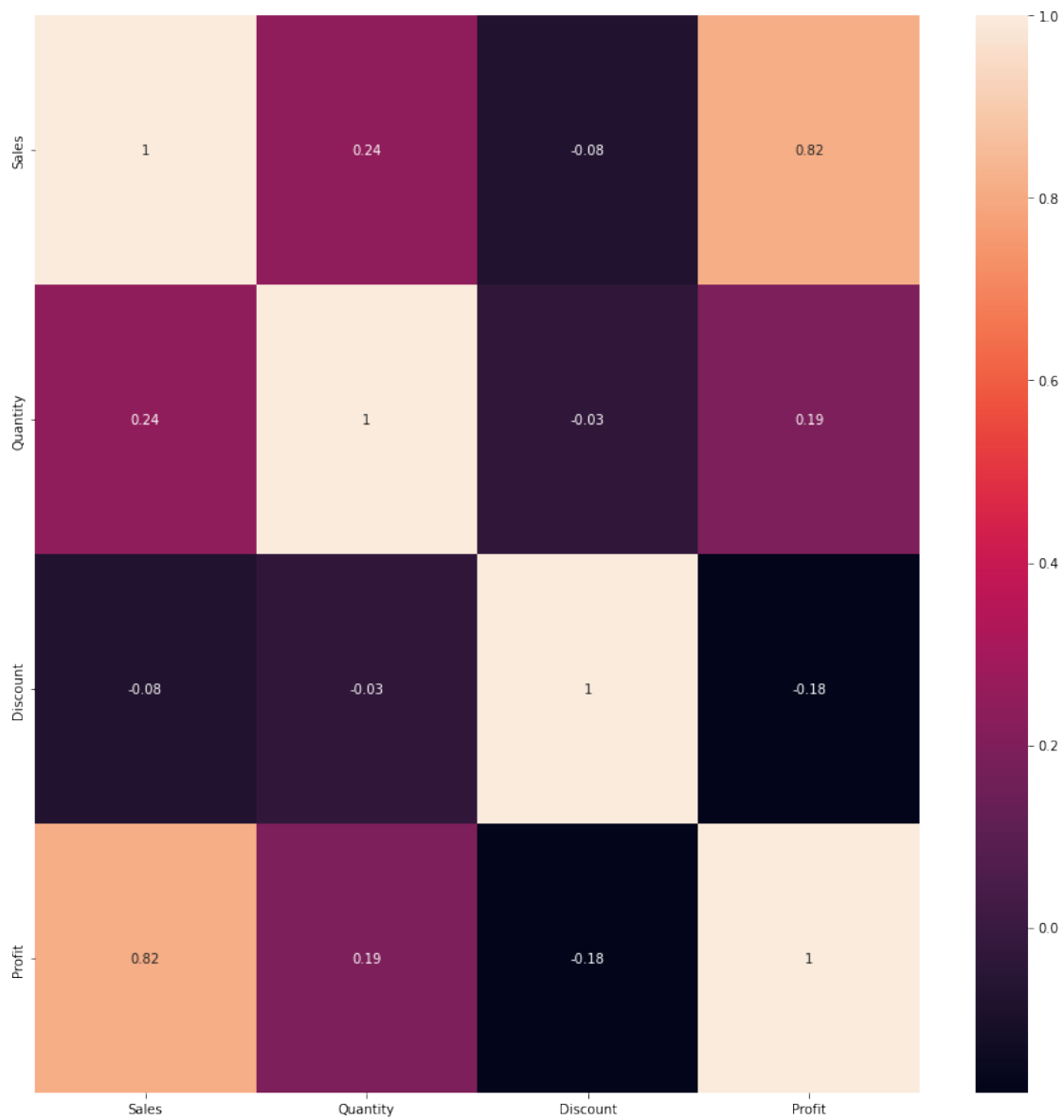
```
for i, (category, correlation_matrix) in enumerate(category_correlations.  
    ↪items()):  
    sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', ax=axes[i])  
    axes[i].set_title(f'Correlation Matrix for {category} Category')  
    axes[i].set_xlabel('Variables')  
    axes[i].set_ylabel('Variables')  
  
plt.tight_layout()  
plt.show()
```



12 Correlation Matrix

```
[150]: #Correlation between columns  
plt.figure(figsize=(15,15))  
sns.heatmap(df[['Sales', 'Quantity', 'Discount', 'Profit']].corr(),annot=True)
```

```
[150]: <AxesSubplot:>
```



```
[151]: df[['Sales', 'Quantity', 'Discount', 'Profit']].corr()
```

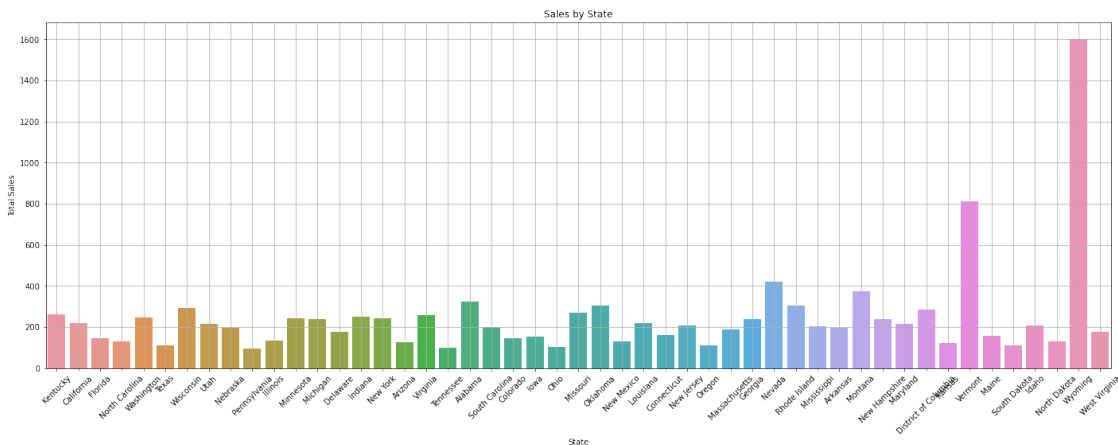
```
[151]:
```

	Sales	Quantity	Discount	Profit
Sales	1.000000	0.243554	-0.079578	0.815014
Quantity	0.243554	1.000000	-0.030481	0.192574
Discount	-0.079578	-0.030481	1.000000	-0.181416
Profit	0.815014	0.192574	-0.181416	1.000000

13 Regional Analysis:

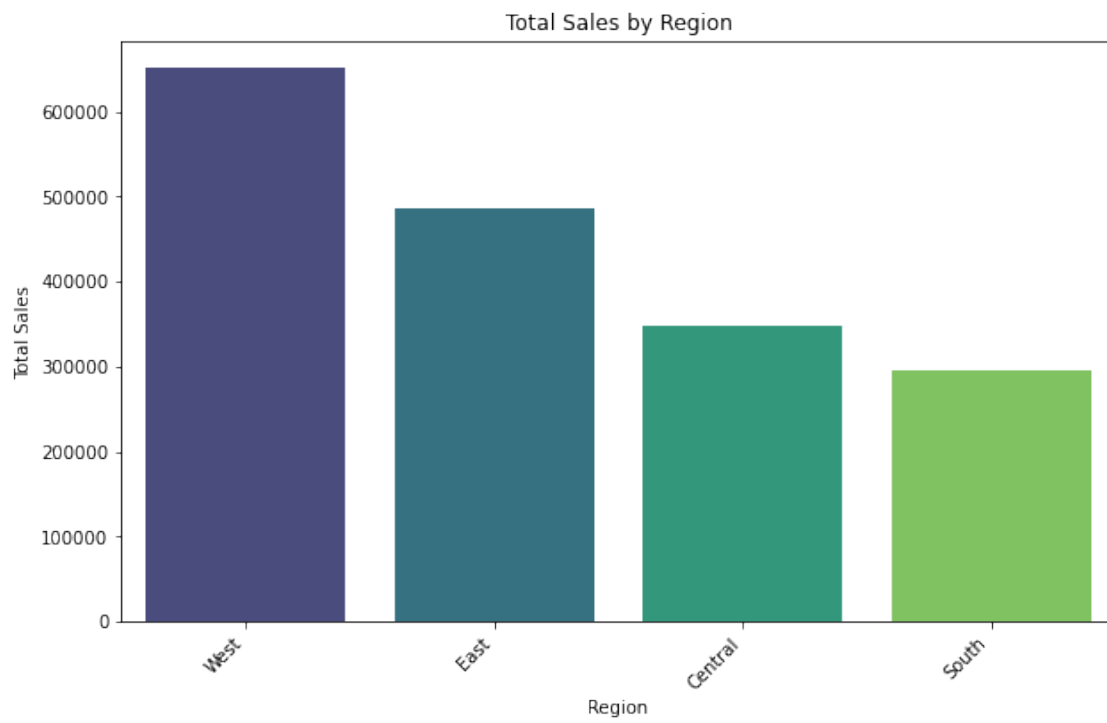
14 Explore regional variations in sales (e.g., by state or city).

```
[152]: # Regional Analysis - Sales by State
plt.figure(figsize=(24, 8))
sns.barplot(x='State', y='Sales', data=df, ci=None)
plt.title('Sales by State')
plt.xlabel('State')
plt.ylabel('Total Sales')
plt.xticks(rotation=45)
plt.grid(True)
plt.show()
```



```
[153]: region_sales = df.groupby('Region')['Sales'].sum().reset_index()
region_sales_sorted = region_sales.sort_values(by='Sales', ascending=False)
plt.figure(figsize=(10, 6))
sns.barplot(x='Region', y='Sales', data=region_sales_sorted, palette='viridis')
plt.title('Total Sales by Region')
plt.xlabel('Region')
plt.ylabel('Total Sales')
plt.xticks(rotation=45, ha='right')
```

```
plt.show()
```



15 Identify factors contributing to sales differences across regions.

16 Exploring Sales Dynamics Across Regions

17 Total Sales by Region:

17.0.1 Analysis of total sales across different regions reveals that the West region boasts the highest cumulative sales figures, closely followed by the East, Central, and South regions.

18 Sales Breakdown by Region and Segment:

18.0.1 Delving deeper into sales segmentation across regions, it becomes evident that the Consumer segment consistently emerges as the primary contributor to sales figures across all regions. This is followed by the Corporate and Home Office segments, albeit with varying degrees of influence depending on the region.

19 Sales Analysis by Region and Category:

19.0.1 An insightful examination of sales performance by product category reveals that Furniture sales exhibit prominence across all regions, with particularly robust figures observed in the West and East regions.

19.0.2 Additionally, Office Supplies and Technology categories emerge as significant contributors to overall sales, displaying noteworthy traction across regions.

20 Monthly Sales Trends Across Regions:

20.0.1 A comprehensive review of monthly sales data uncovers nuanced variations in sales patterns across different regions, hinting at potential seasonal trends or divergent demand dynamics throughout the year.

20.0.2 For instance, the West region showcases a distinct surge in sales during December, suggesting heightened consumer activity aligned with the holiday season festivities.

21 Identifying Factors Driving Regional Sales Disparities:

21.0.1 Regional Economic Landscape and Consumer Preferences: Variances in regional economic conditions and consumer preferences play a pivotal role in shaping sales disparities. Factors such as income levels, employment rates, and cultural inclinations influence consumer spending patterns across regions.

21.0.2 Population Dynamics and Demographics: Population density, age distribution, and demographic composition contribute to regional sales disparities, with densely populated regions often exhibiting heightened consumer activity.

21.0.3 Tailored Marketing and Distribution Strategies: The efficacy of region-specific marketing campaigns and distribution strategies tailored to local preferences and needs significantly impacts sales performance across regions.

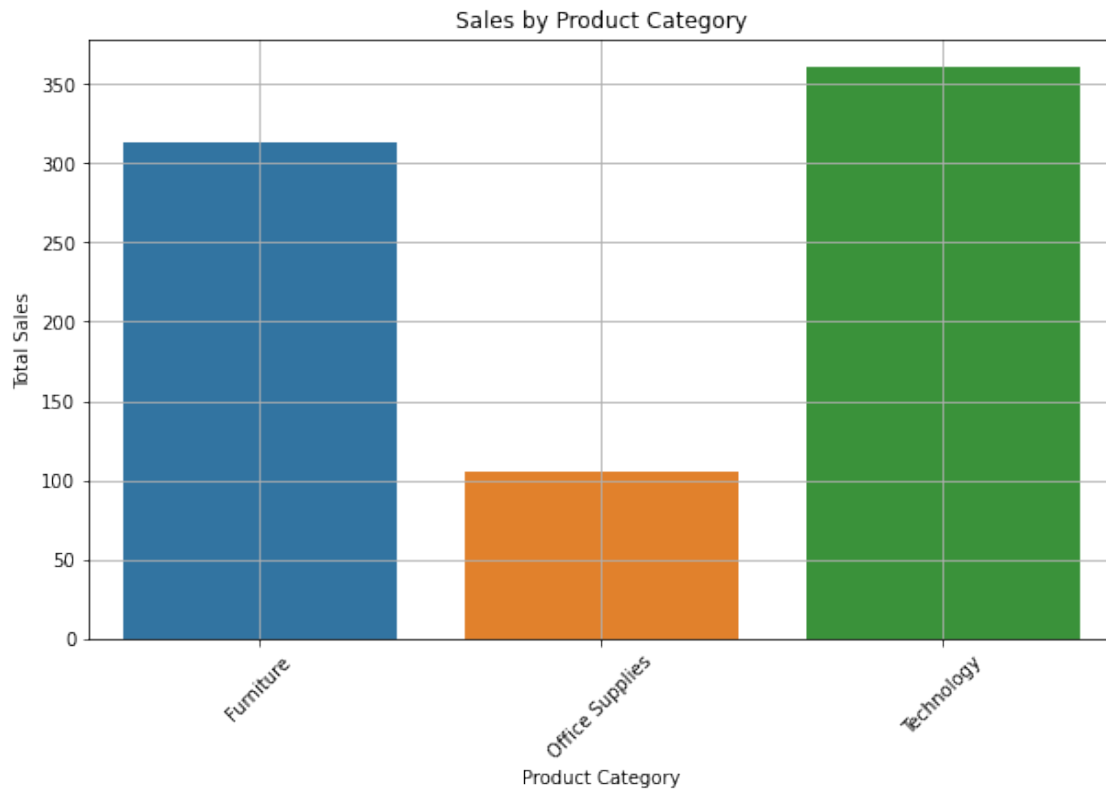
21.0.4 Seasonal Dynamics and Cultural Influences: Seasonal fluctuations and cultural events exert a profound influence on consumer behavior, leading to fluctuations in sales volumes across regions. Understanding and adapting to these seasonal dynamics are essential for optimizing regional sales strategies.

```
[154]: monthly_sales = df.groupby(['order_year', 'order_month'])['Sales'].sum().
        ↪reset_index()
plt.figure(figsize=(12, 6))
sns.lineplot(x='order_month', y='Sales', hue='order_year', data=monthly_sales)
plt.title('Monthly Sales Over the Years')
plt.xlabel('Month')
plt.ylabel('Total Sales')
plt.legend(title='Year')
plt.grid(True)
plt.show()
```



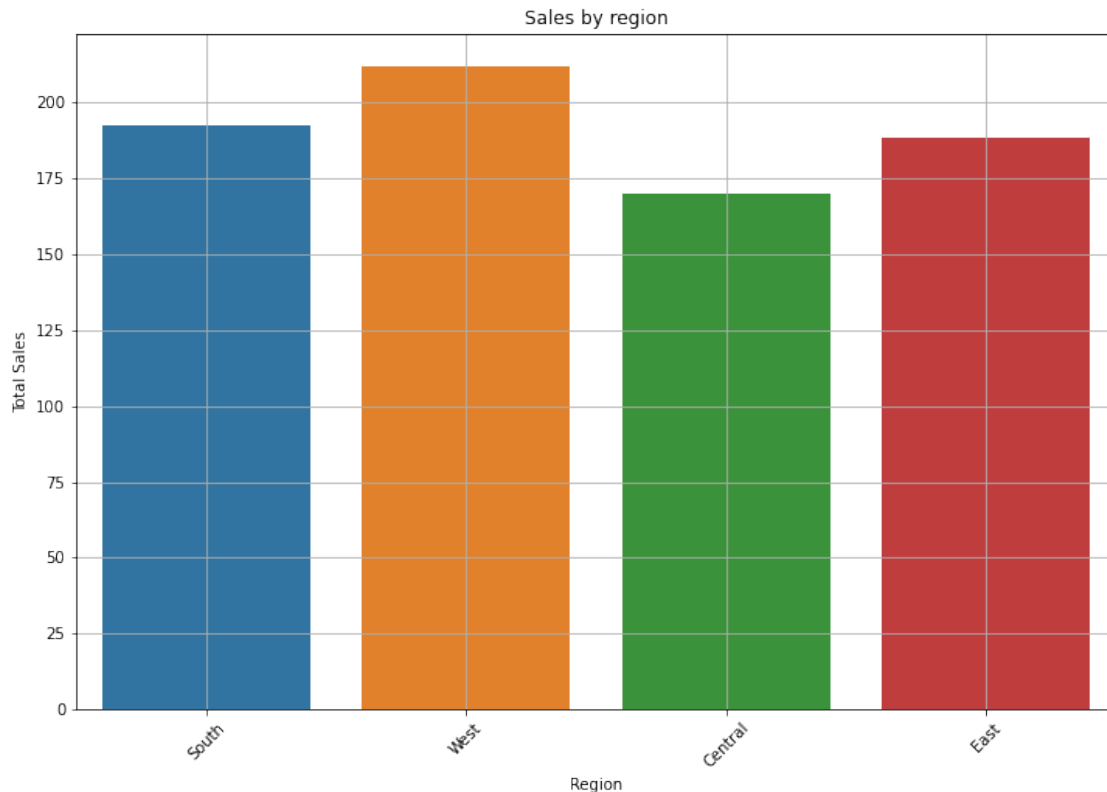
24 Create visualizations to represent product category insights

```
[155]: plt.figure(figsize=(10, 6))
sns.barplot(x='Category', y='Sales', data=df, ci=None)
plt.title('Sales by Product Category')
plt.xlabel('Product Category')
plt.ylabel('Total Sales')
plt.xticks(rotation=45)
plt.grid(True)
plt.show()
```



25 Create visualizations to represent regional variations.

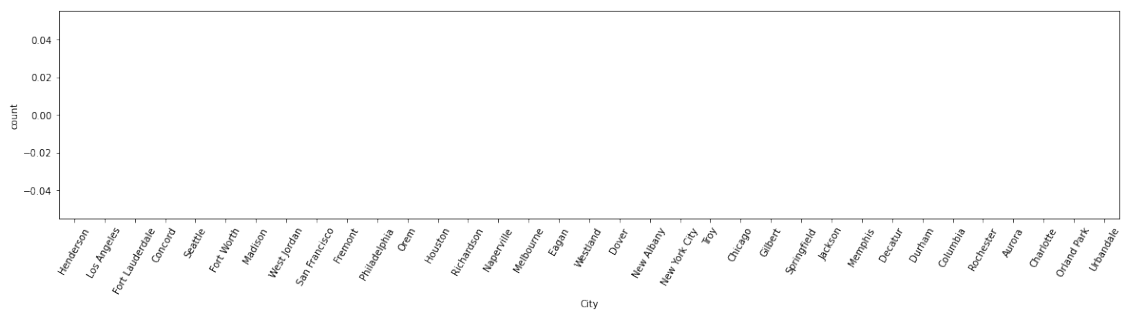
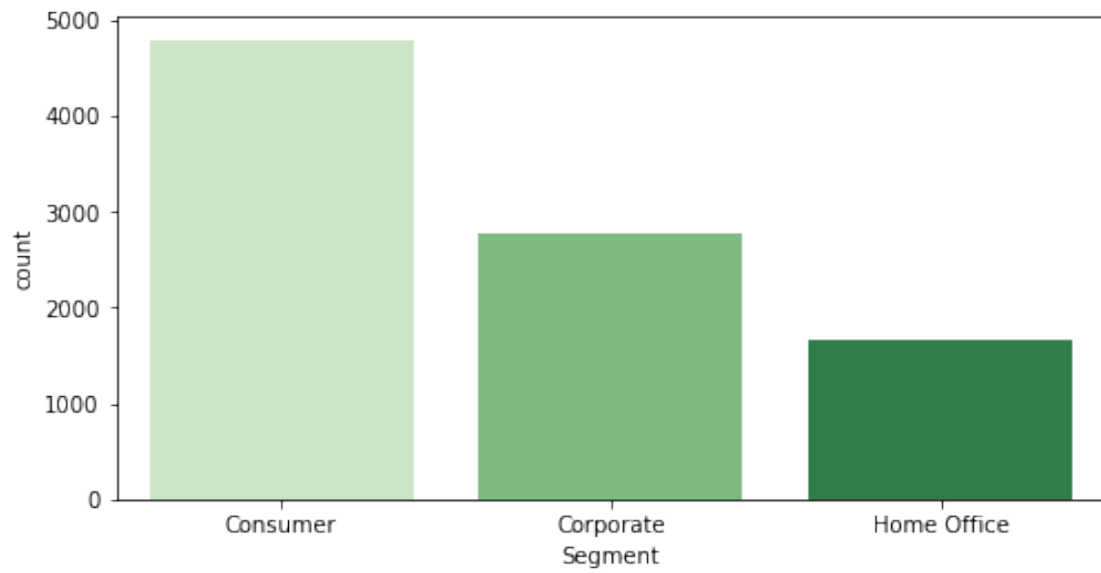
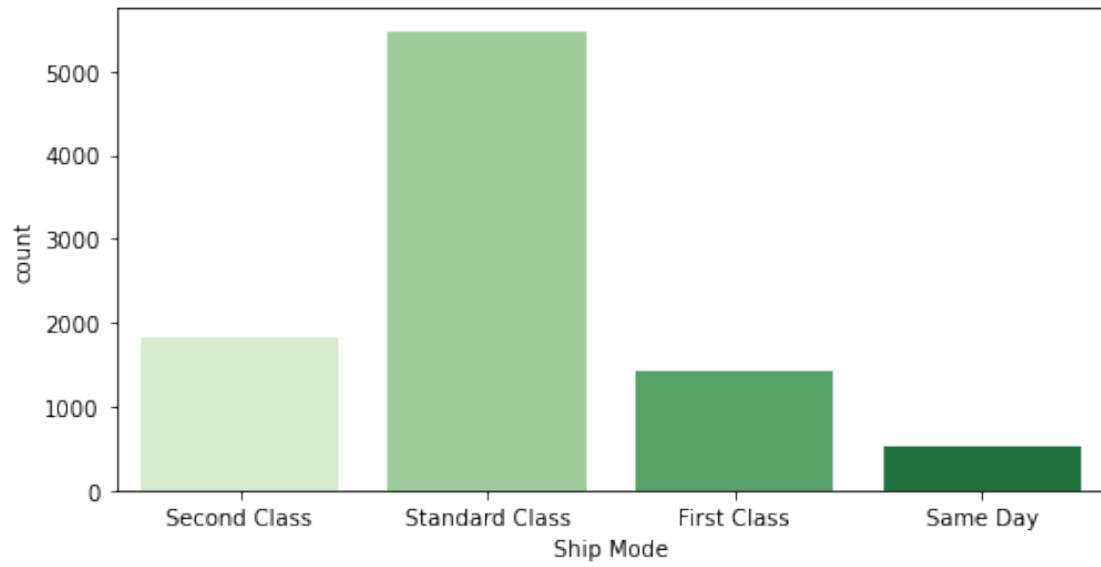
```
[156]: plt.figure(figsize=(12, 8))
sns.barplot(x='Region', y='Sales', data=df, ci=None)
plt.title('Sales by region')
plt.xlabel('Region')
plt.ylabel('Total Sales')
plt.xticks(rotation=45)
plt.grid(True)
plt.show()
```

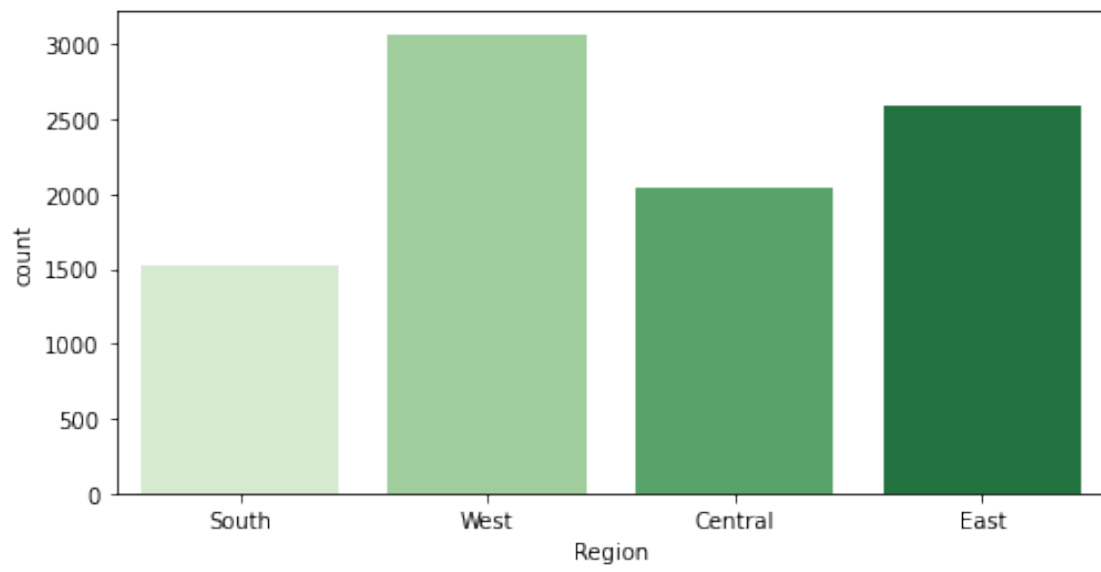
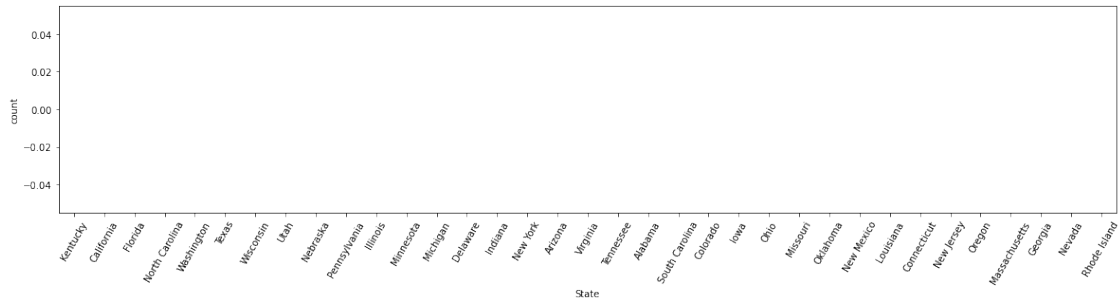


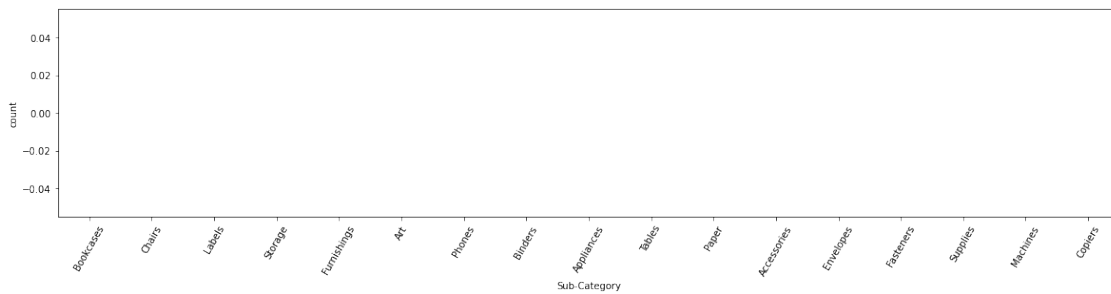
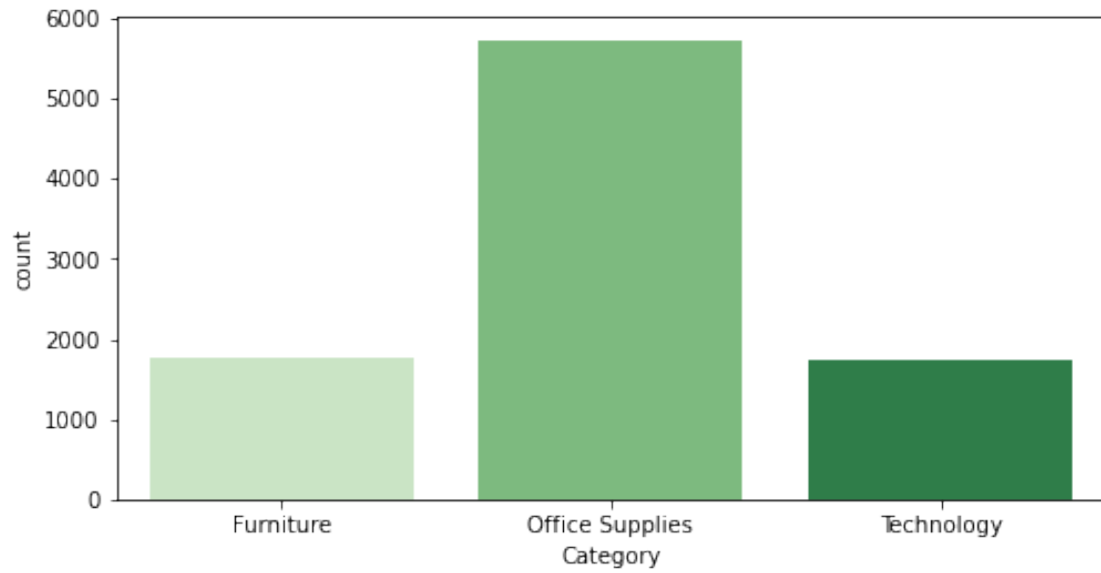
26 Use visualizations to communicate findings effectively.

```
[157]: imp_categorical_features=df[['Ship Mode','Segment','Region','Category']]

for col in imp_categorical_features:
    width=8
    n=len(df[col].unique())
    if n > 15:
        width=20
        plt.figure(figsize=(width,4))
        temp=pd.DataFrame(df[col].value_counts(), index=df[col].
        ↪unique(),columns=['count']).sort_values('count', ascending=False)[:35]
        sns.barplot(data=temp,x=temp.index,y='count',palette= 'Blues')
        plt.ylabel("count")
        plt.xlabel(col)
        plt.xticks(rotation=60)
    else:
        plt.figure(figsize=(width,4))
        sns.countplot(data=df,x=col,palette='Greens')
plt.show()
```







27 Scaling

```
[158]: #Scaling all numeric columns for efficiency,performance,normalization etc
from sklearn.preprocessing import StandardScaler, MinMaxScaler,RobustScaler

numeric_features = ["Quantity","Discount","order_year", "order_month",
                    ↪"order_date","Ship_year", "Ship_month", "Ship_date"]
RobustScaler = RobustScaler()
df[numeric_features] = RobustScaler.fit_transform(df[numeric_features])
```

28 One Hot Encoding(Nominal Data)

```
[159]: #donce encoding to label data(machine can understand only numbers)
df=pd.get_dummies(data=df,columns=['Ship_
↳Mode', 'Segment', 'City', 'State', 'Region', 'Category', 'Sub-Category', 'Country', 'Product_
↳Name'])
df.columns
```

```
[159]: Index(['Sales', 'Quantity', 'Discount', 'Profit', 'order_year', 'order_month',
'order_date', 'Ship_year', 'Ship_month', 'Ship_date',
...
'Product Name_Zebra ZM400 Thermal Label Printer',
'Product Name_Zebra Zazzle Fluorescent Highlighters',
'Product Name_Zipper Ring Binder Pockets',
'Product Name_i.Sound Portable Power - 8000 mAh',
'Product Name_iHome FM Clock Radio with Lightning Dock',
'Product Name_iKross Bluetooth Portable Keyboard + Cell Phone Stand
Holder + Brush for Apple iPhone 5S 5C 5, 4S 4',
'Product Name_iOttie HLCRI0102 Car Mount',
'Product Name_iOttie XL Car Mount',
'Product Name_invisibleSHIELD by ZAGG Smudge-Free Screen Protector',
'Product Name_netTALK DUO VoIP Telephone Service'],
dtype='object', length=2440)
```

```
[160]: # drop columns Preventing Multicollinearity
df.drop(columns=["Product Name_Bush Somerset Collection_
↳Bookcase", "City_Henderson", "Ship Mode_First_
↳Class", "Segment_Consumer", "State_Alabama", "Region_Central", "Category_Furniture", "Sub-Catego
```

```
[161]: df.sample(10) #random 5 samples
```

```
[161]:
```

	Sales	Quantity	Discount	Profit	order_year	order_month	\
8379	182.220	0.000000	0.0	45.5550	0.5	0.500000	
6929	227.460	1.000000	0.0	65.9634	0.5	-1.166667	
5207	49.120	0.333333	0.0	23.0864	0.0	0.500000	
2852	49.650	0.666667	0.0	20.8530	-1.0	-1.000000	
2039	23.680	-0.333333	1.0	8.8800	-1.0	0.333333	
37	113.328	1.666667	1.0	35.4150	-0.5	0.500000	
1005	99.990	-0.666667	0.0	37.9962	-1.0	-0.666667	
5227	27.920	0.333333	0.0	8.0968	0.5	-1.333333	
1780	31.320	1.666667	2.5	-25.0560	0.0	0.500000	
8283	40.480	-0.333333	0.0	17.4064	0.5	0.166667	

	order_date	Ship_year	Ship_month	Ship_date	Transformed_Profit	...	\
8379	-0.733333	0.5	0.500000	-0.3750	3.818920	...	
6929	-0.600000	0.5	-1.166667	-0.2500	4.189100	...	
5207	-0.200000	0.0	0.500000	0.1250	3.139244	...	

2852	1.000000	-1.0	-0.833333	-0.8125	3.037498	...
2039	0.533333	-1.0	0.333333	0.8125	2.183802	...
37	0.800000	-0.5	0.500000	0.9375	3.567135	...
1005	0.800000	-1.0	-0.500000	-0.9375	3.637486	...
5227	0.066667	0.5	-1.333333	0.0000	2.091469	...
1780	0.466667	0.0	0.500000	0.7500	-3.221113	...
8283	-0.200000	0.5	0.166667	0.0000	2.856838	...

Product Name_Zebra GX420t Direct Thermal/Thermal Transfer Printer \

8379	0
6929	0
5207	0
2852	0
2039	0
37	0
1005	0
5227	0
1780	0
8283	0

Product Name_Zebra ZM400 Thermal Label Printer \

8379	0
6929	0
5207	0
2852	0
2039	0
37	0
1005	0
5227	0
1780	0
8283	0

Product Name_Zebra Zazzle Fluorescent Highlighters \

8379	0
6929	0
5207	0
2852	0
2039	0
37	0
1005	0
5227	0
1780	0
8283	0

Product Name_Zipper Ring Binder Pockets \

8379	0
6929	0

5207	0
2852	0
2039	0
37	0
1005	0
5227	0
1780	0
8283	0

Product Name_i.Sound Portable Power - 8000 mAh \

8379	0
6929	0
5207	0
2852	0
2039	0
37	0
1005	0
5227	0
1780	0
8283	0

Product Name_iHome FM Clock Radio with Lightning Dock \

8379	0
6929	0
5207	0
2852	0
2039	0
37	0
1005	0
5227	0
1780	0
8283	0

Product Name_iKross Bluetooth Portable Keyboard + Cell Phone Stand Holder
+ Brush for Apple iPhone 5S 5C 5, 4S 4 \

8379	0
6929	0
5207	0
2852	0
2039	0
37	0
1005	0
5227	0
1780	0
8283	0

Product Name_iOttie HLCRIO102 Car Mount \

8379	0
6929	0
5207	0
2852	0
2039	0
37	0
1005	0
5227	0
1780	0
8283	0

Product Name_iOttie XL Car Mount \

8379	0
6929	0
5207	0
2852	0
2039	0
37	0
1005	0
5227	0
1780	0
8283	0

Product Name_invisibleSHIELD by ZAGG Smudge-Free Screen Protector \

8379	0
6929	0
5207	0
2852	0
2039	0
37	0
1005	0
5227	0
1780	0
8283	0

Product Name_netTALK DUO VoIP Telephone Service

8379	0
6929	0
5207	0
2852	0
2039	0
37	0
1005	0
5227	0
1780	0
8283	0

[10 rows x 2432 columns]

```
[162]: df.dtypes #datatypes of all columns
```

```
[162]: Sales
float64
Quantity
float64
Discount
float64
Profit
float64
order_year
float64
...
Product Name_iKross Bluetooth Portable Keyboard + Cell Phone Stand Holder +
Brush for Apple iPhone 5S 5C 5, 4S 4      uint8
Product Name_iOttie HLCRI0102 Car Mount
uint8
Product Name_iOttie XL Car Mount
uint8
Product Name_invisibleSHIELD by ZAGG Smudge-Free Screen Protector
uint8
Product Name_netTALK DUO VoIP Telephone Service
uint8
Length: 2432, dtype: object
```

29 Splitting

```
[163]: #Splitting the data into dependent variable(y) and independent variables(X)
X=df.drop(columns=["Profit","Sales","Transformed_Sales"],axis=1)
y=df["Transformed_Sales"]
```

```
[164]: X
```

```
[164]:
```

	Quantity	Discount	order_year	order_month	order_date	Ship_year	\
0	-0.333333	0.0	0.0	0.333333	-0.466667	0.0	
1	0.000000	0.0	0.0	0.333333	-0.466667	0.0	
2	-0.333333	0.0	0.0	-0.500000	-0.200000	0.0	
4	-0.333333	1.0	-0.5	0.166667	-0.266667	-0.5	
5	1.333333	0.0	-1.0	-0.500000	-0.400000	-1.0	
...	
9867	0.333333	0.0	0.0	0.000000	-0.933333	0.0	
9868	-0.333333	0.0	0.0	0.000000	-0.933333	0.0	
9869	0.666667	0.0	0.0	0.000000	-0.933333	0.0	
9870	0.666667	1.0	0.5	0.333333	-0.266667	0.5	

9871	0.000000	0.0	-0.5	0.333333	-0.466667	-0.5
------	----------	-----	------	----------	-----------	------

	Ship_month	Ship_date	Transformed_Profit	Ship Mode_Same Day	\
0	0.333333	-0.3125	3.735610		0
1	0.333333	-0.3125	5.391726		0
2	-0.500000	0.0000	1.927368		0
4	0.166667	0.1250	0.922829		0
5	-0.500000	-0.1250	2.651085		0
...	
9867	0.000000	-0.6875	2.671718		0
9868	0.000000	-0.6875	1.095273		0
9869	0.000000	-0.6875	2.539395		0
9870	0.333333	0.0000	4.869751		0
9871	0.333333	-0.1875	2.926414		0

	Ship Mode_Second Class	...	\
0		1	...
1		1	...
2		1	...
4		0	...
5		0	...
...	
9867		0	...
9868		0	...
9869		0	...
9870		1	...
9871		1	...

	Product Name_Zebra GX420t Direct Thermal/Thermal Transfer Printer	\
0		0
1		0
2		0
4		0
5		0
...	...	
9867		0
9868		0
9869		0
9870		0
9871		0

	Product Name_Zebra ZM400 Thermal Label Printer	\
0		0
1		0
2		0
4		0
5		0

...	...	
9867		0
9868		0
9869		0
9870		0
9871		0

	Product Name_Zebra Zazzle Fluorescent Highlighters \	
0		0
1		0
2		0
4		0
5		0
...	...	
9867		0
9868		0
9869		0
9870		0
9871		0

	Product Name_Zipper Ring Binder Pockets \	
0		0
1		0
2		0
4		0
5		0
...	...	
9867		0
9868		0
9869		0
9870		0
9871		0

	Product Name_i.Sound Portable Power - 8000 mAh \	
0		0
1		0
2		0
4		0
5		0
...	...	
9867		0
9868		0
9869		0
9870		0
9871		0

	Product Name_iHome FM Clock Radio with Lightning Dock \
--	---

0	0
1	0
2	0
4	0
5	0
...	...
9867	0
9868	0
9869	0
9870	0
9871	0

Product Name_iKross Bluetooth Portable Keyboard + Cell Phone Stand Holder
+ Brush for Apple iPhone 5S 5C 5, 4S 4 \

0	0
1	0
2	0
4	0
5	0
...	...
9867	0
9868	0
9869	0
9870	0
9871	0

Product Name_iOttie HLCRI0102 Car Mount \

0	0
1	0
2	0
4	0
5	0
...	...
9867	0
9868	0
9869	0
9870	0
9871	0

Product Name_iOttie XL Car Mount \

0	0
1	0
2	0
4	0
5	0
...	...
9867	0

9868	0
9869	0
9870	0
9871	0

	Product Name_invisibleSHIELD by ZAGG Smudge-Free Screen Protector \
0	0
1	0
2	0
4	0
5	0
...	...
9867	0
9868	0
9869	0
9870	0
9871	0

	Product Name_netTALK DUO VoIP Telephone Service
0	0
1	0
2	0
4	0
5	0
...	...
9867	0
9868	0
9869	0
9870	0
9871	0

[9231 rows x 2429 columns]

30 Feature Selection

```
[165]: best_features = SelectKBest(score_func=f_regression, k=10)
fit = best_features.fit(X, y)
```

```
C:\Users\ganesh\anaconda3\lib\site-
packages\sklearn\feature_selection\_univariate_selection.py:289: RuntimeWarning:
divide by zero encountered in true_divide
    correlation_coefficient /= X_norms
C:\Users\ganesh\anaconda3\lib\site-
packages\sklearn\feature_selection\_univariate_selection.py:358: RuntimeWarning:
invalid value encountered in true_divide
    f_statistic = corr_coef_squared / (1 - corr_coef_squared) * deg_of_freedom
```

```
[166]: dfscores=pd.DataFrame(fit.scores_)
dfcolumns=pd.DataFrame(X.columns)
```

```
[167]: featureScores=pd.concat([dfcolumns,dfscores],axis=1)
featureScores.columns=["Specs","Score"]
```

```
[168]: featureScores
```

```
[168]:
```

	Specs	Score
0	Quantity	1108.203578
1	Discount	491.504260
2	order_year	0.029984
3	order_month	1.162049
4	order_date	0.044602
...
2424	Product Name_iKross Bluetooth Portable Keyboar...	0.207501
2425	Product Name_iOttie HLCRIO102 Car Mount	0.609101
2426	Product Name_iOttie XL Car Mount	0.354925
2427	Product Name_invisibleSHIELD by ZAGG Smudge-Fr...	0.017594
2428	Product Name_netTALK DUO VoIP Telephone Service	2.372775

[2429 rows x 2 columns]

```
[169]: print(featureScores.nlargest(55,"Score"))
```

	Specs	Score
8	Transformed_Profit	5061.240835
589	Category_Office Supplies	2307.897944
0	Quantity	1108.203578
590	Category_Technology	1026.866451
595	Sub-Category_Chairs	784.610347
593	Sub-Category_Binders	669.970751
603	Sub-Category_Phones	526.775250
1	Discount	491.504260
592	Sub-Category_Art	449.609784
606	Sub-Category_Tables	263.951495
598	Sub-Category_Fasteners	248.071964
594	Sub-Category_Bookcases	221.948901
596	Sub-Category_Copiers	203.365874
600	Sub-Category_Labels	201.657948
604	Sub-Category_Storage	186.114097
602	Sub-Category_Paper	185.667940
601	Sub-Category_Machines	152.133133
578	State_Texas	105.131382
573	State_Pennsylvania	62.265101
381	City_Philadelphia	57.783160
548	State_Illinois	51.301009

591	Sub-Category_Appliances	49.581887
588	Region_West	43.703354
540	State_California	42.470469
92	City_Chicago	39.925238
2210	Product Name_Wilson Jones Easy Flow II Sheet L...	39.046504
1243	Product Name_Fellowes PB500 Electric Punch Pla...	38.661972
218	City_Houston	37.970459
2081	Product Name_Staples	35.995389
1945	Product Name_SAFCO Arco Folding Chair	35.892068
1279	Product Name_GBC DocuBind TL300 Electric Bindi...	35.437552
570	State_Ohio	35.106957
1395	Product Name_Hewlett Packard LaserJet 3310 Copier	32.785323
1360	Product Name_Global Troy Executive Leather Low...	32.779995
2130	Product Name_Tennsco 6- and 18-Compartment Loc...	32.463255
1427	Product Name_Hon Deluxe Fabric Upholstered Sta...	32.005188
1447	Product Name_Honeywell Enviracaire Portable HE...	30.940444
1852	Product Name_Plantronics CS510 - Over-the-Head...	30.419267
867	Product Name_Avery Non-Stick Binders	30.276044
2421	Product Name_Zipper Ring Binder Pockets	29.927885
1970	Product Name_Samsung Galaxy Mega 6.3	29.919507
1377	Product Name_HON 5400 Series Task Chairs for B...	29.040657
605	Sub-Category_Supplies	27.857567
1611	Product Name_Logitech P710e Mobile Speakerphone	27.245337
1859	Product Name_Plantronics Savi W720 Multi-Devic...	26.991260
874	Product Name_Avery Reinforcements for Hole-Pun...	26.177325
1394	Product Name_Hewlett Packard 610 Color Digital...	25.964879
1376	Product Name_GuestStacker Chair with Chrome Fi...	25.746257
720	Product Name_Adjustable Depth Letter/Legal Cart	25.739577
1460	Product Name_Hot File 7-Pocket, Floor Stand	25.221180
689	Product Name_Acco Suede Grain Vinyl Round Ring...	24.974089
844	Product Name_Avery Durable Binders	24.892038
1795	Product Name_Office Star - Professional Matrix...	24.598514
1335	Product Name_Global Deluxe High-Back Manager's...	24.473559
1867	Product Name_Poly String Tie Envelopes	24.317337

31 Data Splitting

```
[170]: #splitting the data into train and test
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
```

32 Training Model

```
[171]: #Training model on RandomForestRegressor

model = RandomForestRegressor(random_state=42)
# Train the model
model.fit(X_train, y_train)

# Make predictions
y_pred = model.predict(X_test)

# Calculate evaluation metrics
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

# Print evaluation metrics
print("Mean Squared Error:", mse)
print("R-squared:", r2)
```

Mean Squared Error: 0.14796312980709456
R-squared: 0.9434364088363466

```
[172]: #Training model on KNeighborsRegressor

modelk = KNeighborsRegressor()

#model fit and train
modelk.fit(X_train, y_train)

# Make predictions
y_pred = modelk.predict(X_test)

# Calculate evaluation metrics
msek = mean_squared_error(y_test, y_pred)
r2k = r2_score(y_test, y_pred)

# Print evaluation metrics
print("Mean Squared Error:", msek)
print("R-squared:", r2k)
```

Mean Squared Error: 0.4503792650396285
R-squared: 0.8278282660720876

```
[173]: #Training model on MLPRegressor

model = MLPRegressor(random_state=42)
```

```

#model fit and train
model.fit(X_train, y_train)

# Make predictions
y_pred = model.predict(X_test)

# Calculate evaluation metrics
mse = mean_squared_error(y_test, y_pred)
r2m = r2_score(y_test, y_pred)
mae = mean_absolute_error(y_test, y_pred)
rms=np.sqrt(mse)

print("R-squared:", r2m)

```

R-squared: 0.9836333726161193

33 Evaluate the model's performance using appropriate metrics for MLPRegressor

```

[174]: print("rmse :",rms)
print("Mean Squared Error:", mse)
print("R-squared:", r2m)
print("Mean Absolute Error:", mae)

```

```

rmse : 0.38465975849716144
Mean Squared Error: 0.042813006782030524
R-squared: 0.9836333726161193
Mean Absolute Error: 0.12647042964582925

```

34 Model Prediction

```

[175]: #model Prediction
ll=model.predict(X_test)
ll[:20]

```

```

[175]: array([2.91184998, 8.43196963, 5.01888031, 7.34349144, 4.04138845,
            4.10060321, 3.11776551, 2.93605332, 5.70730388, 5.83653792,
            2.48335608, 2.40176275, 3.56379378, 3.37635933, 4.44597523,
            3.35233614, 5.99622523, 4.62667132, 3.18045318, 4.78959973])

```

```

[176]: #Actual Data
y_test[:20]

```

```

[176]: 1263    2.874242
      1111    7.716229
      9861    5.042780

```

```

977      7.144825
9090     4.065945
1998     4.051855
82       3.062924
4883     2.967333
678      5.691609
7696     5.727291
8996     2.329422
9020     2.474856
6219     3.540959
6609     3.395850
3717     4.472781
9780     3.108436
2274     5.966070
52       4.499699
6624     3.726946
3502     4.663250
Name: Transformed_Sales, dtype: float64

```

```

[177]: dfpre = pd.DataFrame({'y_pred': ll, 'y_test': y_test})
dfpre.head(30)

```

```

[177]:
   y_pred  y_test
1263  2.911850  2.874242
1111  8.431970  7.716229
9861  5.018880  5.042780
977   7.343491  7.144825
9090  4.041388  4.065945
1998  4.100603  4.051855
82    3.117766  3.062924
4883  2.936053  2.967333
678   5.707304  5.691609
7696  5.836538  5.727291
8996  2.483356  2.329422
9020  2.401763  2.474856
6219  3.563794  3.540959
6609  3.376359  3.395850
3717  4.445975  4.472781
9780  3.352336  3.108436
2274  5.996225  5.966070
52    4.626671  4.499699
6624  3.180453  3.726946
3502  4.789600  4.663250
6010  2.537331  2.501518
7490  5.472728  5.529207
1741  3.655699  3.497719
7297  3.136132  3.061520

```


7446	2.259446	2.232163
4960	3.820677	3.866816
5118	5.644951	5.710295
1515	3.613258	3.565355
3928	2.841434	2.601207
4994	4.931324	5.297497

```
[178]: #back transformation on predicted data
count=0
a=[]
for i in ll:
    if count == 1847:
        break
    else:
        print(round(np.exp(i),2))
        a.append(round(np.exp(i),2))
        count += 1
```

18.39
4591.53
151.24
1546.1
56.91
60.38
22.6
18.84
301.06
342.59
11.98
11.04
35.3
29.26
85.28
28.57
401.91
102.17
24.06
120.25
12.65
238.11
38.69
23.01
9.58
45.64
282.86
37.09
17.14

138.56
87.65
7.03
61.31
184.53
19.06
9.3
20.92
46.99
34.31
7.08
39.76
1699.09
552.45
71.66
18.72
52.74
19.22
12.22
1054.54
270.43
18.27
6.11
81.85
42.19
54.92
14.23
9.8
61.96
57.62
43.3
1213.88
189.15
345.21
217.26
35.85
15.18
1083.33
168.02
14.85
8.69
24.85
56.2
91.78
3.44
14.61
98.36
322.04

2.46
43.98
119.53
12.22
2.49
7.95
148.23
63.72
7.53
49.97
63.01
18.47
43.17
687.87
132.2
172.78
507.98
574.43
10.15
82.07
10.38
1612.02
12.49
19.15
60.59
20.8
146.34
16.05
5.01
12.89
64.08
9.51
27.45
12.07
42.51
12.08
72.93
20.06
6.02
10.37
20.58
13.52
10.42
352.81
10.7
41.56
21.75
5.41

21.33
14.6
818.78
227.48
30.89
3.54
86.76
571.83
2.36
17.09
1600.77
22.7
60.13
480.73
237.42
1436.44
66.49
118.62
27.74
18.88
45.68
122.3
24.58
115.27
2.04
42.29
308.96
161.42
17.57
21.38
7.05
33.3
12.43
45.01
6.26
15.73
112.87
100.69
578.51
110.44
32.2
12.93
20.11
33.54
398.39
86.1
55.22
1555.48

114.13
39.83
5.12
16.52
44.84
6.87
8.28
12.17
172.88
1958.13
50.27
9.93
150.36
219.23
13.86
505.08
28.82
170.65
3.41
45.65
213.97
13.35
28.13
1001.2
49.27
3.59
265.04
5.9
25.33
13.15
426.67
269.34
300.81
7.12
61.44
6.96
3.48
73.65
6.87
45.36
120.86
688.14
316.65
19.27
21.13
26.23
8.05
1465.98

3.45
11.56
265.1
163.7
16.63
24.69
24.37
262.54
84.83
32.43
9.31
129.41
3.12
6.45
5.61
101.37
71.37
13.26
264.4
421.47
3.22
11.78
45.42
154.45
451.25
723.53
175.43
10.14
19.03
89.27
142.56
21.21
2.03
371.62
34.77
604.73
8.45
18.99
18.8
46.56
10.61
10.19
81.3
5.19
124.45
18.72
66.14
2.79

3.23
14.11
84.0
24.44
53.93
8.03
3.94
7.8
19.35
180.4
82.43
42.28
36.75
8.41
127.67
130.13
17.57
133.72
595.74
379.29
30.41
15.76
3.42
20.04
1688.96
164.44
50.44
55.27
38.26
3.19
2.89
51.89
1767.2
3.07
2.71
58.34
13.89
41.86
19.66
30.79
1.85
3.14
7.09
64.24
12.55
17.9
40.62
228.27

53.88
25.08
212.98
188.3
48.89
1099.55
188.35
4.67
91.41
36.85
17.19
26.0
7.6
1861.38
17.72
68.14
529.97
340.33
28.67
264.1
62.96
27.71
160.43
81.35
587.1
7.22
8.64
480.92
50.42
552.56
511.37
296.76
7.98
10.25
168.81
7.51
10.75
4.37
171.63
46.91
20.73
355.52
849.65
95.1
268.12
109.64
81.65
249.74

6.16
174.64
92.56
68.72
4.86
3269.83
12.28
86.39
51.43
49.73
7.2
357.65
89.89
12.23
524.11
12.15
79.16
15.06
9.4
53.81
9.72
13.43
271.22
64.96
365.72
758.82
31.03
77.39
31.54
12.17
7.11
16.47
128.62
270.9
567.91
7.99
533.84
57.11
11.85
217.35
35.74
4.81
43.13
1050.42
5.87
21.12
1891.19
16.58

7.14
143.21
11.08
33.34
359.0
20.19
17.15
11.11
39.21
421.95
43.16
27.45
8.79
2.3
438.52
2.74
179.96
543.55
242.0
19.09
3640.35
28.8
151.86
136.02
150.42
116.2
16.7
19.36
4.8
2.79
20.03
1803.78
25.48
94.86
1386.51
14.66
230.88
14.9
250.41
859.46
49.13
42.36
38.28
253.57
104.21
45.09
61.66
1.29

7.87
158.0
21.25
64.15
34.72
73.06
9.37
63.33
435.27
7.85
32.13
24.87
489.24
229.82
86.42
31.96
288.07
14.35
126.85
84.83
12.66
694.97
172.94
54.05
6.84
204.45
8.08
495.29
100.67
6.76
8.48
714.85
7.4
1.7
8.64
74.76
61.88
4.58
269.74
1052.5
5.37
31.56
21.68
82.52
7.31
60.07
9.95
389.74

39.6
9.73
7.45
974.04
816.44
760.01
17.68
9.5
62.63
8.5
161.85
973.69
106.46
18.62
58.64
35.1
86.25
27.59
159.81
54.0
8.6
64.48
30.54
10.23
17.36
4.77
48.18
2.66
40.78
20.25
33.21
8.96
25.55
11.08
2.06
1.59
42.85
197.25
155.26
912.19
43.36
486.06
13.2
15.46
87.01
62.37
592.83
50.3

5.96
104.85
324.79
604.22
74.34
336.67
3.83
8.97
13.79
665.05
987.77
46.88
102.05
310.11
167.09
106.72
124.4
62.43
26.75
63.78
977.09
152.5
71.68
13.63
42.48
11.89
16.85
211.97
1275.51
212.51
120.55
14.47
168.43
67.92
26.72
81.74
225.3
183.6
209.72
943.16
550.85
66.78
12.58
178.44
6.1
42.39
9.18
12.69

87.76
30.63
249.34
428.35
14.56
105.34
64.12
7.87
29.86
671.55
25.17
65.98
196.59
88.02
218.45
12.44
637.85
32.94
123.62
32.48
5.91
7.5
6.84
7.99
80.76
28.44
146.12
5.3
13.71
37.63
992.96
4.9
35.97
29.63
1298.79
121.8
10.72
14.62
1639.14
121.79
8.95
20.97
17.45
6.13
20.8
67.51
28.42
99.45

56.39
1070.93
74.81
23.68
45.24
37.28
55.31
22.29
451.53
136.08
5.91
44.29
68.77
21.52
4.5
91.08
751.53
95.59
3.87
669.99
36.12
124.3
12.19
14.76
292.59
3.24
22.81
52.84
29.98
15.51
917.64
8.99
325.21
2.73
355.63
25.87
324.06
1389.94
319.84
28.24
61.41
102.62
440.07
19.27
35.33
82.74
5.79
13.99

12.46
36.15
85.9
43.75
6.84
518.31
52.68
148.84
32.09
39.9
10.17
21.35
7.31
29.56
39.46
334.71
22.59
16.03
41.37
405.26
3.71
10.4
92.24
2.98
15.62
32.04
4.61
6.52
1885.24
138.56
446.66
17.58
119.45
34.77
263.36
221.83
84.8
361.65
29.42
39.26
1796.93
7.11
17.88
212.41
9.39
384.98
55.67
80.16

11.4
292.59
42.92
25.9
151.88
9.27
119.09
195.83
3.86
641.05
41.14
42.86
3.36
104.06
19.17
4.45
70.42
1875.66
58.61
6.84
16.82
101.81
10.4
210.65
57.6
91.13
58.99
32.16
5.98
6.06
10.25
23.35
285.96
839.72
181.13
4.81
9.86
22.24
4.66
438.7
8.77
26.33
25.69
66.6
6.8
297.51
176.69
83.63

57.59
89.32
38.07
82.01
108.79
110.21
19.61
21.78
3.66
239.16
178.26
592.55
108.84
494.27
157.61
299.14
64.64
20.73
277.85
315.21
12.02
77.82
2.03
96.59
33.23
154.56
569.31
16.05
11.98
860.74
198.62
225.24
15.23
44.11
47.05
783.49
34.96
54.6
1.73
3.99
14.25
2.69
100.88
28.35
187.14
181.56
253.07
462.08

1103.87
287.96
89.06
93.83
49.86
24.5
200.64
493.86
74.24
26.74
22.62
34.95
159.67
83.39
3.46
336.84
42.48
144.23
46.21
7.12
9.73
27.33
795.17
26.54
8.79
57.02
26.18
57.36
27.24
252.84
43.94
69.77
21.49
382.77
62.68
24.76
16.32
1.68
122.6
33.75
188.93
40.07
8.65
5.14
10.59
14.94
5.05
1876.9

39.51
62.6
13.99
228.16
138.04
34.61
159.43
18.91
264.43
7.74
11.03
15.82
476.08
460.78
4.7
4.16
2.75
11.76
28.37
2.2
879.48
135.78
32.86
11.26
21.18
38.27
712.37
71.83
9.5
151.1
427.96
774.31
27.78
4.91
145.53
1594.92
136.1
313.11
109.35
31.3
27.27
20.26
32.18
3.04
25.4
214.5
1247.81
106.47

3.44
240.68
49.52
446.07
129.89
276.15
265.48
14.25
19.83
202.46
75.68
8.9
17.62
13.6
20.32
5.79
36.37
138.69
111.98
73.35
261.54
6.64
106.48
478.4
134.92
849.81
1582.31
143.09
47.39
92.68
58.71
61.08
194.98
27.27
21.09
6.85
4.62
25.67
21.78
867.05
481.15
242.7
21.2
17.1
48.9
7.37
15.02
306.75

69.26
173.22
14.63
82.12
47.76
26.9
17.13
38.09
3.99
7.45
46.87
5123.06
22.53
1.65
22.71
1408.01
152.8
12.64
335.52
80.92
22.66
26.4
19.26
232.23
388.71
4.78
307.12
50.21
252.09
19.98
818.79
26.96
15.51
8.23
7.1
30.37
156.24
897.83
47.91
12.53
161.73
630.79
14.47
297.56
1046.98
36.47
63.55
84.42

12.54
356.98
133.69
65.7
9.35
10.52
139.79
43.15
11.83
531.32
83.48
24.81
202.17
37.93
243.53
292.7
12.56
20.92
22.17
143.1
54.1
5.82
111.29
117.55
244.62
142.06
13.43
708.18
14.87
23.37
24.62
757.01
16.47
48.72
56.86
14.43
103.38
293.0
328.56
2.81
15.09
972.45
26.31
27.07
11.08
10.21
2671.63
82.45

325.12
185.53
121.43
49.5
5.63
267.45
36.66
100.76
239.7
17.26
321.11
78.14
4.22
4.09
54.43
76.01
39.23
3.13
162.02
39.98
18.78
254.1
16.62
18.42
148.91
80.03
81.88
98.86
599.84
7.72
138.1
13.66
135.92
840.2
10.03
35.81
17.7
29.07
63.64
28.76
1597.34
17.44
156.14
204.92
34.85
32.59
193.24
208.72

55.06
10.18
259.14
1322.97
6.48
12.91
24.04
8.81
24.62
19.5
1.84
124.59
1031.6
13.37
10.79
459.51
101.23
108.73
29.59
270.99
4.44
7.15
22.43
249.32
23.8
113.49
136.08
353.1
266.29
20.98
127.39
172.56
148.11
173.72
73.91
13.82
107.89
28.37
177.68
54.83
27.84
76.67
61.41
211.68
13.9
4.15
306.85
59.97

165.97
157.96
235.41
127.35
5.09
18.6
186.05
5.08
2.4
112.16
153.35
206.14
48.78
7.36
102.81
22.94
15.05
195.96
734.94
91.36
35.69
48.78
23.03
91.33
34.28
16.61
49.61
1174.01
239.1
21.72
15.23
24.95
6.45
27.64
9.08
95.84
8.48
101.68
294.06
297.0
26.65
172.14
452.77
31.93
24.58
34.54
8.81
55.07

12.76
37.35
92.03
36.09
72.56
3.98
75.04
432.77
407.63
39.23
22.78
38.65
364.42
11.31
139.62
6.32
39.67
25.94
175.78
2458.84
11.36
129.42
809.82
42.83
741.7
124.1
97.93
723.12
10.09
30.28
373.19
94.66
31.8
146.87
781.07
141.54
22.28
205.63
65.42
175.62
725.64
34.78
206.29
64.64
7.8
48.71
723.9
347.28

366.86
15.92
155.07
15.9
18.98
585.14
74.12
917.43
133.1
1983.39
199.77
129.64
51.15
56.2
3.47
57.85
386.47
21.0
14.0
261.07
4.77
11.29
118.08
76.68
12.78
35.23
149.67
48.87
16.27
900.64
104.7
14.14
5.61
9.24
32.11
82.19
322.11
165.88
21.64
79.07
60.42
36.79
21.39
388.3
18.2
283.99
150.46
101.9

445.77
41.17
522.16
195.1
90.05
13.05
54.36
210.78
50.4
116.85
9.1
9.94
14.11
56.03
22.13
304.1
11.71
1303.78
71.21
8.56
60.74
7.34
22.64
62.89
643.93
74.78
117.86
413.32
13.01
14.1
14.49
4.2
44.47
9.97
67.85
436.05
33.26
103.34
35.87
18.49
36.13
207.72
238.51
5222.01
46.25
81.87
160.61
880.47

174.49
42.61
4.35
121.84
6.81
75.96
409.11
64.87
52.09
40.77
721.38
1096.95
29.65
32.88
15.8
888.49
130.38
188.99
56.96
75.35
32.1
4.98
44.62
107.39
107.97
3.89
35.88
1191.54
119.37
436.9
36.18
7.26
422.03
105.04
160.63
14.78
28.48
218.62
19.61
41.43
6.76
25.18
7.1
320.5
108.02
149.49
217.51
457.98

165.01
4.16
10.11
65.4
2.84
56.28
95.32
16.48
43.12
443.33
5.72
789.08
44.94
44.5
81.91
90.56
498.16
2.36
86.25
179.26
38.27
540.19
15.95
286.17
5.42
26.3
10.12
77.24
14.23
18.43
5.2
4.28
69.73
1519.9
102.95
7.98
51.4
10.64
86.86
14.77
10.67
46.59
9.96
34.83
16.03
29.63
56.1
50.19

16.79
36.48
96.48
20.06
5.46
3.39
13.65
125.05
274.97
187.31
18.1
229.72
706.79
3518.63
1238.34
2941.16
23.33
12.66
49.07
10.9
18.43
72.55
36.01
44.33
90.54
442.49
12.85
300.47
510.12
161.28
10.25
462.02
56.64
6.96
228.01
21.98
312.92
924.49
28.51
14.82
114.09
38.89
39.3
79.8
75.93
352.91
39.25
4.72

12.88
20.86
351.3
2439.43
95.58
9.45
159.61
5.31
333.67
3.88
339.83
15.92
948.26
13.9
14.14
33.14
87.55
88.09
96.7
28.96
478.47
628.29
28.68
13.01
16.91
48.94
34.71
6.94
286.73
111.22
8.39
242.99
83.85
13.52
385.71
33.7
22.24
70.53
66.25
18.02
51.15
5.47
270.83
295.37
3.97
11.68
10.98
89.98

48.76
421.92
206.97
38.77
33.84
11.31
16.06
110.62
5.02
264.35
6.74
17.82
48.37
12.17
12.36
397.86
82.22
11.71
154.68
11.79
35.08
68.15
276.34
46.5
53.43
521.87
139.64
703.87
7.2
120.4
13.21
16.04
441.98
38.05
9.85
245.12
84.19
21.77
6.33
90.17
178.54
14.43
7.22
178.92
50.01
25.38
635.26
89.39

21.89
889.22
31.4
5.92
814.7
22.98
170.68
15.42
1175.14
36.26
72.76
1552.36
23.89
6.47
3.52
107.36
13.5
738.51
207.18
114.83
192.46
7.46
503.03
76.5
111.48
1684.92
37.71
44.44
8.49
160.56
247.42
244.89
103.99
283.27
568.23
2.67
4.54
31.39
578.24
49.86
13.75
75.52
26.67
227.62
13.54
22.38
8.53
21.05

7.87
19.5
11.14
17.15
144.37
6.87
569.45
5.38
275.12
8.52
11.07
362.64
12.34
18.86
53.68
10.09
269.49
9.22
6.19
322.3
18.75
2.77
80.91
171.65
23.88
4.56
44.01
7.11
125.8
114.4
226.28
598.3
898.69
93.42
152.35
222.94
6.7
28.96
15.57
4.15
20.36
57.73
21.03
17.54
196.98
59.46
195.78
31.44

11.78
24.25
14.27
9.68
5.91
5.08
13.63
404.24
297.22
34.23
20.42
1680.22
205.93
2.1
339.47
669.76
43.65
38.28
50.59
35.17
14.19
196.93
71.86
170.52
63.39
94.57
46.12
12.21
29.11
55.17
49.35
15.29
6.62
10.74
149.92
41.57
41.93
97.8
126.72
5.07
9.99
94.51
8.55
475.73
616.21
44.88
29.38
43.61

1066.3
219.26
101.61
227.22
4.85
74.93
63.29
30.58
167.43
444.4
33.2
56.1
9.84
433.38
925.69
27.63
20.52
135.71
853.55
55.31
7.22
351.77
19.78
4.96
55.74
535.69
41.42
32.86
68.95
1279.2
198.5
14.11
46.39
5.25
33.56
5.77
51.93
17.7
475.33
7.45
64.58
18.51
31.89
45.54
237.79
58.02
773.35
268.23

56.67
130.25
957.17
39.87
112.21
32.08
78.02
169.32
25.47
15.73
45.9
30.24
2.71
171.26
853.03
3.8
10.91
562.21
22.13
59.07
6.55
12.04
42.0
36.96
3.97
652.3
1509.18
122.31
136.95
185.9
204.95
505.91
81.35
13.12
723.7
19.97
2.56
664.02
317.27
48.24
73.12
76.62

```
[179]: #back transformation on test data  
count=0  
b=[]  
for i in y_test:
```

```
if count == 1847:
    break
else:
    print(round(np.exp(i),2))
    b.append(round(np.exp(i),2))
    count += 1
```

17.71
2244.48
154.9
1267.53
58.32
57.5
21.39
19.44
296.37
307.14
10.27
11.88
34.5
29.84
87.6
22.39
389.97
89.99
41.55
105.98
12.2
251.94
33.04
21.36
9.32
47.79
301.96
35.35
13.48
199.84
66.96
7.87
67.9
153.57
19.82
9.08
20.7
53.4
36.78
7.96

33.94
1128.39
681.41
63.98
17.0
50.12
18.97
11.95
1665.62
244.77
18.28
5.87
72.9
41.38
53.25
14.07
7.69
37.75
61.19
44.95
1247.64
169.99
1007.98
239.98
34.85
13.98
1292.94
201.58
15.48
8.74
25.03
62.19
79.14
3.49
11.65
102.34
354.9
1.34
39.72
107.44
11.12
1.34
9.81
146.35
61.4
7.64
50.11
71.98

19.92
46.2
706.86
125.76
174.3
546.66
535.41
10.58
80.88
10.56
1676.88
12.96
19.44
58.24
24.82
141.96
14.52
5.18
11.34
59.52
9.33
19.92
15.14
43.6
12.39
75.96
15.14
6.16
6.36
17.47
13.13
8.8
470.38
7.86
38.34
18.53
5.18
18.84
12.22
899.91
239.5
29.16
3.54
88.75
467.97
1.64
17.18
1573.49

12.54
64.94
636.41
245.88
1799.75
67.92
145.76
27.72
11.95
46.2
124.25
27.92
128.85
1.48
40.03
393.54
166.44
17.9
27.92
11.18
30.0
12.84
42.78
6.03
15.75
106.34
93.02
540.57
109.92
30.02
13.9
20.58
35.91
408.74
71.0
59.97
1194.16
95.84
45.57
5.82
15.23
47.98
7.31
8.34
10.71
173.66
2054.27
52.44

8.86
197.72
176.04
14.88
695.7
29.52
258.48
4.36
45.36
246.38
12.39
35.0
897.15
49.12
2.9
279.9
6.53
24.96
13.78
419.4
302.94
607.52
7.87
113.37
7.38
3.05
68.6
6.48
45.4
117.14
662.88
311.98
18.96
19.99
25.96
7.97
1215.92
3.68
11.56
239.96
177.57
17.04
30.77
18.69
311.98
74.11
40.1
10.37

122.97
2.31
6.85
5.34
165.6
79.96
8.67
209.94
355.36
3.64
12.26
51.56
165.28
431.16
723.92
168.1
10.78
22.61
200.8
143.7
21.12
1.63
368.91
34.66
675.06
8.64
23.38
26.46
48.67
8.38
11.57
79.92
7.24
113.6
18.9
72.64
2.26
3.54
11.7
65.97
25.2
61.68
6.87
4.24
7.78
19.44
182.72
75.79

39.92
31.74
7.31
112.65
129.93
17.45
121.6
587.97
390.75
29.12
15.92
2.47
18.16
1793.98
343.2
54.37
60.12
38.08
2.62
0.85
55.99
1575.14
2.67
3.6
63.92
17.24
40.99
18.06
32.34
1.8
3.88
6.48
91.18
12.96
18.7
41.96
232.88
47.94
25.86
197.97
171.96
41.91
1049.93
199.96
1.5
92.94
37.76
15.55

27.36
5.56
1871.88
16.02
100.79
555.96
369.58
31.4
271.9
55.42
31.16
166.44
79.98
698.35
7.82
8.72
515.88
47.36
533.94
383.98
261.96
6.9
10.78
174.95
7.71
14.85
5.25
139.94
54.22
18.24
339.96
895.92
95.94
269.98
105.52
84.09
296.71
8.75
159.98
100.7
69.5
7.28
2518.29
9.66
77.55
108.4
49.85
7.98

296.85
153.58
11.82
490.32
12.96
95.98
15.38
9.25
44.91
9.54
12.54
629.06
85.22
400.8
799.92
22.2
71.92
32.4
11.52
9.34
17.57
118.78
317.06
657.55
7.83
671.98
56.98
13.9
229.94
38.98
4.37
41.04
1123.13
5.84
26.38
1586.69
15.7
8.7
132.16
14.37
25.16
359.5
17.15
17.12
10.9
50.8
539.66
51.75

25.06
10.43
2.5
481.57
1.68
171.04
1669.6
263.96
18.09
3359.95
26.4
149.95
127.76
127.37
123.92
17.46
19.68
3.98
2.39
20.74
2799.96
26.64
91.01
1363.96
16.72
225.57
19.14
271.76
786.48
45.0
16.7
46.87
255.76
122.33
45.84
37.21
1.44
7.58
230.38
23.1
62.02
35.34
71.98
8.35
60.72
561.57
8.26
35.0

25.06
487.98
243.99
74.24
33.4
287.97
15.12
149.9
109.59
10.78
697.16
169.45
51.15
6.57
364.08
7.97
405.86
91.36
6.19
10.37
563.94
8.22
1.82
8.26
77.52
111.96
3.82
263.88
1082.48
5.25
47.62
22.83
79.12
7.71
59.71
10.68
393.25
43.8
10.74
19.15
965.85
699.98
1127.98
16.59
9.21
52.2
8.56
135.99

975.92
115.3
19.14
58.41
32.9
89.54
28.08
262.11
53.88
9.66
75.18
29.34
6.67
17.9
4.54
43.1
2.61
39.15
19.14
36.0
6.98
25.83
11.36
1.73
1.98
36.78
209.88
151.62
946.34
49.62
485.88
12.96
12.39
94.85
56.16
539.97
53.72
5.16
104.85
265.17
619.95
70.56
360.71
4.98
6.64
16.28
756.8
968.74

51.97
105.58
319.96
240.74
90.99
129.98
67.78
20.34
59.52
1119.89
307.67
56.56
14.8
45.7
12.72
15.55
222.38
1117.92
207.0
105.96
17.47
105.55
67.56
24.05
67.54
303.92
186.54
239.98
1079.85
665.88
94.99
12.96
201.58
5.76
32.4
9.14
13.12
88.78
29.9
110.11
377.97
15.24
95.1
63.77
7.3
32.4
658.75
27.02

60.72
306.2
102.62
242.94
12.82
629.1
33.87
132.52
36.48
5.68
9.16
6.26
8.26
81.92
27.38
164.88
4.94
12.13
32.67
1158.12
5.18
34.68
27.81
1879.96
127.87
8.78
13.71
1499.95
125.13
8.26
19.68
15.55
7.23
26.18
74.45
22.96
90.64
47.79
979.95
76.64
19.87
43.13
40.78
51.17
23.04
479.94
138.56
5.18

44.75
66.69
21.56
4.75
85.23
1035.8
94.68
3.01
587.97
38.82
114.2
13.39
12.53
273.57
3.28
15.59
46.96
27.6
17.22
900.08
9.96
349.95
3.13
299.97
25.4
323.37
1879.96
322.59
29.8
59.52
110.97
479.97
18.18
35.17
74.35
5.76
14.0
20.96
36.56
83.84
33.36
6.68
447.94
50.35
121.6
33.9
43.18
10.36

22.32
7.3
26.98
37.68
1295.78
22.14
38.09
42.6
411.33
2.63
9.84
87.71
3.21
17.48
31.98
4.73
5.95
2003.92
143.95
408.01
18.0
116.28
35.92
242.9
271.96
79.36
341.96
32.4
43.3
1919.98
9.24
19.96
315.2
9.16
371.97
58.68
85.25
11.52
271.96
42.76
24.85
167.94
11.67
111.79
192.16
4.54
854.35
41.4

41.24
3.15
103.94
21.24
5.21
63.92
1999.96
59.7
7.56
14.9
91.2
13.34
155.98
55.94
83.17
56.52
47.58
8.94
5.88
9.48
24.85
441.96
680.01
219.18
7.76
10.27
21.38
5.04
407.98
8.02
20.93
26.98
180.02
6.72
283.92
218.35
90.48
59.94
96.36
37.06
77.88
207.18
100.8
19.46
21.93
3.64
216.4
164.22

599.99
106.96
659.9
167.86
300.77
65.08
29.97
269.49
280.79
12.22
82.95
2.37
92.52
32.07
143.64
516.96
18.53
10.16
2279.96
196.75
207.98
14.62
45.36
47.04
791.88
38.98
60.42
1.41
6.05
12.67
2.52
110.96
24.9
186.54
196.78
199.98
499.17
1056.86
329.58
78.26
84.42
40.68
37.84
242.14
447.86
68.7
29.33
23.13

30.34
252.8
89.34
3.32
303.25
47.81
107.98
44.78
7.66
8.64
29.22
853.93
26.18
8.45
40.68
27.15
86.35
27.18
278.4
47.98
46.38
17.92
383.61
62.94
26.0
20.32
1.24
123.14
32.06
177.54
37.4
8.96
4.5
10.02
14.94
6.0
2621.32
41.86
124.79
25.98
291.17
123.92
36.4
220.98
18.69
287.98
5.98
11.65

14.28
436.0
488.65
5.88
4.71
2.82
10.82
26.72
1.19
887.1
129.98
33.29
10.95
22.51
39.92
699.98
85.9
9.43
170.88
356.79
696.42
26.34
4.93
140.74
1633.14
146.14
294.93
121.68
31.5
24.67
18.56
33.96
2.95
21.21
193.95
2430.08
113.33
1.64
208.56
43.31
459.88
132.7
261.74
247.44
13.36
20.24
259.92
75.0

10.37
17.96
17.46
20.74
7.92
38.88
110.98
115.02
70.08
272.94
7.16
160.98
447.84
241.18
862.34
1685.88
142.78
40.75
88.92
45.48
15.88
182.11
25.92
14.94
6.41
2.89
26.96
20.56
638.29
519.68
339.14
20.74
15.98
46.35
7.76
12.54
423.65
63.82
212.8
14.94
119.45
45.36
31.98
16.75
37.94
4.66
6.48
43.86

4548.81
23.08
1.59
20.1
1496.16
213.48
12.96
302.38
71.6
22.75
25.92
18.9
304.78
389.06
3.1
211.96
51.97
249.95
21.36
844.12
26.2
17.3
6.08
7.31
35.98
156.51
799.56
47.98
13.57
131.98
619.95
13.85
318.08
1049.2
37.32
62.92
64.78
13.38
375.34
122.92
59.52
9.82
9.57
145.76
46.67
11.2
541.24
85.96

48.78
209.97
11.91
254.35
291.1
8.62
15.53
17.22
142.86
61.06
5.72
104.23
131.88
221.16
207.14
14.62
713.88
14.35
23.92
21.38
671.93
14.76
46.0
59.1
14.62
87.36
348.93
359.32
2.78
16.46
777.21
25.06
23.47
12.06
11.33
2254.41
64.78
340.7
186.14
128.06
43.18
6.9
248.57
38.04
100.0
254.97
28.49
344.91

77.24
3.17
4.34
53.72
69.9
56.7
3.88
199.76
39.96
16.52
255.98
16.52
24.0
143.7
149.23
85.14
99.98
559.92
6.6
125.7
13.86
172.7
785.88
11.52
44.46
17.31
30.53
60.14
30.84
1799.97
18.24
177.0
182.72
33.96
24.22
206.43
219.9
58.48
9.68
288.0
1395.54
5.78
11.26
24.27
11.36
25.92
19.44
0.56

114.9
787.53
12.32
8.54
341.96
96.96
113.89
29.47
240.78
3.81
6.26
21.24
235.94
28.67
106.08
136.92
359.97
262.24
19.44
118.99
191.97
122.14
179.97
78.3
15.01
122.35
32.78
182.22
47.99
27.86
120.77
67.84
209.88
15.55
4.45
263.96
53.94
201.58
137.94
236.0
122.12
4.13
20.7
199.99
5.1
2.9
108.78
146.82

191.98
47.36
7.82
104.88
22.7
14.94
187.98
783.96
83.88
46.38
42.28
22.58
94.74
30.53
17.54
68.72
1294.75
241.44
20.96
14.2
27.06
5.18
28.16
11.52
95.97
11.76
97.57
255.11
301.96
25.63
159.96
449.91
32.4
25.3
40.78
8.82
52.99
12.96
33.79
84.78
36.11
89.97
5.02
65.57
377.97
361.96
38.52
38.09

39.96
686.4
10.37
149.35
5.9
38.9
28.28
167.54
3610.85
14.7
120.78
772.47
42.98
720.76
133.38
152.69
639.97
8.26
28.4
557.59
98.16
31.56
184.7
641.96
147.92
27.17
173.94
61.12
152.94
617.97
32.4
237.1
135.52
7.97
45.58
629.1
427.64
541.44
20.86
173.24
23.34
18.9
619.15
70.95
2575.94
159.75
2152.78
191.08

126.56
47.95
53.04
3.15
57.42
449.15
19.92
14.02
310.74
5.34
10.65
209.99
81.54
12.96
33.92
166.84
49.85
15.76
751.92
99.99
14.38
4.89
9.82
30.56
81.98
361.92
161.82
21.2
72.6
47.18
41.9
20.04
380.86
15.55
285.55
143.86
209.67
359.98
42.62
479.95
181.35
47.19
12.58
51.97
225.3
52.59
113.82
9.64

11.66
21.12
67.36
19.73
289.24
11.23
1564.29
79.99
9.45
60.12
10.92
22.96
60.77
653.55
73.98
70.69
447.94
12.78
11.3
18.94
4.28
44.13
11.56
59.52
449.97
31.1
99.95
34.65
18.27
36.29
214.7
234.36
5083.96
47.32
79.9
136.46
866.65
182.94
40.46
6.29
131.88
8.32
77.73
452.94
62.91
55.98
39.89
1336.83

1347.52
27.93
31.86
16.68
896.99
155.88
127.3
55.36
164.39
31.4
4.72
41.92
107.98
113.94
3.9
34.6
1087.94
142.18
381.36
35.91
7.56
371.98
98.35
148.48
15.56
28.4
255.85
20.46
43.92
6.69
30.32
8.34
305.01
113.89
140.74
999.43
479.98
163.96
2.67
8.94
64.17
2.89
77.95
115.36
18.65
37.43
575.92
3.33

866.4
43.19
39.96
95.84
77.6
474.95
2.51
104.9
158.37
35.45
542.94
14.98
281.37
6.1
23.88
11.65
162.6
15.98
17.54
4.51
4.13
72.48
1259.93
95.97
8.83
45.84
10.69
88.83
13.49
7.07
49.12
10.5
34.44
14.52
31.78
51.84
47.98
14.98
45.24
71.25
20.74
5.76
3.26
14.62
130.71
302.67
180.96
17.64

249.58
863.88
3812.97
1673.18
3266.38
22.64
13.98
49.12
10.27
20.8
72.42
36.24
45.98
88.8
464.0
13.49
273.92
511.5
144.78
10.37
586.4
57.4
6.48
323.98
26.16
436.0
906.68
27.89
15.96
108.96
38.52
35.56
74.35
87.92
354.9
43.58
4.62
12.96
21.86
408.74
2548.56
98.38
8.7
166.44
5.48
378.0
3.75
469.99

19.75
833.94
13.13
13.48
32.4
86.35
116.31
92.96
29.9
481.32
599.29
29.8
12.96
19.8
49.12
33.09
8.9
265.86
127.88
10.27
244.55
100.49
13.22
350.98
40.78
23.12
70.88
61.96
18.84
52.51
4.93
272.65
286.79
3.2
9.42
11.68
93.15
48.82
430.88
208.16
38.88
35.28
11.96
18.84
104.85
4.82
150.41
5.31

19.83
51.75
10.95
13.9
508.7
102.13
10.86
148.7
11.96
42.38
60.89
387.14
46.8
47.82
563.81
185.38
614.27
8.28
122.94
12.96
22.75
361.76
42.76
10.37
264.18
85.52
21.98
6.63
106.5
270.72
13.27
7.5
261.96
45.98
36.88
599.97
74.94
16.18
1059.12
44.38
5.76
826.11
26.4
333.0
15.0
842.72
37.41
106.32

1454.9
21.44
7.99
4.26
104.9
13.98
686.32
176.8
122.97
205.33
9.35
754.45
84.78
103.92
1454.49
35.17
41.47
7.04
142.78
179.97
241.96
117.46
424.12
561.58
2.96
4.28
31.68
595.38
46.24
13.9
35.49
25.16
239.67
12.96
22.96
7.58
20.7
8.64
21.98
11.76
15.55
144.12
5.55
863.64
5.18
301.96
8.45
15.96

368.91
11.34
20.04
64.96
10.08
271.76
9.3
6.26
314.35
19.04
5.18
94.2
189.0
21.56
5.47
42.85
6.73
145.57
118.0
279.86
601.54
824.97
89.95
155.88
212.94
9.52
27.42
15.75
5.18
13.39
51.84
20.74
17.14
209.79
59.97
205.67
31.1
11.76
24.64
8.32
8.86
7.04
4.73
14.73
362.94
307.98
61.96
19.44

3991.98
206.1
1.87
372.14
844.12
15.88
38.88
50.94
37.31
15.55
300.93
76.12
201.58
88.8
89.97
44.67
12.42
57.59
52.4
48.66
17.04
5.23
11.76
155.37
101.99
41.94
81.42
143.96
5.46
9.96
119.83
9.14
526.45
572.76
42.6
26.49
38.88
959.98
221.06
87.84
206.0
4.96
60.64
58.05
26.35
160.0
491.55
32.45

64.4
11.07
383.96
895.92
23.74
21.56
160.72
801.96
53.57
8.73
342.37
19.99
3.01
59.76
549.99
55.62
30.69
63.97
1415.76
167.84
13.9
45.36
4.98
31.44
5.34
52.51
18.94
603.92
6.56
64.96
18.09
36.05
40.18
212.94
61.93
801.57
299.97
63.96
129.3
917.92
46.86
114.52
32.9
72.78
154.9
28.85
15.26
35.17

31.92
4.3
166.5
755.96
3.48
10.37
556.66
24.45
89.58
10.76
11.16
41.4
38.43
3.65
703.97
1591.02
194.53
159.98
211.96
209.57
636.86
89.97
15.55
909.12
19.46
1.19
523.26
298.12
49.12
78.35
79.96

```
[180]: mse=mean_squared_error(b,a)
      rms=np.sqrt(mse)
      print("mse :",mse)
      print("rms :",rms)
```

```
mse : 16339.518886085545
rms  : 127.82612755648019
```

```
[181]: dfpred = pd.DataFrame({'y_preds': a, 'y_tests': b})
      dfpred.sample(30)
```

```
[181]:
```

	y_preds	y_tests
795	176.69	218.35
1407	160.63	148.48
470	7.85	8.26

643	1639.14	1499.95
1260	94.66	98.16
1731	71.86	76.12
1381	52.09	55.98
1679	6.19	6.26
81	2.49	1.34
79	119.53	107.44
1215	9.08	11.52
586	212.51	207.00
1236	432.77	377.97
827	198.62	196.75
1458	10.64	10.69
1332	210.78	225.30
1492	44.33	45.98
428	2.74	1.68
1558	5.47	4.93
713	7.31	7.30
417	359.00	359.50
1499	10.25	10.37
1041	9.35	9.82
476	31.96	33.40
17	102.17	89.99
1438	2.36	2.51
125	21.33	18.84
336	264.10	271.90
899	159.43	220.98
133	2.36	1.64

```
[182]: mae = mean_absolute_error(b,a)
mae
```

```
[182]: 29.550157011369787
```

```
[183]: for i in y_pred:
        if round(np.exp(i),2)<0:
            print(round(np.exp(i),2))
```

35 Giving a summary over interpretation of data and Suggesting pros and cons of methods used.

36 Interpretation of Data:

36.1 Handling Missing Values and Errors:

36.1.1 The identification and handling of missing values and errors are crucial steps in ensuring the reliability and accuracy of the analysis. By applying appropriate imputation techniques and dropping remaining missing values, the integrity of the dataset is preserved.

36.1.2 Understanding the nature and distribution of missing values can provide insights into data collection processes and potential areas for improvement in data collection protocols.

36.2 Transformation for Outliers:

36.2.1 Outliers can significantly influence model performance and may distort the results of statistical analyses. By applying transformations to ‘profit’ and ‘sales’ columns, the impact of outliers is mitigated, leading to more robust and reliable results.

36.2.2 Exploring various outlier trimming techniques allows for a comprehensive understanding of the data distribution and the identification of potential anomalies that may require further investigation.

36.3 Modeling Approach:

36.3.1 Recognizing the limitations of linear models for complex data structures highlights the importance of selecting appropriate modeling techniques that can capture nonlinear relationships and interactions among variables.

36.3.2 The application of feature selection techniques enhances model interpretability by identifying the most influential predictors of sales. This allows for more focused analysis and facilitates actionable insights for decision-makers.

36.3.3 While achieving a high accuracy of 98% using MLPRegressor is commendable, it’s essential to consider the trade-offs between model complexity, computational resources, and interpretability when selecting modeling approaches.

36.4 Data Quality Assessment:

36.4.1 Conducting a comprehensive assessment of data quality, including identifying and addressing missing values, errors, and outliers, demonstrates a commitment to ensuring the reliability and validity of the analysis results.

36.4.2 Documenting the data preprocessing steps undertaken provides transparency and reproducibility, enabling others to replicate the analysis and validate the findings.

36.5 Continuous Improvement:

36.5.1 The process of data analysis is iterative, and there is always room for improvement. Regularly reviewing and refining data preprocessing techniques, modeling approaches, and interpretation methods can lead to more accurate and actionable insights over time.

36.5.2 Soliciting feedback from stakeholders and incorporating domain knowledge can further enhance the quality and relevance of the analysis results.

[]: