```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
from sklearn.linear_model import LinearRegression


data = pd.read_csv("/content/SampleSuperstore.csv")


data.shape
```

```
(9994, 13)
```

```python
data.head(10)
```

| | Ship Mode | Segment | Country | City | State | Postal Code | Region | Category | Cate |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Second Class | Consumer | United States | Henderson | Kentucky | 42420 | South | Furniture | Booko |
| 1 | Second Class | Consumer | United States | Henderson | Kentucky | 42420 | South | Furniture | C |
| 2 | Second Class | Corporate | United States | Los Angeles | California | 90036 | West | Office Supplies | L |
| 3 | Standard Class | Consumer | United States | Fort Lauderdale | Florida | 33311 | South | Furniture | T |
| 4 | Standard Class | Consumer | United States | Fort Lauderdale | Florida | 33311 | South | Office Supplies | Sto |
| 5 | Standard Class | Consumer | United States | Los Angeles | California | 90032 | West | Furniture | Furnish |

```python
data.isnull().sum()
```

```
Ship Mode       0
Segment         0
Country         0
City            0
State           0
Postal Code     0
Region          0
Category        0
Sub-Category    0
Sales           0
Quantity        0
Discount        0
Profit          0
dtype: int64
```

```python
data.isnull()
```

| | Ship Mode | Segment | Country | City | State | Postal Code | Region | Category | Sub-Category | Sales | Quantity | Discount | Profit |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | False | False | False | False | False | False | False | False | False | False | False | False | False |
| 1 | False | False | False | False | False | False | False | False | False | False | False | False | False |
| 2 | False | False | False | False | False | False | False | False | False | False | False | False | False |
| 3 | False | False | False | False | False | False | False | False | False | False | False | False | False |
| 4 | False | False | False | False | False | False | False | False | False | False | False | False | False |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 9989 | False | False | False | False | False | False | False | False | False | False | False | False | False |
| 9990 | False | False | False | False | False | False | False | False | False | False | False | False | False |
| 9991 | False | False | False | False | False | False | False | False | False | False | False | False | False |
| 9992 | False | False | False | False | False | False | False | False | False | False | False | False | False |
| 9993 | False | False | False | False | False | False | False | False | False | False | False | False | False |

9994 rows × 13 columns

```python
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9994 entries, 0 to 9993
```

```
Data columns (total 13 columns):
 #   Column        Non-Null Count  Dtype
---  ------        --------------  -----
 0   Ship Mode     9994 non-null   object
 1   Segment       9994 non-null   object
 2   Country       9994 non-null   object
 3   City          9994 non-null   object
 4   State         9994 non-null   object
 5   Postal Code   9994 non-null   int64
 6   Region        9994 non-null   object
 7   Category      9994 non-null   object
 8   Sub-Category  9994 non-null   object
 9   Sales         9994 non-null   float64
 10  Quantity      9994 non-null   int64
 11  Discount      9994 non-null   float64
 12  Profit        9994 non-null   float64
dtypes: float64(3), int64(2), object(8)
memory usage: 1015.1+ KB
```

data.info

```
<bound method DataFrame.info of          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
...                ...          ...            ...              ...          ...
9989      Second Class     Consumer  United States            Miami      Florida
9990    Standard Class     Consumer  United States       Costa Mesa   California
9991    Standard Class     Consumer  United States       Costa Mesa   California
9992    Standard Class     Consumer  United States       Costa Mesa   California
9993      Second Class     Consumer  United States      Westminster   California

      Postal Code Region         Category Sub-Category      Sales  Quantity  \
0           42420  South        Furniture    Bookcases   261.9600         2
1           42420  South        Furniture       Chairs   731.9400         3
2           90036   West  Office Supplies       Labels    14.6200         2
3           33311  South        Furniture       Tables   957.5775         5
4           33311  South  Office Supplies      Storage    22.3680         2
...           ...    ...              ...          ...        ...       ...
9989        33180  South        Furniture  Furnishings    25.2480         3
9990        92627   West        Furniture  Furnishings    91.9600         2
9991        92627   West       Technology       Phones   258.5760         2
9992        92627   West  Office Supplies        Paper    29.6000         4
9993        92683   West  Office Supplies   Appliances   243.1600         2

      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      0.20     4.1028
9990      0.00    15.6332
9991      0.20    19.3932
9992      0.00    13.3200
9993      0.00    72.9480

[9994 rows x 13 columns]>
```

data.describe

```
<bound method NDFrame.describe of          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
...                ...          ...            ...              ...          ...
9989      Second Class     Consumer  United States            Miami      Florida
9990    Standard Class     Consumer  United States       Costa Mesa   California
9991    Standard Class     Consumer  United States       Costa Mesa   California
9992    Standard Class     Consumer  United States       Costa Mesa   California
9993      Second Class     Consumer  United States      Westminster   California

      Postal Code Region         Category Sub-Category      Sales  Quantity  \
0           42420  South        Furniture    Bookcases   261.9600         2
1           42420  South        Furniture       Chairs   731.9400         3
2           90036   West  Office Supplies       Labels    14.6200         2
3           33311  South        Furniture       Tables   957.5775         5
4           33311  South  Office Supplies      Storage    22.3680         2
...           ...    ...              ...          ...        ...       ...
9989        33180  South        Furniture  Furnishings    25.2480         3
9990        92627   West        Furniture  Furnishings    91.9600         2
9991        92627   West       Technology       Phones   258.5760         2
9992        92627   West  Office Supplies        Paper    29.6000         4
```

```
9993        92683   West  Office Supplies   Appliances  243.1600            2

      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      0.20     4.1028
9990      0.00    15.6332
9991      0.20    19.3932
9992      0.00    13.3200
9993      0.00    72.9480

[9994 rows x 13 columns]>
```

```
data.describe()
```

|       | Postal Code   | Sales         | Quantity     | Discount     | Profit        |
|-------|---------------|---------------|--------------|--------------|---------------|
| count | 9994.000000   | 9994.000000   | 9994.000000  | 9994.000000  | 9994.000000   |
| mean  | 55190.379428  | 229.858001    | 3.789574     | 0.156203     | 28.656896     |
| std   | 32063.693350  | 623.245101    | 2.225110     | 0.206452     | 234.260108    |
| min   | 1040.000000   | 0.444000      | 1.000000     | 0.000000     | -6599.978000  |
| 25%   | 23223.000000  | 17.280000     | 2.000000     | 0.000000     | 1.728750      |
| 50%   | 56430.500000  | 54.490000     | 3.000000     | 0.200000     | 8.666500      |
| 75%   | 90008.000000  | 209.940000    | 5.000000     | 0.200000     | 29.364000     |
| max   | 99301.000000  | 22638.480000  | 14.000000    | 0.800000     | 8399.976000   |

```
data.isna().any()
```

```
Ship Mode       False
Segment         False
Country         False
City            False
State           False
Postal Code     False
Region          False
Category        False
Sub-Category    False
Sales           False
Quantity        False
Discount        False
Profit          False
dtype: bool
```
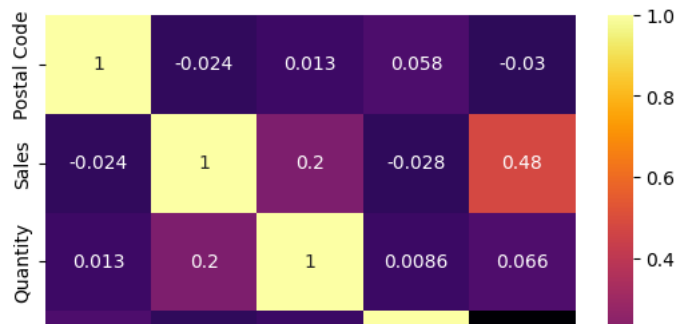
```
correlation=data.corr()
correlation
```

```
<ipython-input-15-d7a18ccdee06>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a f
  correlation=data.corr()
```

|             | Postal Code | Sales     | Quantity  | Discount  | Profit    |
|-------------|-------------|-----------|-----------|-----------|-----------|
| Postal Code | 1.000000    | -0.023854 | 0.012761  | 0.058443  | -0.029961 |
| Sales       | -0.023854   | 1.000000  | 0.200795  | -0.028190 | 0.479064  |
| Quantity    | 0.012761    | 0.200795  | 1.000000  | 0.008623  | 0.066253  |
| Discount    | 0.058443    | -0.028190 | 0.008623  | 1.000000  | -0.219487 |
| Profit      | -0.029961   | 0.479064  | 0.066253  | -0.219487 | 1.000000  |

```
import seaborn as sns
sns.heatmap(data.corr(),annot=True,cmap='inferno')
```

<ipython-input-16-bdc76ecc106c>:2: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a f
  sns.heatmap(data.corr(),annot=True,cmap='inferno')
<Axes: >



```
data['City'].value_counts()
```
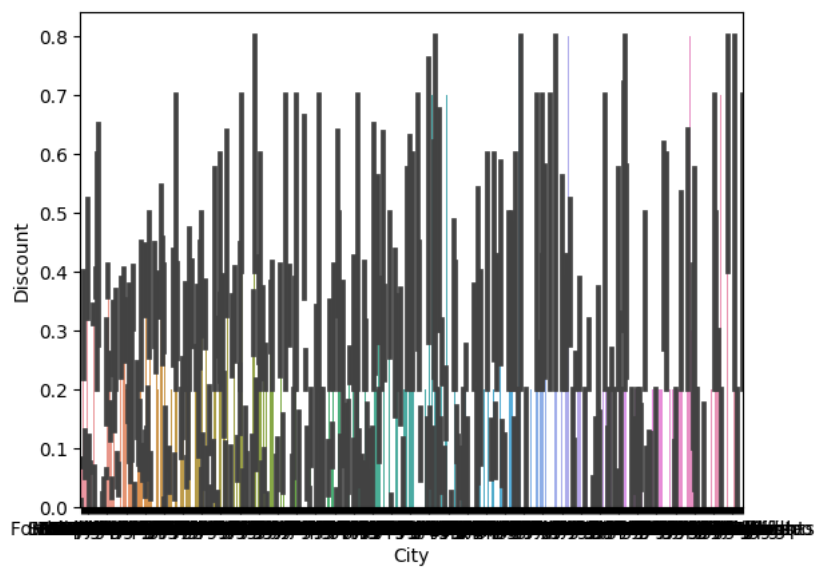
```
New York City      915
Los Angeles        747
Philadelphia       537
San Francisco      510
Seattle            428
                  ...
Glenview             1
Missouri City        1
Rochester Hills      1
Palatine             1
Manhattan            1
Name: City, Length: 531, dtype: int64
```

```
print(data["City"].unique())
```
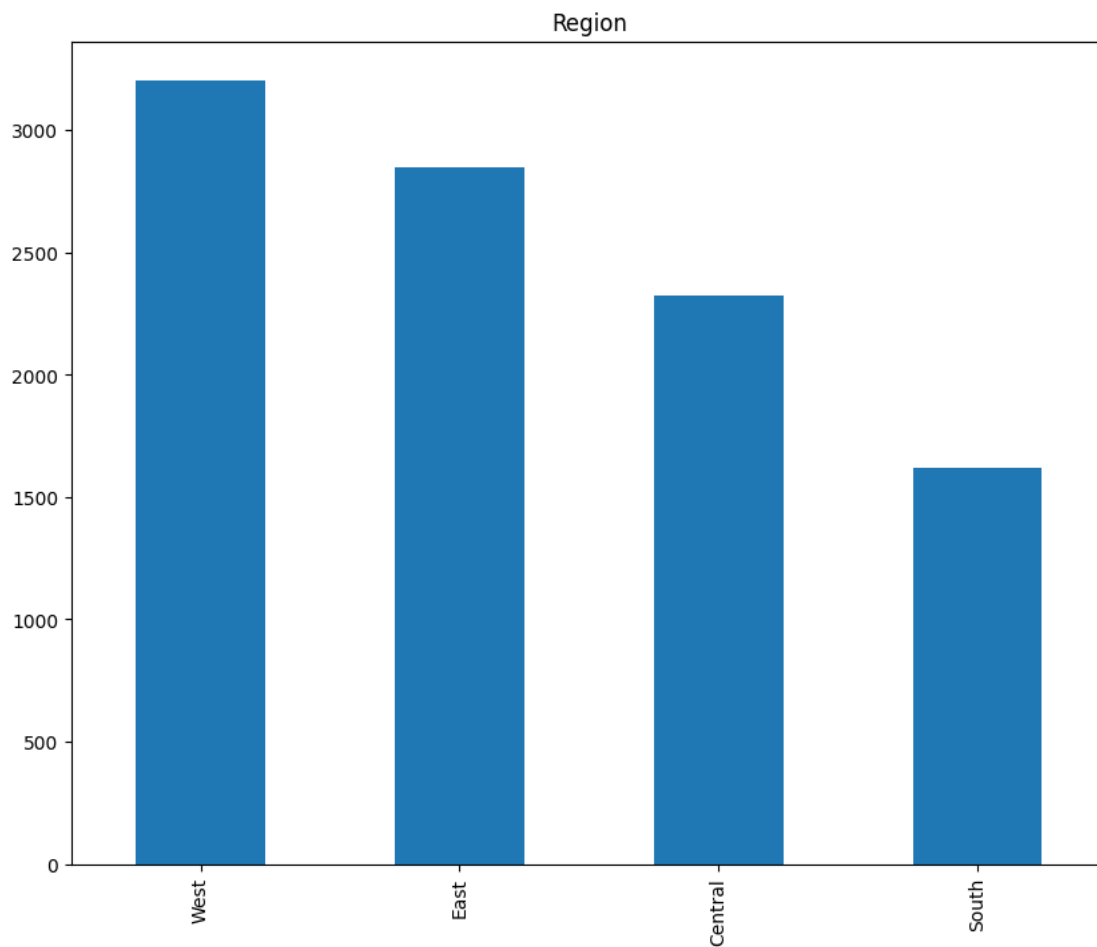
```
['Henderson' 'Los Angeles' 'Fort Lauderdale' 'Concord' 'Seattle'
 'Fort Worth' 'Madison' 'West Jordan' 'San Francisco' 'Fremont'
 'Philadelphia' 'Orem' 'Houston' 'Richardson' 'Naperville' 'Melbourne'
 'Eagan' 'Westland' 'Dover' 'New Albany' 'New York City' 'Troy' 'Chicago'
 'Gilbert' 'Springfield' 'Jackson' 'Memphis' 'Decatur' 'Durham' 'Columbia'
 'Rochester' 'Minneapolis' 'Portland' 'Saint Paul' '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' 'Redlands' 'Hamilton' 'Westfield' 'Akron'
 'Denver' 'Dallas' 'Whittier' 'Saginaw' 'Medina' 'Dublin' 'Detroit'
 'Tampa' 'Santa Clara' 'Lakeville' 'San Diego' 'Brentwood' 'Chapel Hill'
 'Morristown' 'Cincinnati' 'Inglewood' 'Tamarac' 'Colorado Springs'
 'Belleville' 'Taylor' 'Lakewood' 'Arlington' 'Arvada' 'Hackensack'
 'Saint Petersburg' 'Long Beach' 'Hesperia' 'Murfreesboro' 'Layton'
 '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' 'Norman' 'Park Ridge'
 'Amarillo' 'Lindenhurst' 'Huntsville' 'Fayetteville' 'Costa Mesa'
 'Parker' 'Atlanta' 'Gladstone' 'Great Falls' 'Lakeland' 'Montgomery'
 'Mesa' 'Green Bay' 'Anaheim' 'Marysville' 'Salem' 'Laredo' 'Grove City'
 'Dearborn' 'Warner Robins' 'Vallejo' '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' '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'
 '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' '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' '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'
```

```python
data['Profit'].value_counts()
```

```
        0.0000    65
        6.2208    43
        9.3312    38
        5.4432    32
        3.6288    32
                  ..
        83.2508    1
        16.1096    1
        7.1988     1
        1.6510     1
        72.9480    1
        Name: Profit, Length: 7287, dtype: int64
```

```python
data['Ship Mode'].value_counts()
```

```
        Standard Class    5968
        Second Class      1945
        First Class       1538
        Same Day           543
        Name: Ship Mode, dtype: int64
```

```python
sns.barplot(x = 'Ship Mode',y = 'Quantity',data = data)
plt.show()
```



```python
sns.barplot(x = 'City',y = 'Discount',data = data)
plt.show()
```

```
t1 = data['Region'].value_counts()[:150]
t1.plot(kind='bar',figsize=(10,8))
plt.title('Region')
```

Text(0.5, 1.0, 'Region')



```
print(data["Country"].unique())
```

['United States']

```
print(data["Segment"].unique())
```

['Consumer' 'Corporate' 'Home Office']

```
data['Segment'].value_counts()
```

```
Consumer      5191
Corporate     3020
Home Office   1783
Name: Segment, dtype: int64
```

```
s1 = data['Segment'].value_counts()[:150]
s1.plot(kind='pie',figsize=(10,8))
plt.title('Segment')
```

Text(0.5, 1.0, 'Segment')

# Segment

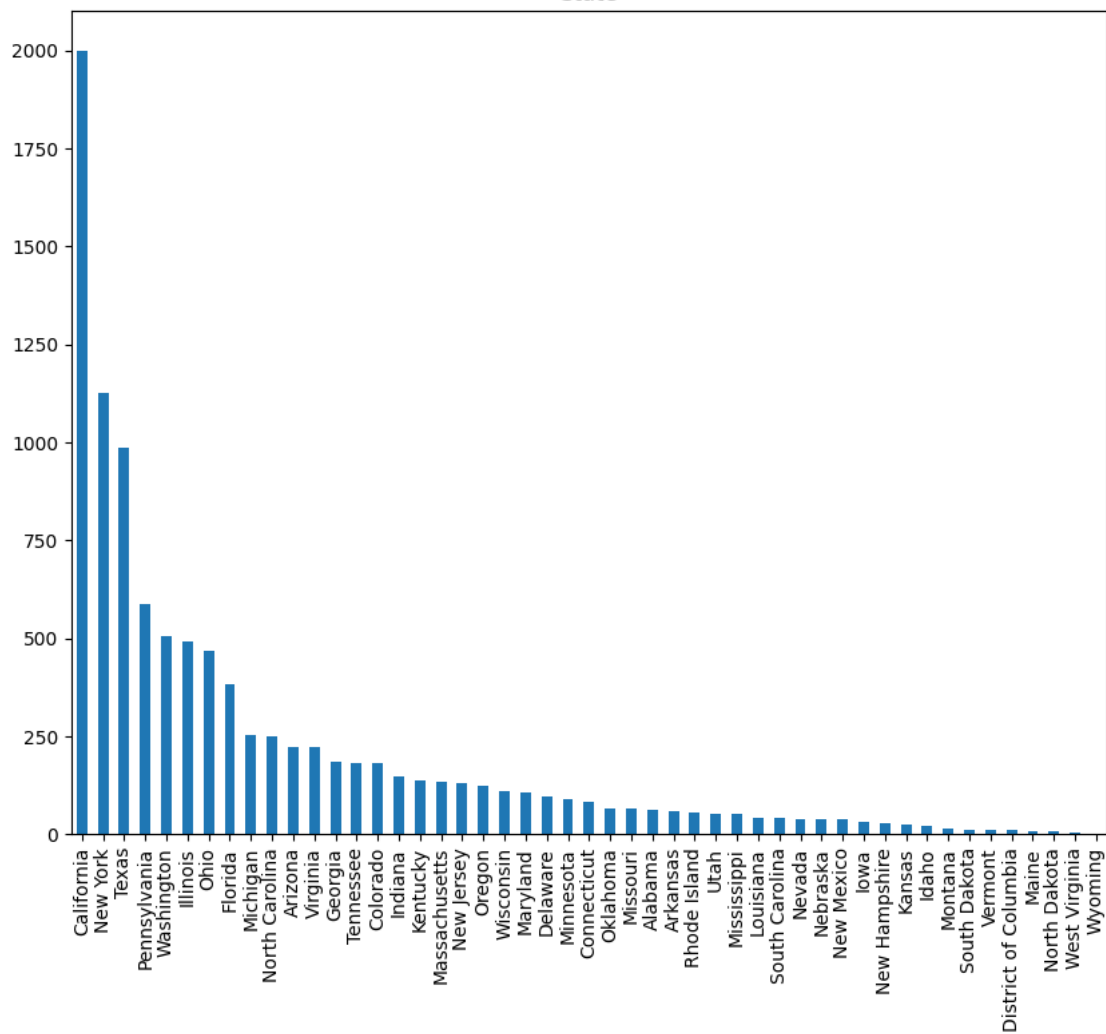### Consumer



```
R1 = data['State'].value_counts()[:150]
R1.plot(kind='bar',figsize=(10,8))
plt.title('State')
```
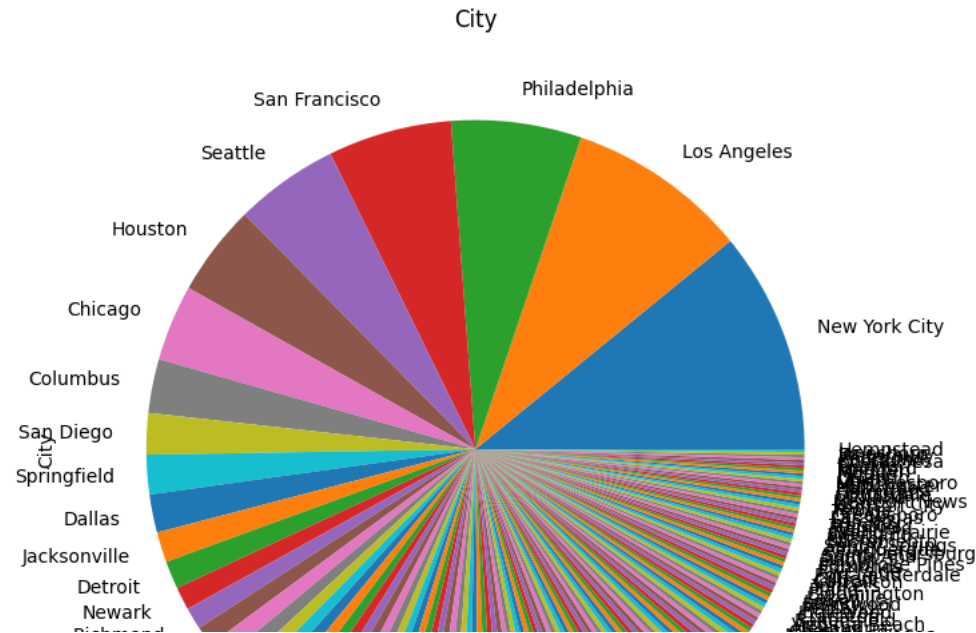
Text(0.5, 1.0, 'State')

# State



```
c1 = data['City'].value_counts()[:150]
c1.plot(kind='pie',figsize=(10,8))
plt.title('City')
```

Text(0.5, 1.0, 'City')

## City



```
data['Sub-Category'].value_counts()
```

```
    Binders        1523
    Paper          1370
    Furnishings     957
    Phones          889
    Storage         846
    Art             796
    Accessories     775
    Chairs          617
    Appliances      466
    Labels          364
    Tables          319
    Envelopes       254
    Bookcases       228
    Fasteners       217
    Supplies        190
    Machines        115
    Copiers          68
    Name: Sub-Category, dtype: int64
```

```
data['Sub-Category'].value_counts().sum()
```

```
    9994
```

```
data.cov()
```

```
    <ipython-input-32-72e63cb34c7c>:1: FutureWarning: The default value of numeric_only in DataFrame.cov is deprecated. In a fu
      data.cov()
```

|  | Postal Code | Sales | Quantity | Discount | Profit |
|---|---|---|---|---|---|
| **Postal Code** | 1.028080e+09 | -476682.766590 | 910.415885 | 386.870404 | -225045.849445 |
| **Sales** | -4.766828e+05 | 388434.455308 | 278.459923 | -3.627228 | 69944.096586 |
| **Quantity** | 9.104159e+02 | 278.459923 | 4.951113 | 0.003961 | 34.534769 |
| **Discount** | 3.868704e+02 | -3.627228 | 0.003961 | 0.042622 | -10.615173 |
| **Profit** | -2.250458e+05 | 69944.096586 | 34.534769 | -10.615173 | 54877.798055 |

```
data['Sub-Category'].value_counts().mean()
```

```
    587.8823529411765
```

```
data['Segment'].value_counts().mean()
```

```
    3331.3333333333335
```