n [160 ut[160	Ship Mode Segment Country City State Postal Code Region Category Sub-Category Sales Quantity Discount Profit Second Class Consumer United States Henderson Kentucky 42420 South Furniture Bookcases 261.9600 2 0.00 41.9136 Second Class Consumer United States Henderson Kentucky 42420 South Furniture Chairs 731.9400 3 0.00 219.5820 Second Class Corporate United States Los Angeles California 90036 West Office Supplies Labels 14.6200 2 0.00 6.8714 Standard Class Consumer United States Fort Lauderdale Florida 33311 South Furniture Tables 957.5775 5 0.45 -383.0310 data.tail() Ship Mode Segment Country City State Postal Code Region Category Sub-Category Sub-Category Sales Quantity Discount Profit
n [161	9989 Second Class Consumer United States Miami Florida 33180 South Furniture Furnishings 25.248 3 0.2 4.1028 9990 Standard Class Consumer United States Costa Mesa California 92627 West Technology Phones 258.576 2 0.2 19.3932 9991 Standard Class Consumer United States Costa Mesa California 92627 West Technology Phones 258.576 2 0.2 19.3932 9992 Standard Class Consumer United States Costa Mesa California 92627 West Office Supplies Paper 29.600 4 0.0 13.3200 9993 Second Class Consumer United States Westminster California 92683 West Office Supplies Appliances 243.160 2 0.0 72.9480 data info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 9994 entries, 0 to 9993 Data Columns (total 13 columns): # Column Non-Null Count Dtype</class>
	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 object 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
n [162 ut[162	Ship Mode Segment Country City State Postal Code Region Category Sub-Category Sales Quantity Discount Profit
n [163 ut[163	9990 True True <th< td=""></th<>
	1FalseFalseFalseFalseFalseFalseFalseFalseFalseFalseFalseFalse2FalseFalseFalseFalseFalseFalseFalseFalseFalseFalseFalse3FalseFalseFalseFalseFalseFalseFalseFalseFalseFalseFalse4FalseFalseFalseFalseFalseFalseFalseFalseFalseFalse4FalseFalseFalseFalseFalseFalseFalseFalseFalse5FalseFalseFalseFalseFalseFalseFalseFalse6FalseFalseFalseFalseFalseFalseFalse7FalseFalseFalseFalseFalseFalseFalse8FalseFalseFalseFalseFalseFalseFalse992FalseFalseFalseFalseFalseFalseFalseFalse993FalseFalseFalseFalseFalseFalseFalseFalse993FalseFalseFalseFalseFalseFalseFalseFalse
n [164 ut[164 n [165	data.columns Index(['Ship Mode', 'Segment', 'Country', 'City', 'State', 'Postal Code',
	Segment 3 Country 1 City 531 State 49 Postal Code 631 Region 4 Category 3 Sub-Category 17 Sales 5825 Quantity 14 Discount 12 Profit 7287 dtype: int64 Exploratory Data Analysis In order to clean our data we should drop redundant and irrelevant columns. Here, we should drop Country and Postal Code
n [166 n [167 n [168	newdata = data.drop(['Country', 'Postal Code'], axis=1) newdata['Cost'] = newdata['Sales'] - newdata['Profit'] newdata['Profit%'] = (newdata['Profit']/newdata['Cost'])*100 newdata.head() Ship Mode Segment City State Region Category Sub-Category Sales Quantity Discount Profit Cost Profit% 0 Second Class Consumer Henderson Kentucky South Furniture Bookcases 261.9600 2 0.00 41.9136 220.0464 19.047619
n [169	1 Second Class Consumer Henderson Kentucky South Furniture Chairs 731.9400 3 0.00 219.5820 512.3580 42.857143 2 Second Class Corporate Los Angeles California West Office Supplies Labels 14.6200 2 0.00 6.8714 7.7486 88.679245 3 Standard Class Consumer Fort Lauderdale Florida South Furniture Tables 957.5775 5 0.45 -383.0310 1340.6085 -28.571429 4 Standard Class Consumer Fort Lauderdale Florida South Office Supplies Storage 22.3680 2 0.20 2.5164 19.8516 12.676056 Analyzing the relationship plt.figure(figsize=(8,5)) sns.kdeplot(data['Sales'], color='yellow', label='Sales', shade=True, bw=0.5) sns.kdeplot(data['Profit'], color='Pink', label='Profit', shade=True, bw=0.5) plt.xlim([-100, 1000]) plt.legend()
ut[169	<pre><matplotlib.legend.legend 0x21591d25340="" at=""></matplotlib.legend.legend></pre> 0.0030 0.0025 0.0020 0.0010
n [170 n [171	Heatmap It's used to find the correlation between column elements correlation = newdata.corr() sns.heatmap(newdata.corr(), cmap='viridis', annot=True) <axessubplot:></axessubplot:>
ut[171	Sales - 1
	Sales Quantity Discount Profit Cost Profit% Heatmap Inference Sales and profit are moderately correlated. Quantity and profit are moderately correlated. Discount and profit are negatively correlated. Dividing the dataset further into:prof(profit) and los(loss) prof=newdata.loc[newdata['Profit%']>0]
ut[173 n [174	Ship Mode Segment City State Region Category Sub-Category Sales Quantity Discount Profit Cost Profit Profit Cost Profit Profit Cost Profit Profit Cost Profit
n [175 ut[175 n [176	Ship Mode Segment City State Region Category Sub-Category Sales Quantity Discount Profit Cost Profit%
ut[176	Ship Mode Segment City State Region Category Sub-Category Sales Quantity Discount Profit Cost Profit Profit Cost Profit Pr
n [177 ut[177	Standard Class Consumer Costa Mesa California West Technology Phones 258.576 2 0.2 19.3932 239.1828 8.108108
	1 Standard Class Home Office Fort Worth Texas Central Office Supplies Appliances 68.8100 5 0.80 -123.8580 192.6680 -64.285714 2 Standard Class Home Office Fort Worth Texas Central Office Supplies Binders 2.5440 3 0.80 -3.8160 6.3600 -60.00000 3 Second Class Consumer Philadelphia Pennsylvania East Furniture Bookcases 308.34300 7 0.50 -1665.0522 4748.4822 -35.064935
n [178	1871 rows × 13 columns prof.shape (8058, 13) los.shape (1871, 13)
ut[180	Ship Mode Segment City State Region Category Office Supplies Paper 19.440 3 0.0 9.3312 10.1088 92.307692 Corporate Houston Texas Central Office Supplies Binders 7.184 2 0.2 5.4432 10.1088 53.846154 Corporate Los Angeles California West Office Supplies Binders 7.184 2 0.2 2.2450 4.9390 45.454545 Standard Class Consumer Dallas Texas Central Office Supplies Labels 4.928 2 0.2 1.7248 3.2032 53.846154 Second Class Corporate Jacksonville Florida South Office Supplies Labels 6.264 3 0.2 2.0358 4.2282 48.148148
n [181… ut[181…	1 1 Art 5.248 2 0.2 0.5904 4.6576 12.676056 1 Appliances 20.808 3 0.2 1.8207 18.9873 9.589041 1 Standard Class Home Office Yonkers New York East Technology Machines 52.440 4 0.0 24.1224 28.3176 85.185185 1 Length: 8014, dtype: int64 Ship Mode Segment City State Region Category Sub-Category Sales Quantity Discount Profit Cost Profit Key Profit Cost Pro
	Standard Class Consumer San Francisco California West Furniture Bookcases 359.499 3 0.15 -29.6058 389.1048 -7.8696 2 Second Class Corporate Chicago Illinois Central Office Supplies Binders 3.564 3 0.80 -6.2370 9.8010 -63.36364 2 Standard Class Home Office Columbus Ohio East Furniture Chairs 281.372 2 0.30 -12.0588 293.4308 -4.9589 2 Consumer Chicago Illinois Central Office Supplies Storage 35.168 2 0.20 -8.3524 43.5204 -19.91919 2 New York City New York East Furniture Bookcases 353.568 2 0.20 -44.1960 397.7640 -19.911111 2 Second Class Corporate Mcallen Texas Central Furniture Chairs 56.686 1 0.30 -20.2450 76.9310 -26.91780 -49.91780 1.188 2 0.70 -0.9900 2.1780 -49.91780 -49.9188 2 0.70 -0.9900 2.1780 -49.91888 2 0.70 -0.9900 2.1780 -0.9900
n [182 ut[182	Los Angeles California West Furniture Tables 1322.352 3 0.20 -99.1764 1421.5284 -6.6744 1 71.088 2 0.20 -1.7772 72.8652 -2.9024 1 Standard Class Home Office Yuma Arizona West Technology Machines 599.985 5 0.70 -479.9880 1079.9730 -44.4444 1 Length: 1865, dtype: int64 ltrs noteworthy that the business is overall yielding profits as the company is making profits in 8058 entries and losses in 1871 entries.Let's analyse this further. Prof['City'].value_counts() New York City 868
	San Francisco 482 Seattle 406 Philadelphia 283 Orland Park 1 Whittier 1 Holyoke 1 Richardson 1 Mason 1 Name: City, Length: 513, dtype: int64 los['City'].value_counts() Philadelphia 250 Houston 185
	Chicago 155 Dallas 73 Columbus 51 Harlingen 1 Champaign 1 North Las Vegas 1 Altoona 1 San Bernardino 1 Name: City, Length: 229, dtype: int64 We can hypothesize that the most profit is generated in New York City while the most loss is generated in Philadelphia. prof.groupby('Ship Mode')['Profit%'].mean().plot.bar() <axessubplot:xlabel='ship mode'=""></axessubplot:xlabel='ship>
n [186… ut[186…	los.groupby('Ship Mode')['Profit%'].mean().plot.bar() <axessubplot:xlabel='ship mode'=""></axessubplot:xlabel='ship>
	-15202530
n [187 ut[187	<pre>prof.groupby('Category')['Profit%'].mean().plot.bar() <axessubplot:xlabel='category'> 50 -</axessubplot:xlabel='category'></pre>
n [188 ut[188	los.groupby('Category')['Profit%'].mean().plot.bar() <axessubplot:xlabel='category'></axessubplot:xlabel='category'>
	$\begin{array}{cccccccccccccccccccccccccccccccccccc$
n [189 ut[189	We can see that the most profits as well as losses have been incurred in the Office Supplies Category. We need to analyse the sub-categories as well. los.groupby('Sub-Category')['Profit%'].mean().plot.bar() <axessubplot:xlabel='sub-category'></axessubplot:xlabel='sub-category'>
	Accessories Accessories And actions a series of the case of the
n [190… ut[190…	prof.groupby('Sub-Category')['Profit%'].mean().plot.bar() <axessubplot:xlabel='sub-category'> ### AxesSubplot:xlabel='Sub-Category'> ### AxesSubplot:xlabel='Sub-Category'></axessubplot:xlabel='sub-category'>
n [191	Appliances and Binders have suffered most losses while paper and labels have brought out the most profit. prof. groupby('Segment')['Profit%'].mean().plot.bar()
ut[191	<pre><axessubplot:xlabel='segment'> 50 - 40 - 20 - 10 - </axessubplot:xlabel='segment'></pre>
n [192 ut[192	los.groupby('Segment')['Profit%'].mean().plot.bar() <axessubplot:xlabel='segment'></axessubplot:xlabel='segment'>
	-1015202530 - We have the second of the s
n [194… ut[194…	<pre>prof.groupby('Region')['Profit%'].mean().plot.pie(autopct="%.1f%%") <axessubplot:ylabel='profit%'> East</axessubplot:ylabel='profit%'></pre>
n [198 n [199 ut[199	graph=los['Ship Mode'].value_counts() graph.plot.pie(autopct="%.1f%%") <axessubplot:ylabel='ship mode'=""> Standard Class</axessubplot:ylabel='ship>
n [200	prof.groupby('State')['Profit%'].mean()
	State Alabama 57.035981 Arizona 29.596523 Arkansas 66.037987 California 49.296384 Colorado 26.723697 Connecticut 59.535868 Delaware 61.867211 District of Columbia 74.969477 Florida 29.573372 Georgia 58.758111 Idaho 45.027213 Illinois 31.171962 Indiana 60.042150 Iowa 71.072453
	Kansas55.323261Kentucky56.021568Louisiana57.003986Maine60.659094Maryland60.162794Massachusetts60.291566Michigan56.013903Minnesota62.775354Mississippi60.65238Missouri59.800716Montana50.835528Nebraska55.741688Nevada55.741688New Hampshire61.055753New Jersey58.659983
	New Mexico 50.973856 New York 52.961044 North Carolina 27.811624 North Dakota 57.396117 Ohio 31.453295 Oklahoma 55.598073 Oregon 30.513179 Pennsylvania 31.261299 Rhode Island 58.975207 South Carolina 55.013330 South Dakota 65.246039 Tennessee 30.783915 Texas 33.678226 Utah 50.474091 Vermont 55.873440 Virginia 56.222189
n [201… ut[201…	Washington 47.318674 West Virginia 96.078431 Wisconsin 51.819017 Wyoming 6.666667 Name: Profit%, dtype: float64 los.groupby('State')['Profit%'].mean() State Arizona -30.265361 California -9.274100 Colorado -34.971742 Connecticut -6.650071 Delaware -14.621449
	Florida -29.034221 Illinois -40.464033 Maryland -8.492823 Massachusetts -13.895860 Nevada -13.978495 New Hampshire -9.090909 New Jersey -10.824408 New Mexico -19.191919 New York -17.066147 North Carolina -30.815065 Ohio -27.484667 Oregon -30.983764 Pennsylvania -26.750433 Rhode Island -10.264431 Tennessee -27.695211 Texas -38.113559 Washington -10.466661
	Washington -10.468661
	West Virginia -10.256410 Name: Profit%, dtype: float64 Most profits have come from West Virginia while the most losses have come from Illinois. Conclusion By dividing the data and analysing it from various perspectives using different fields, we can conclude that though our business is doing well, improvement is necessary in the Home Office Segment. Furthermore, we need to improve our reach in the southern region.