[164	from sklearn.tree import DecisionTreeClassifier from sklearn.model_selection import train_test_split from io import StringIO import pydotplus import seaborn as sns import matplotlib.image as mpimg from sklearn.tree import export_graphviz from sklearn.tree import DecisionTreeClassifier as DT from sklearn.metrics import accuracy_score from sklearn.metrics import confusion_matrix from sklearn import tree from IPython.display import Image Data Collection com_data=pd.read_csv('Company_Data.csv',sep=',') com_data
	Sales CompPrice Income Advertising Population Price ShelveLoc Age Education Urban US 0 9.50 138 73 11 276 120 Bad 42 17 Yes Yes 1 11.22 111 48 16 260 83 Good 65 10 Yes Yes 2 10.06 113 35 10 269 80 Medium 59 12 Yes Yes 3 7.40 117 100 4 4 466 97 Medium 55 14 Yes Yes 4 4.15 141 64 3 3 340 128 Bad 38 13 Yes No
[165 [166 [167	Com_data.isna().sum()
[168	8 Education 400 non-null int64 9 Urban 400 non-null object 10 US 400 non-null object dtypes: float64(1), int64(7), object(3) memory usage: 34.5+ KB Com_data.head() Sales CompPrice Income Advertising Population Price ShelveLoc Age Education Urban US 0 9.50 138 73 11 276 120 Bad 42 17 Yes Yes 1 11.22 111 48 16 260 83 Good 65 10 Yes Yes 2 10.06 113 35 10 269 80 Medium 59 12 Yes Yes
[169	3 7.40 117 100 4 466 97 Medium 55 14 Yes Yes 4 4.15 141 64 3 340 128 Bad 38 13 Yes No com_data.describe().T \[\begin{array}{c c c c c c c c c c c c c c c c c c c
[170	Advertising 400.0 6.655000 6.650364 0.0 0.0 0.5.00 12.00 29.00 Population 400.0 264.840000 147.376436 10.0 139.00 272.00 398.50 509.00 Price 400.0 115.795000 23.676664 24.0 100.00 117.00 131.00 191.00 Age 400.0 53.322500 16.200297 25.0 39.75 54.50 66.00 80.00 Education 400.0 13.900000 2.620528 10.0 12.00 14.00 16.00 18.00 Outliers Check import warnings warnings.filterwarnings('ignore') data1=sns.boxplot(com_data['Sales'])
[172	The data has 2 outliers plt.rcParams['figure.figsize']=9,4 plt.figure(figsize=(15,5)) print("skew: {}".format(com.data['Sales'].skew())) print("Kurtosis: {}".format(com.data['Sales'].kurtosis())) data1 = sns.kdeplot(com.data['Sales'].kurtosis())) plt.xticks([i for i in range(0,20,1]]) plt.show() Skew: 0.18556036318721578 Kurtosis: -0.08087736743346197
	0.12 - 0.10 - 0.08 - 0.06 - 0.04 - 0.02 - 0.02 - 0.02 - 0.02 - 0.02 - 0.03 - 0.04 - 0.02 - 0.04 - 0.
[174	Obj_colum = com_data.select_dtypes(include='object').columns.tolist() pll.figure(flgsize=(16,18)) for i.col in enumerate(obj_colum,1): plt.subjoc(2,2,1) sns.countplot(data=com.data,y=col) plt.subjoc(2,2,1) com_data(col).value_counts(normalize=True).plot.bar() plt.tight(2,0); plt.tight(2,0); plt.tight(2); condition of the first state of the
[176	No-
	Dit.show() num.data = com.data[num_columns] pd.DataFrame(data=[num.data.skew(), num.data.kurtosis()], index=['skewness', 'kurtosis']) 012 013 008 009 009 009 009 009 009 009 009 009
177	120
178 179 180 181	<pre>corr=com_data.corr() com_data = pd.get_dummies(com_data, columns = ['ShelveLoc', 'Urban', 'US']) corr=com_data.corr() plt.figure(figsize=(15,15)) sns.heatmap(corr,annot=True) <axessubplot:> Sales - 1 0064 015 027 005 4.044 4.023 4.052 0.39 05 4.074 0.015 4.018 0.18 CompPrice - 0.064 1 4.081 4.095 0.58 4.1 0.025 4.035 0.026 0.0087 4.067 0.067 0.017 0.017 Income - 0.15 4.081 1 0.059 4.0079 4.0077 4.0057 0.072 4.0013 4.038 0.038 4.099 0.09</axessubplot:></pre>
	Advertising - 027
182 183 185	Decision treeModel com_data["sales"]="small" com_data.loc[com_data["Sales"]>7.49, "sales"]="large" com_data.drop(["Sales"], axis=1, inplace=True) X = com_data.iloc[:,0:14] y = com_data.iloc[:,14] x_train,x_test,y_train,y_test = train_test_split(X,y,test_size = 0.2, stratify = y) y_train.value_counts() small 161
186 187 187	<pre>large 159 Name: sales, dtype: int64 model = DT(criterion='entropy') model.fit(x_train,y_train) DecisionTreeClassifier(criterion='entropy') pred_train = model.predict(x_train) accuracy_score(y_train,pred_train) 1.0</pre>
[189 [190	<pre>confusion_matrix(y_train, pred_train) array([[159, 0],</pre>
191 192 193	<pre>confusion_matrix(y_test,pred_test) array([[32, 8],</pre>
	80 large large 288 small small 265 small small 89 large large 28 small small 36 large large 196 small small 348 large large
[194 [195 [197 [198	248 small small 80 rows × 2 columns cols = list(com_data.columns) predictors = cols[0:14] target = cols[14] dot_data = StringIO() export_graphviz(model,out_file = dot_data, filled =True, rounded = True, feature_names =predictors, class_names = target, impurity = False) graph = pydotplus.graph_from_dot_data(dot_data.getvalue()) graph.write_png('company_full.png')
200	True img = mpimg.imread('company_full.png') plt.imshow(img) <pre> cmatplotlib.image.AxesImage at 0x2cefb3fe760> 000 1000 1000 1000 1000 1000 1000 1</pre>
[202 [203	0 500 1000 1500 2000 2500 3000 3500 model.feature_importances_
[]:	6 Education 0.020485 11 Urban_Yes 0.012885 9 ShelveLoc_Medium 0.012370 12 US_No 0.008609 7 ShelveLoc_Bad 0.000000 10 Urban_No 0.000000 13 US_Yes 0.000000