import numpy as np import pandas as pd import matplotlib.pyplot as plt import scipy.stats as stats import statsmodels.api as sm import statsmodels.formula.api as smf import seaborn as sns from sklearn.model selection import train test split from sklearn.preprocessing import StandardScaler from sklearn.neighbors import KNeighborsClassifier from sklearn import metrics **Predicting Market Value Using Linear Regression** ## Download the dataset GlenCove ## Upload the dataset using pandas In [4]: glen file = 'C:\\Users\\Joey\\Documents\\GlenCove.xlsx' GlenCove = pd.read excel(glen file) ## View the first 10 rows of the data using the head method print (GlenCove.head (n=10)) Address Fair Market Value(\$000) Property Size (acres) Age 0 9 Sycamore Road 522.9 0.2297 0.2192 61 21 Jefferson St 425.0 2 38 Hitching Post Lane 539.2 0.1630 39 4 Poppy Lane 628.2 0.4608 28 4 5 Daniel Drive 490.4 0.2549 56 0.2290 15 Francis Terrace 5 487.7 98 23 Guilfoy Street 6 370.3 0.1808 58 17 Carlyle Drive 777.9 0.5015 17 8 8 Craft Avenue 347.1 0.2229 62 9 22 Beechwood Ct. 0.1300 756.8 House Size (square feet) Rooms Baths Garage 0 2448 3.5 7 2.5 1 1942 2 2073 5 3.0 8 2.5 3 2707 4 2042 7 1.5 7 2.0 5 2089 1433 2.0 6 7 0 2991 9 2.5 8 1008 5 1.0 0 2.5 3202 ## Rename the columns GlenCove = GlenCove.rename(columns={'Fair Market Value(\$000)':'fair market value', 'Property Size (acres)':'pro In [9]: ## Recode garage values as binary and remove the age column garage dummies = pd.get dummies(GlenCove['Garage']).iloc[:, 1:] garage_dummies = garage_dummies[1] + garage_dummies[2] GlenCove final = GlenCove.drop('Age', 1) GlenCove_final['Garage'] = garage_dummies GlenCove final.head(n=10) Address fair_market_value property_size house_size Rooms Baths Garage 0 9 Sycamore Road 522.9 0.2297 2448 7 3.5 1 21 Jefferson St 425.0 0.2192 1942 2 38 Hitching Post Lane 539.2 2073 0.1630 3.0 1 3 4 Poppy Lane 628.2 0.4608 2707 5 487.7 0.2290 2089 7 2.0 15 Francis Terrace 23 Guilfoy Street 6 370.3 0.1808 1433 7 2.0 0 7 17 Carlyle Drive 777.9 0.5015 2991 2.5 8 8 Craft Avenue 347.1 0.2229 1008 5 1.0 0 9 22 Beechwood Ct. 756.8 0.1300 3202 8 2.5 In [11]: ## Print the summary statistics for all the numerical columns GlenCove_numeric = GlenCove_final[['fair_market_value', 'property_size', 'house_size', 'Rooms', 'Baths']] GlenCove_sum = GlenCove_numeric.describe() print(GlenCove_sum) fair_market_value property_size house_size Rooms Baths
20 000000 30.000000 30.000000 30.000000
30.000000 30.000000 30.000000 30.000000 30.000000 0.275960 0.236445 count 550.875234 1.487496 0.694808 144.312669 std 0.085200 1008.000000 5.000000 1.000000 310.200000 min 374.675000 0.152650 1622.250000 6.000000 1.625000 50% 431.200000 0.204050 1992.000000 7.000000 2.000000 0.270350 2205.500000 7.750000 2.500000 1.310000 3202 000000 11 000000 2.500000 75% 535.125000 889.000000 1.310000 3202.000000 11.000000 3.500000 In [13]: ### Print the correlation matrix for property_size, house size, Rooms, and Baths In [14]: print(GlenCove final[['property size','house size','Rooms','Baths']].corr()) property size house size Rooms Baths property_size 1.000000 0.211664 0.052963 0.053408 0.211664 1.000000 0.399573 0.521314 house_size

 0.052963
 0.399573
 1.000000
 0.133457

 0.053408
 0.521314
 0.133457
 1.000000

 Rooms Baths In [15]: | ## Type out the direction and strength of the correlation between ## each pair of variables in your matrix ## Property Size: ## weak positive correlations with house size, Rooms, and Baths ## weak positive correlations with property_size and Rooms, ## moderate positive correlation with Baths ## weak positive correlations with property size, house size, and Baths ## Baths: ## weak positive correlations with property_size and Rooms, ## moderate positive correlation with house size ## Using the statsmodels library, fit the data to a linear regression model In [18]: results = smf.ols('fair_market_value ~ property_size + house_size + Rooms + Baths + Garage', data = GlenCove fi ## Print the regression results summary In [19]: print(results.summary()) OLS Regression Results ______ Dep. Variable: fair_market_value R-squared: 0.815 OLS Adj. R-squared:
Least Squares F-statistic:
Tue, 18 May 2021 Prob (F-statistic): Model: 0.776 21.09 Method: 4.52e-08 Date: 01:46:51 Log-Likelihood: -165.9430 AIC: 343.9 No. Observations: Df Residuals: 24 BIC: 352.3 Df Model: 5 Covariance Type: nonrobust coef std err t P>|t| [0.025 0.975] Intercept 78.0523 67.779 1.152 0.261 -61.837 217.942 property_size 359.3104 56.919 6.313 0.000 241.836 476.785 house_size 0.1025 0.033 3.117 0.005 0.035 0.170 Rooms 0.7385 9.667 0.076 0.940 -19.214 20.691 Baths 36.8445 21.904 1.682 0.106 -8.363 82.052 Garage 28.3802 31.716 0.895 0.380 -37.078 93.839 ______ 4.370 Durbin-Watson: 1.619 Omnibus: 0.112 Jarque-Bera (JB): 2.767 Prob(Omnibus): 0.569 Prob(JB): 0.251 Skew: 3.958 Cond. No. Kurtosis: 1.15e+04 ______ Notes: [1] Standard Errors assume that the covariance matrix of the errors is correctly specified. [2] The condition number is large, 1.15e+04. This might indicate that there are strong multicollinearity or other numerical problems. In [21]: ## Create and print an ANOVA table In [22]: aov_table = sm.stats.anova_lm(results, typ=2) print(aov table) sum_sq df F PR(>F)
property_size 185956.980096 1.0 39.850226 0.000002 house_size 45335.829171 1.0 9.715382 0.004693 Rooms 27.229339 1.0 0.005835 0.939743 Baths 13203.326340 1.0 2.829448 0.105510 13203.326340 1.0 2.829448 0.105518

 Baths
 13203.326340
 1.0
 2.829448

 Garage
 3736.437085
 1.0
 0.800711

 Residual
 111993.529961
 24.0
 NaN

 3736.437085 1.0 0.800711 0.379768 In [23]: ## Explain whether or not you think this is a good model and why. ## You must use at least 2 measurements from your model to explain. In [24]: | ## I believe this is a decent model. According to the R-squared value, the model accounts for about 82% of ## the variability in fair market value, and the p-value for the F-statistic is close to 0, indicating that ##the results are likely significant. However, when looking through the p-values for each independent ## variable, it seems that only two (property_size and house_size) are significant, and it may be worthwile ## to consider removing these insignificant variables. (this is all assuming a significance level of 0.05) In [25]: | ## Type out the regression equation with the coefficients from the model output In [26]: $\#\#\ Y = 78.0523 + (359.3104)$ property size $+\ (0.1025)$ house size $+\ (0.7385)$ Rooms $+\ (36.8445)$ Baths ## + (28.3802) Garage In [27]: | ## Predict the fair market value of a house that is 2000 sqft, on 0.30 acres, with 6 rooms, 2 bathrooms, and prediction = results.predict(exog=dict(house size = 2000, property size = 0.30, Rooms = 6, Baths = 2, Garage = prediction_string = prediction.to_string() print(prediction_string[:11], "Thousand Dollars") 0 497.40 Thousand Dollars Plotting the Distribution of NBA Franchise Values In [29]: ## Download the dataset NBAValues.xlsx In [30]: ## Upload the dataset into your jupyter notebook using pandas In [31]: nba_file = 'C:\\Users\\Joey\\Documents\\NBAValues.xlsx' nba = pd.read excel(nba file) In [32]: | ## Display the entire dataframe without using the print function In [33]: nba Team Name Team Code Revenue (\$mil) Current Value (\$mil) Atlanta Hawks 1 **Boston Celtics** BOS 173 1700 2 **Brooklyn Nets** BKN 212 1500 3 **Charlotte Bobcats** CHA 130 725 4 Chicago Bulls CHI 201 2000 915 5 Cleveland Cavaliers CLE 149 6 Dallas Mavericks DAL 168 1150 7 **Denver Nuggets** DEN 855 136 8 **Detroit Pistons** DET 144 810 9 Golden State Warriors **GSW** 168 1300 10 **Houston Rockets** HOU 175 1250 11 Indiana Pacers IND 149 830 12 Los Angeles Clippers LAC 146 1600 13 Los Angeles Lakers LAL 293 2600 14 Memphis Grizzlies MEM 135 750 15 Miami Heat MIA 188 1175 Milwaukee Bucks 16 MIL 110 600 17 Minnesota Timberwolves MIN 128 625 18 **New Orleans Pelicans** NOH 131 650 19 New York Knicks NYK 278 2500 Oklahoma City Thunder 20 OKC 152 930 21 Orlando Magic ORL 143 875 22 Philadelphia 76ers PHI 125 700 23 **Phoenix Suns** PHX 145 910 Portland Trail Blazers 24 POR 153 940 25 Sacramento Kings SAC 125 800 San Antonio Spurs 26 SAS 174 1000 27 **Toronto Raptors TOR** 920 151 28 Utah Jazz UTA 142 830 29 Washington Wizards WAS 143 900 In [34]: ## Create at least 3 histograms of the Current Value column and change the number of bins and the color ## for each graph. ## Each histogram should have an x-axis label, y-axis label, and title. The axis labels should have size ## 10 font and the title should be size 16. fig, ax = plt.subplots() plt.hist(nba['Current Value (\$mil) '], bins = 4, color = 'powderblue') plt.title('Value of NBA Teams', fontsize = 16) ax.set_xlabel('Value (Millions of Dollars)', fontsize = 10) ax.set_ylabel('Number of Teams', fontsize = 10) plt.xticks(np.arange(600, 2550, 500)) Out[36]: ([<matplotlib.axis.XTick at 0x29072233fa0>, <matplotlib.axis.XTick at 0x29072233f70>, <matplotlib.axis.XTick at 0x29071b90910>, <matplotlib.axis.XTick at 0x290722842e0>], [Text(0, 0, ''), Text(0, 0, ''), Text(0, 0, ''), Text(0, 0, '')]) Value of NBA Teams 20.0 17.5 15.0 Number of Teams 12.5 10.0 7.5 5.0 2.5 0.0 600 1100 1600 2100 Value (Millions of Dollars) fig, ax = plt.subplots() plt.hist(nba['Current Value (\$mil) '], bins = 5, color = 'indianred') plt.title('Value of NBA Teams', fontsize = 16) ax.set xlabel('Value (Millions of Dollars)', fontsize = 10) ax.set ylabel('Number of Teams', fontsize = 10) plt.xticks(np.arange(600, 2600, 400)) Out[37]: ([<matplotlib.axis.XTick at 0x29072305d60>, <matplotlib.axis.XTick at 0x29072305d30>, <matplotlib.axis.XTick at 0x290722fa9d0>, <matplotlib.axis.XTick at 0x2907233a370>, <matplotlib.axis.XTick at 0x2907233a880>], [Text(0, 0, ''), Text(0, 0, ''), Text(0, 0, ''), Text(0, 0, ''), Text(0, 0, '')]) Value of NBA Teams 17.5 15.0 Number of Teams 12.5 10.0 7.5 5.0 2.5 1000 1400 1800 2200 600 Value (Millions of Dollars) fig, ax = plt.subplots() plt.hist(nba['Current Value (\$mil) '], bins = 10, color = 'gainsboro') plt.title('Value of NBA Teams', fontsize = 16) ax.set xlabel('Value (Millions of Dollars)', fontsize = 10) ax.set ylabel('Number of Teams', fontsize = 10) plt.xticks(np.arange(600, 2600, 200)) Out[38]: ([<matplotlib.axis.XTick at 0x2907236d4c0>, <matplotlib.axis.XTick at 0x2907236d490>, <matplotlib.axis.XTick at 0x29072367490>, <matplotlib.axis.XTick at 0x290723aa8b0>, <matplotlib.axis.XTick at 0x290723aadc0>, <matplotlib.axis.XTick at 0x290723b3310>, <matplotlib.axis.XTick at 0x290723b3820>, <matplotlib.axis.XTick at 0x290723b3d30>, <matplotlib.axis.XTick at 0x290723bb280>, <matplotlib.axis.XTick at 0x290723bb790>], [Text(0, 0, ''),Text(0, 0, ''), Text(0, 0, '')]) Value of NBA Teams 12 10 Number of Teams 6 2 800 1000 1200 1400 1600 1800 2000 2200 2400 Value (Millions of Dollars) ## Specify what you think is the best number of bins for this dataset ## The best number of bins is 5 as it limits the number of white space while also providing a smooth In [40]: ## and detailed distribution curve. ## Describe the shape of the distribution (e.g. normal, bimodal, right-skewed, left-skewed, etc.) In [42]: ## This is a right-skewed distribution as the right side (tail) is longer than the left side. Plotting the Distribution of MLB Franchise Values ## Download the dataset baseball values.csv In [44]: ## Upload the dataset using pandas mlb file = 'C:\\Users\\Joey\\Documents\\baseball values.csv' In [45]: mlb = pd.read csv(mlb file) ## Create a boxplot of the value column. Include an x-axis label, title, and change the color In [46]: print(mlb) Team Revenue Value 0 Baltimore 245 1 370 2100 Boston 227 2 975 Chicago White Sox 825 Cleveland 207 3 254 1125 Detroit 5 175 Houston 231 700 6 Kansas City 304 7 1300 Los Angeles Angels 8 Minnesota 223 895 9 New York Yankees 508 3200 10 202 Oakland 725 250 11 1100 Seattle Tampa Bay 188 625 13 266 1220 Texas 226 14 Toronto 870 15 211 840 Arizona 16 Atlanta 267 1150 17 Chicago Cubs 302 1800 227 18 Cincinnati 885 19 Colorado 214 855 20 Los Angeles Dodgers 403 2400 21 188 650 Miami 226 22 Milwaukee 875 263 23 New York Mets 1350 24 Philadelphia 265 1250 25 Pittsburgh 229 1400 26 St. Louis 294 27 225 San Diego 890 San Francisco 387 2000 Washington 287 1280 mlb_plot = sns.boxplot(x = 'Value', data = mlb, color = 'palegoldenrod') mlb_plot.set_title('Value of MLB Teams') mlb_plot.set_xlabel('Value (Millions of Dollars)') Out[47]: Text(0.5, 0, 'Value (Millions of Dollars)') Value of MLB Teams 500 1000 Value (Millions of Dollars) In [48]: ## Change the orientation of the boxplot so it is displayed vertically instead of horizontally. ## (seaborn has since changed syntax for vertical boxplots, need to check documentation to fix) mlb plot = sns.boxplot(x = 'Value', data = mlb, color = 'palegoldenrod', orient = 'v') In [49]: mlb plot.set title('Value of MLB Teams') mlb plot.set xlabel('Value (Millions of Dollars)') C:\Users\Joey\anaconda3\lib\site-packages\seaborn_core.py:1303: UserWarning: Vertical orientation ignored with only `x` specified. warnings.warn(single_var_warning.format("Vertical", "x")) Out[49]: Text(0.5, 0, 'Value (Millions of Dollars)') Value of MLB Teams 1000 1500 2000 2500 500 3000 Value (Millions of Dollars) ## Adjust an optional argument so that the whiskers extend 3X the IQR mlb_plot = sns.boxplot(x = 'Value', data = mlb, color = 'palegoldenrod', whis = 3) mlb_plot.set_title('Value of MLB Teams') mlb_plot.set_xlabel('Value (Millions of Dollars)') Out[51]: Text(0.5, 0, 'Value (Millions of Dollars)') Value of MLB Teams 500 1000 1500 2000 2500 Value (Millions of Dollars) ## Label the outlier with the team name displayed in red above the point $mlb_plot = sns.boxplot(x = 'Value', data = mlb, color = 'palegoldenrod', whis = 3)$ mlb_plot.set_title('Value of MLB Teams') mlb_plot.set_xlabel('Value (Millions of Dollars)') mlb plot.text(x = 3245, y = -.02, s = 'The New York Yankees', color = 'tab:red') Out[53]: Text(3245, -0.02, 'The New York Yankees') Value of MLB Teams The New York Yankees 500 1000 1500 2000 2500 3000 Value (Millions of Dollars) ## Create a violin plot of the value column. Include a title and x-axis label. In [54]: mlb_violin = sns.violinplot(x = 'Value', data = mlb, color = 'lightcyan') mlb_violin.set_title('Value of MLB Teams') mlb_violin.set_xlabel('Value (Millions of Dollars)') Out[55]: Text(0.5, 0, 'Value (Millions of Dollars)') Value of MLB Teams 500 1000 1500 2000 2500 3000 3500 Value (Millions of Dollars) ## Type the approximate mode of the dataset ## The approximate mode is 990 Using k-NN to Classify Types of Fruit ## Download the dataset fruit.tsv ## Upload the dataset using a pandas function that can read in tabular ## separated values fruit_file = 'C:\\Users\\Joey\\Documents\\fruit.tsv' fruit = pd.read_table(fruit_file) ## Display the first 5 rows of the data fruit.head() fruit_label fruit_name fruit_subtype mass width height color_score 0 granny_smith 192 8.4 7.3 0.55 apple apple granny_smith 180 8.0 6.8 0.59 2 granny_smith apple 7.4 7.2 0.60 3 mandarin mandarin 6.2 4.7 0.80 2 mandarin mandarin 84 6.0 4.6 0.79 ## Determine how many observations of each fruit are contained in the dataset In [64]: num_apple = (fruit.fruit_label == 1).sum() num_mandarin = (fruit.fruit_label == 2).sum() num orange = (fruit.fruit label == 3).sum() num_lemon = (fruit.fruit_label == 4).sum() print(num_apple, 'apples,', num_mandarin, 'mandarins,', num_orange, 'oranges,', num_lemon, 'lemons') 19 apples, 5 mandarins, 19 oranges, 16 lemons ## Subset the dataframe and create a new dataframe with just the columns 'fruit name', 'mass', 'width', ## 'height', and 'color score' fruit_final = fruit[['fruit_name', 'mass', 'width', 'height', 'color_score']] fruit_final.head() mass width height color_score fruit_name 0 apple 192 8.4 7.3 0.55 1 180 8.0 6.8 0.59 apple 2 apple 176 7.4 7.2 0.60 3 mandarin 4.7 0.80 86 6.2 0.79 mandarin 6.0 4.6 84 ## Create a scatterplot matrix using the new dataframe and color the dots based on fruit name. ## (Be sure to include title and legend) fruit_matrix = sns.pairplot(fruit_final, hue = 'fruit_name') fruit_matrix.fig.suptitle("Scatterplot Matrix: Fruit Name", y = 1.05, fontsize = 18) Out[70]: Text(0.5, 1.05, 'Scatterplot Matrix: Fruit Name') Scatterplot Matrix: Fruit Name 350 300 250 200 150 100 width fruit_name apple mandarin orange 10 lemon 6 0.9 0.8 0.7 0.6 200 10 0.6 8 width mass height color_score ## From the scatterplot matrix determine whether or not you think mass and height would be could predictor ## variables for the fruit name and leave your answer in a brief comment or markdown cell ## I do not think that mass and height would be good predictor variables, as it seems as though there is a lot ## of overlap in the different types of fruit near the middle of the distribution. ## Assign fruit name as your dependent variable and the rest of the variables in your new dataframe as your ## independent variable x = fruit_final.iloc[:, 1:4].values y = fruit_final.iloc[:, 0].values ## Split the data into training and test splits (70% for training and 30% for testing). Set the random state ## to 0. x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.3, random_state = 0) ## Print out the summary statistics for the independent variables in your training set print(pd.DataFrame(x_train).describe()) 0 1 41.000000 41.000000 41.000000 count mean 160.878049 7.058537 7.682927 53.607926 0.842608 1.236710 std 80.000000 5.800000 4.300000 min 130.000000 25% 6.200000 158.000000 50% 7.200000 7.600000 75% 172.000000 7.500000 8.200000 356.000000 9.200000 10.500000 max ## Print out the summary statistics for the independent variables in your test set print(pd.DataFrame(x_test).describe()) 0 1 18.000000 18.000000 18.000000 count 168.222222 mean 7.211111 7.716667 59.380682 0.767646 76.000000 5.800000 144.000000 7.100000 7.200000 25% 50% 162.000000 7.200000 7.550000 179.000000 7.375000 8.175000 75% 362.000000 9.600000 10.300000 In [81]: ## Create a K-NN classifier model with 5NN classifier = KNeighborsClassifier(n_neighbors=5,) ## Standardize data (if necessary) ## scaler = StandardScaler() ## scaler.fit(x_train) ## x_train = scaler.transform(x_train) ## x_test = scaler.transform(x_test) In [83]: ## Initialize the classifier classifier.fit(x train, y train) Out[83]: KNeighborsClassifier() In [84]: ## Use your test set to make predictions y_pred = classifier.predict(x_test) y_pred Out[85]: array(['orange', 'apple', 'lemon', 'lemon', 'apple', 'apple', 'lemon', 'lemon', 'apple', 'lemon', 'mandarin', 'apple', 'orange', 'apple', 'lemon', 'apple', 'mandarin'], dtype=object) ## Print confusion matrix conf = metrics.confusion_matrix(y_test, y_pred) print(conf) [[3 1 0 0] [1 2 0 0] [0 0 2 0] [4 3 0 2]] In [88]: ## Print the model accuracy rate In [89]: metrics.accuracy_score(y_test, y_pred) Out[89]: 0.5