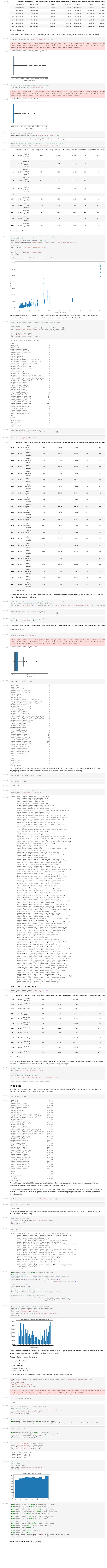
In [1]: #libraries import pandas as pd import numpy as np import seaborn as sns import matplotlib.pyplot as plt In [2]: #updating number of output results pd.set option('display.max rows', 500) **Business Questions** 1. Which regions of the United States have the hottest housing markets? 2. Is there a cutoff for the market hotness index where a significant difference between markets becomes apparent? 3. Does the allure of the certain cities (favorable place to live) play into the cities' assigned market hotness index? 4. Can we predict the market hotness index of cities on a monthly/yearly basis? **Reading Data** Realtor.com Historical Monthly Inventory Data per Metro Area Monthly data updated on March 31, 2022 with data through March 2022. Next update scheduled for May 10, 2022 with data through April 2022. Market trends and monthly statistics on active for-sale listings (including median list price, average list price, luxury list price, median days on market, average days on market, total active listings, new listings, price increases, price reductions). Attribution: cite any full or partial use of the data to the 'realtor.com residential listings database.' In [3]: #core metrics csv file core history = pd.read csv("C:/Users/phill/OneDrive/Desktop/Capstone DSC680/Data/RDC Inventory Core Metrics Met C:\Users\phill\anaconda3\lib\site-packages\IPython\core\interactiveshell.py:3444: DtypeWarning: Columns (0,1,3 8,39) have mixed types. Specify dtype option on import or set low memory=False. exec(code obj, self.user global ns, self.user ns) In [4]: core history.head() month_date_yyyymm cbsa_code Out[4]: cbsa_title HouseholdRank median_listing_price median_listing_price_mm median_listing_price_yy active_li new yorknewark-0 202203 35620 1.0 699000.0 0.0029 0.0779 jersey city, ny-nj-pa los angeles-202203 31080 2.0 949950.0 0.0231 -0.0496 long beachanaheim, ca chicagonaperville-2 202203 16980 3.0 337000.0 0.0307 -0.0369 elgin, il-indallas-fort worth-3 202203 19100 4.0 425000.0 0.0391 0.1759 arlington, tx houston-26420 woodlands-373733.0 0.0310 4 202203 5.0 0.0945 sugar land, tx 5 rows × 41 columns In [5]: core history.month date yyyymm.unique() array([202203, 202202, 202201, 202112, 202111, 202110, 202109, 202108, Out[5]: 202107, 202106, 202105, 202104, 202103, 202102, 202101, 202012, 202011, 202010, 202009, 202008, 202007, 202006, 202005, 202004, 202003, 202002, 202001, 201912, 201911, 201910, 201909, 201908, 201907, 201906, 201905, 201904, 201903, 201902, 201901, 201812, 201811, 201810, 201809, 201808, 201807, 201806, 201805, 201804, 201803, 201802, 201801, 201712, 201711, 201710, '201710', '201709', '201708', '201707', '201706', '201705', '201704', '201703', '201702', '201701', '201612', '201611', '201610', '201609', '201608', '201607', 'quality flag = 1: year-over-year figures may be impacted'], dtype=object) For the variables which are suffixed with "mm" or "yy", they represent the percentage change in the certain data value from the previous month (mm) or previous year (yy). Therefore, they will allow for understanding trends in the values over time. Realtor.com Market Hotness per Metro Area Realtor.com Market Hotness Index: scores and rankings based on days on market (supply index) and realtor.com views per property (demand index). Attribution: cite any full or partial use of the data to the 'realtor.com market hotness index'. In [6]: #hotness metrics csv file hot history = pd.read csv("C:/Users/phill/OneDrive/Desktop/Capstone DSC680/Data/RDC Inventory Hotness Metrics N In [7]: hot history.head() cbsa_title nielsen_hh_rank hotness_rank hotness_rank_mm hotness_rank_yy hotness_score supply_sco Out[7]: month_date_yyyymm cbsa_code las vegas-0 -10.0 202203 29820 33.0 173.0 -4.0 45.652174 87.9598 hendersonparadise, nv elizabethtown-1 202203 21060 283.0 -3.0 69.0 -17.0 67.892977 77.5919 fort knox, ky warner robins, 2 202203 47580 229.0 254.0 -27.0 -122.0 25.919732 27.0903 harrisburg-3 56.020067 202203 25420 91.0 117.0 31.0 -4.0 45.4849 carlisle, pa oklahoma city, 4 202203 36420 41.0 246.0 18.0 -11.0 27.257525 41.1371 ok 5 rows × 24 columns In [8]: hot_history.month_date_yyyymm.unique() array(['202203', '202202', '202201', '202112', '202111', '202110', Out[8]: '202109', '202108', '202107', '202106', '202105', '202104', '202103', '202102', '202101', '202012', '202011', '202010', '202009', '202008', '202007', '202006', '202005', '202004', '202003', '202002', '202001', '201912', '201911', '201910', '201909', '201908', '201907', '201906', '201905', '201904', '201903', '201902', '201901', '201812', '201811', '201810', '201809', '201808', '201807', '201806', '201805', '201804', '201803', '201802', '201801', '201712', '201711', '201710', '201709', '201708', 'quality flag = 1'], dtype=object) Simplemaps: World Cities Database Information on U.S. cities, includes: City Name State Name County Name Latitude Longitude Population Density Timezone In [9]: #us cities information cities pd = pd.read csv("C:/Users/phill/OneDrive/Desktop/Capstone DSC680/Data/uscities.csv") In [10]: cities_pd.head() city city_ascii state_id state_name county_fips county_name Out[10]: lat source military incorporated population density New New York 36061 New York 40.6943 -73.9249 18713220 10715 polygon New York NY False True Los Los CA California 6037 Los Angeles 34.1139 -118.4068 12750807 3276 polygon False True Angeles Angeles Illinois 17031 IL Cook 41.8373 -87.6862 8604203 4574 polygon **2** Chicago Chicago False True Miami Miami FL 12086 Miami-Dade 25.7839 -80.2102 6445545 5019 polygon Florida False True Dallas TX 48113 Dallas 32.7936 -96.7662 5743938 Dallas Texas 1526 polygon False True **Data Cleansing Monthly Inventory Data** In [11]: print("Shape: ", core history.shape) Shape: (63274, 41) In [12]: core_history.dtypes month_date_yyyymm object Out[12]: cbsa_code object cbsa_title object HouseholdRank float64 median listing_price float64 median_listing_price_mm float64 median_listing_price_yy float64 active_listing_count float64 active_listing_count_mm float64 active_listing_count_yy float64 median_days_on_market float64 median_days_on_market_mm float64 median_days_on_market_yy float64 new_listing_count float64 new_listing_count_mm float64 new_listing_count_yy float64 price increased count float64 price_increased_count_mm float64 price_increased_count_yy float64 price_reduced count float64 price_reduced_count_mm float64 price_reduced_count_yy float64 pending_listing_count float64 pending_listing_count_mm float64 pending_listing_count_yy float64 float64 median_listing_price_per_square_foot median_listing_price_per_square_foot_yy float64 median_square_feet float64 median_square_feet_mm float64 median square feet yy float64 average_listing_price float64 average_listing_price_mm float64 average_listing_price_yy float64 total listing_count float64 total listing count mm float64 total_listing_count_yy float64 pending_ratio float64 pending_ratio_mm object pending_ratio_yy object quality_flag float64 dtype: object I am going to convert the variables 'pending_ratio_mm' and 'pending_ratio_yy' to numeric types, since they are currently defined as objects which does not represent the data. In [13]: core history[core history['pending ratio mm']=='#NAME?'] Out[13]: month_date_yyyymm cbsa_code cbsa_title HouseholdRank median_listing_price median_listing_price_mm median_listing_price_yy activ guymon, 202102 1.9179 12811 25100 891.0 350000.0 2.3333 1 rows × 41 columns In [14]: core history[core history['pending ratio yy']=='#NAME?'] Out[14]: month_date_yyyymm cbsa_code cbsa_title HouseholdRank median_listing_price median_listing_price_mm median_listing_price_yy active guymon, 375000.0 2724 202201 25100 891.0 0.3636 2.5714 1 rows × 41 columns In [15]: #replace arbitrary value in pending ratio mm column with NaN core history['pending ratio mm'] = core history['pending ratio mm'].replace({'#NAME?': np.nan}) In [16]: #replace arbitrary value in pending ratio yy column with NaN core_history['pending_ratio_yy'] = core_history['pending_ratio_yy'].replace({'#NAME?': np.nan}) In [17]: #convert pending ratio mm to float64 core history['pending ratio mm'] = core history['pending ratio mm'].astype(float, errors = 'raise') In [18]: #convert pending_ratio_yy to float64 core_history['pending_ratio_yy'] = core_history['pending_ratio_yy'].astype(float, errors = 'raise') In [19]: core history.describe() Out[19]: HouseholdRank median_listing_price median_listing_price_mm median_listing_price_yy active_listing_count active_listing_count_mm active_listing_count_m 52269.000000 63273.000000 6.327300e+04 52269.000000 63273.000000 52268.000000 count mean 459.000000 2.511767e+05 0.007882 0.081275 1047.051381 -0.016822 std 264.717033 1.715830e+05 0.067474 0.158750 3613.741968 0.237067 0.000000 1.000000 1.990000e+04 -0.692500 -0.843400 -1.000000 min 25% 230.000000 1.499000e+05 -0.018200 0.000000 122.000000 -0.075900 -0.019800 50% 459.000000 2.099000e+05 0.000400 0.064100 271.000000 **75**% 688.000000 2.946500e+05 0.028800 0.140600 707.000000 0.032400 917.000000 2.985500 max 3.000000e+06 5.600000 85761.000000 48.000000 8 rows × 38 columns In [20]: #inventory data # Rows containing duplicate data -- none! duplicate rows core = core history[core history.duplicated()] print("number of duplicate rows: ", duplicate rows core.shape) # Finding the null values print("\nNull values") print(core history.isnull().sum()) number of duplicate rows: (0, 41) Null values month date_yyyymm 0 cbsa code 0 cbsa title 1 HouseholdRank 1 median listing price 1 median listing price mm 11005 median listing price yy 11005 active listing count 1 active listing count mm 11006 11006 active listing count yy median days on market 1 11005 median days on market mm median days on market yy 11005 new listing count 1 new listing count mm 11297 11343 new listing count yy price increased count price increased count mm 37735 price increased count yy 37435 price reduced count price reduced count mm 12080 price reduced count yy 11829 pending_listing_count 2708 pending listing count mm 13398 pending listing count yy 14674 median_listing_price_per_square_foot 18 median_listing_price_per_square_foot_mm median listing price per square foot yy 11025 median square feet 18 11011 ${\tt median \ square \ feet \ mm}$ median square feet yy 11025 average listing price average listing price mm 11005 average listing price yy 11005 total listing count total listing count mm 11005 total listing count yy 11005 pending ratio 2708 pending ratio mm 13127 pending ratio yy 14194 quality flag 11005 dtype: int64 In [21]: #row 63273 seems to have more informational data than pertaining to the dataset #i am going to remove it #last row in dataframe --> tail 1 record core history = core history.drop(core history.tail(1).index) In [22]: #looking at rows with missing data core_history[core_history.isnull().any(axis=1)] Out[22]: month_date_yyyymm cbsa_code cbsa_title HouseholdRank median_listing_price median_listing_price_mm median_listing_price_yy acti new 35300 62 202203 haven-63.0 354500.0 0.0292 0.1936 milford, ct peoria, il 129 202203 37900 130.0 110950.0 0.1321 0.0577 anchorage, 139 202203 11260 140.0 425000.0 0.0572 0.1863 utica-161 202203 46540 162.0 183950.0 0.0856 0.0773 rome, ny duluth, 0.1013 0.0989 162 202203 20260 163.0 249950.0 mn-wi 63268 201607 46900 vernon, tx 913.0 139000.0 NaN NaN 63269 201607 18780 914.0 175000.0 NaN NaN craig, co 101900.0 63270 201607 29500 lamesa, tx 915.0 NaN NaN 63271 201607 49820 zapata, tx 916.0 128000.0 NaN NaN 169900.0 63272 201607 37780 917.0 NaN NaN pecos, tx 42673 rows × 41 columns Many of the missing values in this dataset pertain to the variables suffixed with "mm" or "yy", which as explained above indicate the percentage change in the variable from the previous month or year respectively. If I remove the entire row that has missing values, then we will lose all of the data for the specific metro area which I am not particularly interested in doing. Therefore, I am leaning towards defining a default value for the missing data and then in visualizations they can be subsetted out to look at trends for other cities. Boxplot and Histogram of median_listing_price_mm to assess how to fill missing values In [23]: #trying to decide on missing values to fill in sns.boxplot(core_history.median_listing_price_mm) #many outliers #main distribution around 0 C:\Users\phill\anaconda3\lib\site-packages\seaborn_decorators.py:36: FutureWarning: Pass the following variabl e 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(<AxesSubplot:xlabel='median_listing_price_mm'> Out[23]: 0.0 0.5 1.0 1.5 3.0 median_listing_price_mm In [24]: plt.hist(core_history['median_listing_price_mm']) plt.title("Histogram of median_listing_price_mm") plt.show() Histogram of median_listing_price_mm 40000 30000 20000 10000 0 In [25]: core history['median listing price yy'].mode() 0.0 Out[25]: dtype: float64 In [26]: core history['median listing price yy'].median() 0.0641 Out[26]: I am going to come back to handling the missing values, since I am conflicted on whether it will disrupt the solution of my business questions. I also think it will be better to handle missing data after I merge all of the datasets, so we can have a better picture. I do not want to naively fill in data that will lead to a misunderstanding of what each metro area brings to the issue. **Market Hotness** In [27]: print("Shape: ", hot_history.shape) Shape: (16801, 24) In [28]: hot_history.dtypes month date yyyymm object Out[28]: cbsa code object cbsa title object nielsen hh rank float64 hotness rank float64 hotness rank mm float64 hotness rank yy float64 hotness score float64 supply_score float64 demand score float64 median days on market float64 median days on market mm float64 median dom mm day float64 median_days_on_market_yy float64 median_dom_yy_day float64 median dom vs us float64 ldp_unique_viewers_per_property_mmfloat64ldp_unique_viewers_per_property_yyfloat64ldp_unique_viewers_per_property_vs_usfloat64 median_listing_price float64 median_listing_price_mm float64 median_listing_price_yy float64 median_listing_price_vs_us float64 quality_flag float64 dtype: object In [29]: #market hotness data # Rows containing duplicate data -- none! duplicate rows hot = hot history[hot history.duplicated()] print("number of duplicate rows: ", duplicate rows hot.shape) # Finding the null values print("\nNull values") print(hot history.isnull().sum()) number of duplicate rows: (0, 24) Null values month date_yyyymm 0 cbsa code 0 cbsa title 1 nielsen hh rank 1 hotness rank 1 hotness rank mm 3601 hotness rank yy 3601 hotness score 1 supply score 1 demand score 1 median days on market 1 3601 median days on market mm 3601 median dom mm day median_days_on_market_yy 3601 median dom yy day 3601 median dom vs us 1 ldp_unique_viewers_per_property_mm 3601
ldp_unique_viewers_per_property_yy 3601 median listing price median listing price mm 3601 median listing price yy 3601 median listing_price_vs_us 1 quality flag dtype: int64 In [30]: #looking at rows with missing data hot history[hot history.isnull().any(axis=1)] cbsa_title nielsen_hh_rank hotness_rank hotness_rank_mm hotness_rank_yy hotness_score suppl Out[30]: month_date_yyyymm cbsa_code youngstownwarren-13200 201807 49660 94.0 166.0 NaN NaN 45.652174 boardman, oh-201807 21060 elizabethtown-283.0 61.705686 fort knox, ky spartanburg, 43900 13202 201807 159.0 165.0 NaN NaN 45.986622 39 gulfportbiloxi-13203 201807 269.0 10 25060 133.0 NaN NaN 16.889632 pascagoula, new haven-201807 35300 46.822742 45 13204 63.0 158.0 NaN NaN milford, ct 16796 201708 20260 197.0 NaN 37.625418 15 duluth, mn-wi 163.0 NaN 201708 st. louis, mo-il 19.0 51.505017 16797 41180 136.0 NaN NaN 56 16798 201708 17900 67.0 columbia, sc 249.0 NaN NaN 24.414716 41 traverse city, 16799 201708 45900 261.0 223.0 NaN NaN 32.107023 25 Data falls outside of 16800 $quality_flag = 1$ expected NaN NaN NaN NaN NaN NaN range. Please ... 3601 rows × 24 columns In [31]: #row 16800 seems to have more informational data than pertaining to the dataset #i am going to remove it #last row in dataframe --> tail 1 record hot history = hot history.drop(hot history.tail(1).index) Cities Database In [32]: print(cities_pd.shape) (28338, 17)In [33]: cities_pd.dtypes city object Out[33]: city_ascii object state id object state name object county_fips int64 county_name object float64 lng float64 population int64 density int64 source object bool military incorporated bool timezone object ranking int64 zips object int64 dtype: object In [34]: #looking at unique values of some categorical variables cities pd.source.unique() array(['polygon'], dtype=object) Out[34]: In [35]: cities pd.incorporated.unique() array([True, False]) Out[35]: In [36]: #bar chart of city rankings #captures the importance of a city --> from 1-5 cities pd.ranking.value counts().sort values().plot(kind = 'barh') <AxesSubplot:> Out[36]: 2 1 5000 10000 15000 20000 25000 I am going to drop some of the columns which are not of interest to me for the business problem/questions: city_ascii (duplicate information) county fips source (only has one value) timezone (not relevant) zips id In [37]: #dropping above columns cities pd = cities pd.drop(['city ascii','county fips','source','timezone','zips','id'],1) cities pd C:\Users\phill\AppData\Local\Temp/ipykernel 25828/4130845510.py:2: FutureWarning: In a future version of pandas all arguments of DataFrame.drop except for the argument 'labels' will be keyword-only cities pd = cities pd.drop(['city ascii', 'county fips', 'source', 'timezone', 'zips', 'id'],1) Out[37]: Ing population density military incorporated ranking city state_id state_name county_name lat New York -73.9249 New York NY New York 40.6943 18713220 10715 False True 1 Los Angeles California Los Angeles 34.1139 -118.4068 12750807 3276 False True Chicago Illinois Cook 41.8373 -87.6862 8604203 4574 False True Miami-Dade 25.7839 -80.2102 Miami Florida 6445545 5019 True 4 Dallas TX Texas Dallas 32.7936 -96.7662 5743938 1526 False True 1 28333 Gross ΝE Nebraska Boyd 42.9461 -98.5697 2 False True 3 28334 Lotsee Oklahoma Tulsa 36.1334 -96.2091 39 True 28335 The Ranch MN Minnesota Mahnomen 47.3198 -95.6952 2 False True 3 28336 Shamrock Oklahoma Creek 35.9113 False 3 OK -96.5772 True Nebraska -98.3296 28337 Monowi NE Boyd 42.8307 False True 3 28338 rows × 11 columns In [38]: print('Shape of Cities Dataset after Column Removal') print(cities pd.shape) Shape of Cities Dataset after Column Removal In [39]: #city data # Rows containing duplicate data -- none! duplicate rows city = cities pd[cities pd.duplicated()] print("number of duplicate rows: ", duplicate rows city.shape) # Finding the null values -- none! print("\nNull values") print(cities pd.isnull().sum()) number of duplicate rows: (0, 11) Null values city state id 0 0 state name 0 county_name 0 population density military 0 incorporated ranking 0 dtype: int64 The cities database dataset is very clean in terms of duplicates and missing data, of which it has zero! Merging Realtor.com datasets The inventory dataset and the market hotness dataset are both from Realtor.com. They have a few categorical columns in common, and for the purpose of joining, I am going to use the cbsa_code column since it should be the same between the two datasets. I was originally going to use cbsa_title but it is more prone to give incorrect results, since strings can contain different spacing and punctuations. I am going to do an inner join between the DataFrames, so we can capture information for the metro areas that have information from both datasets, instead of having null values with an outer join which will not be useful. This may subset the dataset but it will be more advantageous to the problem at hand. I ended up joining the datasets in Tableau and then exporting that dataset to CSV to read into Python, since the merge using Python functions was not working as expected. In [40]: #merged in external source --> Tableau #python merge was giving incorrect results #reading in merged dataset merged realtor = pd.read csv("C:/Users/phill/OneDrive/Desktop/Capstone DSC680/Data/joined realtor data.csv") In [41]: merged realtor.shape (16800, 60)Out[41]: In [42]: merged realtor.head() Out[42]: Median Active Active Price Price Average **Average** Month Median Active Price Price Listing Cbsa Dom **Listing Listing** Reduced Reduced Qu **Dom Listing** Listing Date **Cbsa Title** Increased Reduced Code Mm Count Count Price Count Count Yyyymm Yy Day Count Price **Count Yy** Count Day Mm Mm Mm Yy new yorknewark-202203 35620 -17.0 -12.5 42816 0.0564 -0.0807 1700187 0.0001 ... 0.3333 5968 0.1477 -0.0744 jersey city, ny-nj-pa angeles-202203 31080 0.0703 -0.2700 2411312 -0.0253 lona beachanaheim, ca chicagonaperville-202203 16980 -6.0 -8.0 11706 0.0225 -0.2399 528104 0.0197 ... -0.1183 2644 0.1338 -0.2205 elgin, il-indallas-fort worth-202203 19100 -6.0 -12.5 0.0891 -0.2105 660783 0.0312 ... -0.4571 1056 0.3679 -0.1456 arlington, houston-202203 26420 woodlands--3.5 -12.0 10603 -0.0384 -0.1722 542520 0.0218 ... -0.1067 2728 0.0912 -0.0870 sugar land, 5 rows × 60 columns In [43]: merged realtor['Month Date Yyyymm'].unique() array([202203, 202202, 202201, 202112, 202111, 202110, 202109, 202108, Out[43]: 202107, 202106, 202105, 202104, 202103, 202102, 202101, 202012, 202011, 202010, 202009, 202008, 202007, 202006, 202005, 202004, 202003, 202002, 202001, 201912, 201911, 201910, 201909, 201908, 201907, 201906, 201905, 201904, 201903, 201902, 201901, 201812, 201811, 201810, 201809, 201808, 201807, 201806, 201805, 201804, 201803, 201802, 201801, 201712, 201711, 201710, 201709, 201708], dtype=int64) In [44]: #merged data # Rows containing duplicate data -- none! duplicate rows merge = merged realtor[merged realtor.duplicated()] print("number of duplicate rows: ", duplicate rows merge.shape) # Finding the null values -- none! print("\nNull values") print(merged realtor.isnull().sum())

number of duplicate rows: (0, 60) Null values Month Date Yyyymm Cbsa Code 0 Cbsa Title 0 Median Dom Mm Day 3600 Median Dom Yy Day 3600 Active Listing Count 0 Active Listing Count Mm Active Listing Count Yy Average Listing Price Average Listing Price Mm Average Listing Price Yy 0 Demand Score Hotness Rank 0 Hotness Rank Mm 3600 Hotness Rank Yy 3600 Hotness Score Household Rank Ldp Unique Viewers Per Property Mm 3600 Ldp Unique Viewers Per Property Vs Us 0 Ldp Unique Viewers Per Property Yy 3600 Median Days On Market 0 Median Days On Market Mm median days on market mm hotness 3600 Median Days On Market Yy 0 median days on market yy hotness 3600 Median Dom Vs Us 0 Median Listing Price Median Listing Price Mm 3600 median listing price mm hotness Median Listing Price Per Square Foot Median Listing Price Per Square Foot Mm Median Listing Price Per Square Foot Yy Median Listing Price Vs Us Median Listing Price Yy median listing price yy hotness 3600 Median Square Feet Median Square Feet Mm Median Square Feet Yy New Listing Count New Listing Count Mm New Listing Count Yy Nielsen Hh Rank 0 Pending Listing Count 164 Pending Listing Count Mm 225 Pending Listing Count Yy 356 Pending Ratio 164 Pending Ratio Mm 210 Pending Ratio Yy 314 Price Increased Count Price Increased Count Mm 1553 Price Increased Count Yy 1479 Price Reduced Count Price Reduced Count Mm Price Reduced Count Yy Quality Flag quality flag hotness Supply Score Total Listing Count Total Listing Count Mm 0 Total Listing Count Yy dtype: int64 In [45]: #looking at rows with missing data merged realtor[merged realtor.isnull().any(axis=1)] Out[45]: Median Active Active **Price Price** Average Price Median Active Average Month Price Cbsa Dom **Listing Listing** Listing Reduced Reduced Listing Date **Dom Listing Cbsa Title** ... Increased Reduced Code Mm Count Count Price Count Count Yyyymm Yy Day Count **Price** Count Yy Count Mm new 202203 35300 haven--20.5 -0.0191 -0.6658 514659 0.0036 -0.50 0.1786 -0.0294 milford, ct -0.0766 -0.4562 202203 37900 0.0000 -0.4286 129 peoria, il -9.0208356 0.1836 -1.00144 anchorage, 202203 11260 0.0176 -0.4627 0.2308 139 -24.5 -20.5 478697 0.0645 ... 0.00 -0.4839 utica-161 202203 46540 -28.0 -0.0493 -0.0816 267986 0.0762 ... NaN 0.0000 -0.1429 rome, ny duluth, 202203 20260 -25.5 385764 162 -26.0 287 -0.0304 -0.3051 0.0570 ... NaN 60 0.8750 0.2500 mn-wi bismarck, 0.0346 0.0802 16795 201708 13900 NaN NaN 318638 0.0055 ... NaN 224 0.0980 0.0769 springfield, 0.2840 -0.0287 16796 201708 44220 NaN 152726 -0.0080 ... 0.40 248 0.1698 0.0508 battle 0.5417 16797 201708 12980 NaN NaN 555 0.0221 173958 0.0137 ... NaN 200 0.1364 1.3810 creek, mi lebanon, 16798 201708 30140 -0.0081 -0.2337 236837 0.0028 ... 0.25 192 0.3714 -0.0400ра 16799 201708 46180 tupelo, ms NaN NaN 532 0.0019 0.0114 213864 -0.0260 ... NaN 188 0.0930 1.3500 5679 rows × 60 columns Metropolitan Wikipedia In [46]: #libraries import requests # library to handle requests from bs4 import BeautifulSoup # library to parse HTML documents In [47]: # get the response in the form of html wikiurl="https://en.wikipedia.org/wiki/Metropolitan_statistical_area" table_class="wikitable sortable jquery-tablesorter" response=requests.get(wikiurl) print(response.status_code) 200 In [48]: # parse data from the html into a beautifulsoup object soup = BeautifulSoup(response.text, 'html.parser') metrotable=soup.find('table', {'class':"wikitable"}) In [49]: df=pd.read html(str(metrotable)) # convert list to dataframe df=pd.DataFrame(df[0]) print(df.head()) print(df.tail()) Metropolitan statistical area 2021 estimate Rank 19768458 0 New York-Newark-Jersey City, NY-NJ-PA MSA 1 Los Angeles-Long Beach-Anaheim, CA MSA 12997353 2 Chicago-Naperville-Elgin, IL-IN-WI MSA 9509934 Dallas-Fort Worth-Arlington, TX MSA 3 7759615 4 Houston-The Woodlands-Sugar Land, TX MSA 2020 Census % change Encompassing combined statistical area 0 20140470 -1.85% New York-Newark, NY-NJ-CT-PA CSA -1.54% 1 13200998 Los Angeles-Long Beach, CA CSA 2 Chicago-Naperville, IL-IN-WI CSA 9618502 -1.13% Dallas-Fort Worth, TX-OK CSA 3 7637387 +1.60% 7122240 +1.19% Houston-The Woodlands, TX CSA Rank Metropolitan statistical area 2021 estimate 2020 Census % change \ 73095 379 Danville, IL MSA 74188 -1.47% 380 381 Lewiston, ID-WA MSA 64851 64375 +0.74% 382 381 Walla Walla, WA MSA 62584 +0.16% 62682 62846 -1.46% Enid, OK MSA 382 383 61926 58639 +0.60% 383 Carson City, NV MSA 58993 384 Encompassing combined statistical area 379 380 Kennewick-Richland-Walla Walla, WA CSA 381 382 Reno-Carson City-Fernley, NV CSA 383 There seem to be control characters in the columns with spaces, so I am going to work to remove them so they are easier to work with and In [50]: df.columns Index(['Rank', 'Metropolitan statistical area', '2021 estimate', '2020 Census', Out[50]: '% change', 'Encompassing combined statistical area'], dtype='object') In [51]: #updating column names with underscores df.columns = [c.replace(' ', ' ') for c in df.columns] In [52]: df.columns = df.columns.str.strip() In [53]: df.columns[5] 'Encompassing_combined_statistical\xa0area' Out[53]: In [54]: #lower case text in column metropolitan area df['Metropolitan\xa0statistical\xa0area'] = df['Metropolitan\xa0statistical\xa0area'].str.lower() In [55]: #replacing metropolitan area name 'MSA' with '' In [56]: #lower case text in column statistical area df['Encompassing combined statistical\xa0area'] = df['Encompassing combined statistical\xa0area'].str.lower() In [57]: #replacing statistical area name 'CSA' with '' df['Encompassing combined statistical\xa0area'] = df['Encompassing combined statistical\xa0area'].str.replace(1 In [58]: #renaming column names with control characters df = df.rename(columns={'Metropolitan\xa0statistical\xa0area': 'metro area', 'Encompassing combined statistical In [59]: df.head() Out[59]: metro_area 2021_estimate 2020_Census %_change combined_statistical_area 0 1 new york-newark-jersey city, ny-nj-pa 19768458 20140470 –1.85% new york-newark, ny-nj-ct-pa 1 los angeles-long beach-anaheim, ca 12997353 13200998 -1.54% los angeles-long beach, ca 2 3 chicago-naperville-elgin, il-in-wi 9509934 9618502 -1.13% chicago-naperville, il-in-wi 3 7759615 7637387 +1.60% dallas-fort worth, tx-ok dallas-fort worth-arlington, tx 5 houston-the woodlands-sugar land, tx 7206841 7122240 +1.19% houston-the woodlands, tx In [60]: df.dtypes Out[60]: metro_area int64 object 2021 estimate int64 2020 Census int64 % change object combined statistical area object dtype: object In [61]: df.describe() Out[61]: Rank 2021_estimate 2020_Census **count** 384.000000 3.840000e+02 3.840000e+02 mean 192.500000 7.460229e+05 7.450639e+05 **std** 110.995495 1.655492e+06 1.670046e+06 5.899300e+04 5.863900e+04 min 1.000000 25% 96.750000 1.476232e+05 1.474762e+05 **50%** 192.500000 2.487975e+05 2.484905e+05 **75%** 288.250000 6.012120e+05 5.956365e+05 **max** 384.000000 1.976846e+07 2.014047e+07 The variable '%_change' is currently of object type, but I'd like it to be converted to a float since it is a numerical measure. Therefore, I am going to work on correcting that. In [62]: print(df['%_change'][0][1:-1]) 1.85 In [63]: #remove first and last characters from the value #first char is + or - currently #last char is % df['% change'] = df['% change'].str[1:-1] In [64]: df['% change'].head(20) 1.85 Out[64]: 1.54 2 1.13 3 1.60 4 1.19 5 0.45 6 0.26 7 0.89 8 0.76 9 2.07 10 0.84 11 1.16 12 2.65 13 0.61 14 0.18 15 0.01 16 0.38 17 1.39 18 0.30 19 0.22 Name: %_change, dtype: object In [65]: #cast object type to float type #convert df['%_change'] to float64 df['% change'] = df['% change'].astype(float, errors = 'raise') In [66]: #checking column types after update df.dtypes int64 Rank Out[66]: metro area object 2021 estimate int64 2020 Census int64 % change float64 combined statistical area object dtype: object Merging Joined Realor Dataset and Metro Population Dataset The Metro Population dataset contains population information per metro area for the year of 2020 and then also provides an estimate for 2021. The merged realtor dataset has information for dates ranging from 2017-2022, so in order to join the two datasets, I am going to subset the realtor one to only focus on the dates in the years of 2020 and 2021 in order to make the metro dataset relevant. In [67]: #convert int month date column to date merged realtor['Month Date Yyyymm'] = pd.to datetime(merged realtor['Month Date Yyyymm'].astype(str), format=' In [68]: #creating new columns for month and year in merged dataset merged realtor['year'] = pd.DatetimeIndex(merged realtor['Month Date Yyyymm']).year In [69]: merged realtor['year'].unique() array([2022, 2021, 2020, 2019, 2018, 2017], dtype=int64) Out[69]: In [70]: #month column merged realtor['month'] = pd.DatetimeIndex(merged realtor['Month Date Yyyymm']).month In [71]: merged realtor['month'].unique() array([3, 2, 1, 12, 11, 10, 9, 8, 7, 6, 5, 4], dtype=int64) Out[71]: In [72]: merged realtor.head() Out[72]: Median **Price** Active Active Price Average Month Median Active **Average** Quality Cbsa Dom Listing Listing Listing Reduced Reduced quality **Cbsa Title** Date Dom Listing Listing Code Mm Count Count Price Count Count Flag flag_hotness **Yyyymm** Yy Day Count Price Day Mm Yy Mm Mm Yy new york-2022-03newark-35620 42816 0.0564 -0.0807 1700187 0.0001 -17.0 -12.5 0.1477 -0.07440 2 jersey city, ny-nj-pa los angeles-2022-03-31080 long -0.5 -5.0 0.0703 -0.2700 2411312 -0.0253 ... 0.1951 -0.3438 0 7 01 beachanaheim, ca chicago-2022-03naperville-16980 -6.0 -8.0 11706 0.0225 -0.2399 528104 0.0197 ... 0.1338 -0.2205 0 5 elgin, il-indallas-fort 2022-03worth-0.0891 -0.2105 0.0312 ... 0.3679 19100 -6.0 -12.5 4524 660783 -0.14560 0 7 01 arlington, tx houstonthe 2022-03-26420 woodlands--3.5 -12.0 10603 -0.0384 -0.1722 542520 0.0218 ... 0.0912 -0.0870 0 3 sugar land, 5 rows × 62 columns In [73]: #subsetting merged dataset for years of 2021 and 2020 realtor years = merged realtor[(merged realtor.year == 2020) | (merged realtor.year == 2021)] In [74]: realtor years.shape #less than half of the original, not too bad! (7200, 62)Out[74]: In [75]: realtor years['year'].unique() array([2021, 2020], dtype=int64) Out[75]: In [76]: #I am going to drop the month date column now since we have individual 'year' and 'month realtor years = realtor years.drop(['Month Date Yyyymm'],1) C:\Users\phill\AppData\Local\Temp/ipykernel 25828/4283728241.py:2: FutureWarning: In a future version of pandas all arguments of DataFrame.drop except for the argument 'labels' will be keyword-only realtor years = realtor years.drop(['Month Date Yyyymm'],1) In [77]: #cleansing strings for first dataset realtor years['Cbsa Title'] = realtor years['Cbsa Title'].apply(lambda x: x.lower().strip()) In [78]: #cleansing strings for second dataset df['metro_area'] = df['metro_area'].apply(lambda x: x.lower().strip()) In [79]: #merge datasets on specified key columns merged metro = pd.merge(realtor years, df, left on='Cbsa Title', right on='metro area') In [80]: merged_metro.head() Out[80]: Median Active Active **Total Total** Average Median Active **Average** Average Cbsa Cbsa Dom Listing Listing Listing Listing Listing year month Rank metro_ard Dom Listing Listing Listing **Title** Code Mm Count Count Price Count Count Yy Day Count **Price Price Yy** Mm Day Mm Yy Mm Yy new yorknew yor newarknewar 35620 -7.0 46763 -0.1382 -0.1824 1631077 0.0473 0.2380 ... -0.0940 -0.1408 2021 12 11.0 jersey jersey cit city, nyny-nj-l nj-pa new yorknew yor newarknewar **1** 35620 4.0 -6.0 54263 -0.0524 -0.1738 1557436 0.0819 0.2129 ... -0.0390 -0.1636 2021 11 jersey jersey cit city, nyny-njnj-pa new yorknew yor newarknewar **2** 35620 57262 0.0060 -0.1657 1439479 10 4.0 0.0207 0.1067 ... -0.0105 -0.1574 2021 jersey cit jersey city, nyny-nj-l nj-pa new yorknew yor newarknewar **3** 35620 56921 -0.0022 -0.1521 1410300 0.0249 0.1403 ... -0.0026 -0.1175 2021 jersey jersey cit city, nyny-nj-լ nj-pa new yorknew yor newarknewar **4** 35620 -1.0 57049 0.0366 -0.1444 1376032 0.0776 0.1214 ... 0.0406 -0.0924 2021 jersey cit jersey city, nyny-njnj-pa 5 rows × 67 columns In [81]: merged_metro.shape (6192, 67)Out[81]: In [82]: merged metro.columns Index(['Cbsa Code', 'Cbsa Title', 'Median Dom Mm Day', 'Median Dom Yy Day', 'Active Listing Count', 'Active Listing Count Mm', 'Active Listing Count Yy', 'Average Listing Price', 'Average Listing Price Mm', 'Average Listing Price Yy', 'Demand Score', 'Hotness Rank', 'Hotness Rank Mm', 'Hotness Rank Yy', 'Hotness Score', 'Household Rank', 'Ldp Unique Viewers Per Property Mm', 'Ldp Unique Viewers Per Property Vs Us', 'Ldp Unique Viewers Per Property Yy', 'Median Days On Market', 'Median Days On Market Mm', 'median days on market mm hotness', 'Median Days On Market Yy', 'median days on market yy hotness', 'Median Dom Vs Us', 'Median Listing Price', 'Median Listing Price Mm', 'median listing price mm hotness', 'Median Listing Price Per Square Foot', 'Median Listing Price Per Square Foot Mm', 'Median Listing Price Per Square Foot Yy', 'Median Listing Price Vs Us', 'Median Listing Price Yy', 'median listing price yy hotness', 'Median Square Feet', 'Median Square Feet Mm', 'Median Square Feet Yy', 'New Listing Count', 'New Listing Count Mm', 'New Listing Count Yy', 'Nielsen Hh Rank', 'Pending Listing Count', 'Pending Listing Count Mm', 'Pending Listing Count Yy', 'Pending Ratio', 'Pending Ratio Mm', 'Pending Ratio Yy', 'Price Increased Count', 'Price Increased Count Mm', 'Price Increased Count Yy', 'Price Reduced Count', 'Price Reduced Count Mm', 'Price Reduced Count Yy', 'Quality Flag', 'quality flag hotness', 'Supply Score', 'Total Listing Count', 'Total Listing Count Mm', 'Total Listing Count Yy', 'year', 'month', 'Rank', 'metro area', '2021 estimate', '2020 Census', '% change', 'combined statistical area'], dtype='object') In [83]: #creating correlation plot for all columns corr = merged metro.corr() corr.style.background_gradient(cmap='coolwarm') Out[83]: Active **Average** Median Median Active Active **Average Average** Hotness Cbsa Listing Listing Demand Hotness Dom Mm Dom Yy Listing Listing Listing Listing Rank Code Count Price Score Rank Day Day Count **Count Yy** Price **Price Yy** Mm Mm Mm 0.007877 0.046950 **Cbsa Code** 1.000000 0.054846 0.021216 0.005642 -0.034003 0.095928 -0.107527 -0.010092 **Median Dom** 1.000000 0.231246 0.006812 -0.136367 0.150442 -0.065573 -0.035854 -0.508700 **Mm Day Median Dom** 0.097786 -0.208965 0.100125 0.054846 0.231246 1.000000 -0.063438 0.050753 -0.088590 -0.155278 0.426696 Yy Day **Active Listing** 0.100460 0.006812 0.097786 1.000000 0.012836 0.182266 -0.043520 -0.268661 0.170073 -0.018193 Count **Active Listing** -0.086072 -0.045502 -0.064947 -0.020810 -0.136367 0.199868 -0.265055 **Count Mm Active Listing** 0.150442 0.100460 0.199868 1.000000 -0.447008 -0.232489 -0.091095 0.426696 0.141616 **Count Yy Average** -0.001151 -0.063438 0.268598 -0.069405 0.182266 -0.004483 -0.057394 1.000000 0.066678 0.017333 0.020239 **Listing Price Average** 1.000000 0.053011 **Listing Price** -0.012544 -0.065753 0.050753 -0.003431 -0.087863 0.303435 -0.013252 0.033772 **Average** -0.043520 **Listing Price** -0.034003 -0.065573 -0.086072 0.268598 0.303435 1.000000 0.149194 -0.130189 0.043415 **Demand Score** -0.035854 -0.088590 -0.045502 -0.232489 0.053011 0.149194 1.000000 -0.798863 0.097320 0.141616 1.000000 **Hotness Rank** -0.1075270.088315 0.100125 0.170073 -0.064947 -0.069405 -0.798863 -0.153072 **Hotness Rank** -0.508700 -0.153072 -0.010092-0.155278 -0.018193 -0.020810 -0.091095 0.033772 0.043415 0.097320 1.000000 **Hotness Rank** 0.064536 -0.151377 -0.529468 -0.074959 -0.006487 0.053134 0.260651 0.278672 -0.342328 0.301847 Yy **Hotness Score** -0.090573 -0.104870 0.066947 -0.145804 0.128789 0.804682 -0.991518 0.143849 Household -0.035628 -0.008775 -0.075846 -0.409578 0.012034 0.046412 -0.175466 0.046874 0.025241 0.011304 0.185555 Rank **Ldp Unique** -0.542419 0.010274 **Viewers Per** 0.130309 0.214428 0.004356 -0.089725 0.162345 0.011132 0.050939 0.165906 **Property Mm Ldp Unique Viewers Per** -0.046203 -0.069325 -0.246850 -0.078881 -0.217639 0.064115 0.138881 0.938262 0.083548 **Property Vs Ldp Unique Viewers Per** -0.040640 -0.407494 -0.787839 0.442456 0.199903 -0.134179 0.116897 0.115644 **Property Yy Median Days** -0.055681 0.153948 0.375128 0.057866 -0.394274 0.158089 -0.09179² 0.092952 -0.208744 -0.075138 **On Market Median Days** -0.084586 On Market 0.006419 0.915465 0.227440 0.000950 -0.138538 0.192826 -0.026950 0.074087 -0.523592 median days -0.084588 on market 0.006425 0.915465 0.227440 0.000951 -0.138538 0.192827 -0.026952 -0.523588 mm_hotness **Median Days** 0.914344 0.085702 0.081933 -0.136499 0.211491 0.399349 -0.105659 0.195082 On Market Yy median days 0.039471 0.211490 0.914343 0.085701 0.399349 -0.071055 0.081933 -0.254413 -0.105658 0.195081 -0.136497 on market yy_hotness **Median Dom** 0.015453 0.0731440.134700 0.125705 -0.126537 0.042649 -0.057349 -0.068610 -0.269711 -0.096764 Vs Us Median -0.000953 -0.030097 0.156813 0.015526 -0.031635 0.919130 0.042457 0.228999 0.001898 -0.110128 0.014031 **Listing Price** Median 0.035158 -0.008548 **Listing Price** -0.118579 -0.143835 -0.053857 0.473054 0.135993 0.054248 -0.028000 0.033390 Mm median listing price -0.007371 -0.118584 0.035162 -0.008547 -0.143847 -0.053851 0.473054 0.135995 0.054246 -0.028005 0.033389 mm_hotness Median **Listing Price** 0.021313 0.152906 -0.043542 -0.044723 -0.028381 0.204490 0.007259 0.012228 **Per Square Foot** Median **Listing Price** 0.004120 -0.058250 0.026056 -0.144580 0.344013 0.179903 0.081734 -0.127338 -0.077201 **Per Square Foot Mm** Median **Listing Price** 0.029458 0.060341 0.078484 -0.231909 0.185810 0.055812 0.495553 **0.128551** -0.201243 -0.005849 **Per Square Foot Yy** Median **Listing Price** 0.015708 0.168915 -0.016053 -0.007117 0.905963 0.051258 0.209440 0.003536 -0.111637 0.014990 Vs Us Median -0.044651 -0.102801 **Listing Price** -0.138386 -0.279214 0.157659 0.161970 0.635176 **0.135945** -0.130773 0.027181 Yy median listing -0.038032 -0.138388 -0.102801 -0.279216 -0.044647 0.157657 0.161972 0.635177 0.135945 -0.130777 0.027183 price yy_hotness Median 0.034138 -0.074154 -0.032637 0.010214 -0.113210 -0.047661 -0.021086 -0.085054 0.127968 0.073567 0.211264 -0.079403 **Square Feet** Median 0.016458 **Square Feet** -0.080911 0.080324 -0.026324 -0.042983 0.047042 0.011052 0.018298 -0.001245 0.332617 Mm Median -0.090681 **Square Feet** -0.028831 -0.075087 0.085977 -0.187572 0.135332 0.278232 0.041044 0.043225 0.035975 Yy **New Listing** 0.021365 0.056507 0.887919 -0.023119 -0.024770 0.060034 0.033825 0.209199 -0.266800 0.099775 -0.012835 Count **New Listing** -0.017257 -0.010998 0.104086 0.011034 0.080912 0.067290 0.044812 -0.022198 0.162981 **Count Mm New Listing** -0.001761 -0.137752 0.322347 -0.016672 0.189905 0.031885 0.041410 -0.021557 0.054800 **Count Yy** Nielsen Hh 0.035628 -0.008775 -0.075846 -0.409578 0.012034 0.046412 -0.175466 0.046874 0.185555 0.025241 0.011304 Rank **Pending** 0.030454 0.047291 -0.003261 0.004717 0.014299 -0.008005 0.196401 -0.272360 0.144815 -0.018939 **Listing Count Pending** -0.067644 **Listing Count** 0.007307 -0.008041 -0.002620 -0.006661 0.083187 0.023728 -0.014792 0.007546 0.033243 Mm **Pending** -0.024914 **Listing Count** 0.010404 0.064077 -0.003654 -0.000967 0.006034 -0.076720 0.080178 -0.010482 Yy **Pending Ratio** -0.039802 0.036089 0.218745 0.173319 -0.180391 0.015799 **Pending Ratio** 0.031963 0.015592 -0.138065 0.140498 -0.006409 -0.594986 -0.189110 0.269536 0.088531 -0.003780 Mm **Pending Ratio** 0.040511 -0.124519 -0.087454 0.062802 0.341528 0.134929 -0.090424 0.021808 Yy Price -0.018363 -0.044697 0.046081 -0.025747 0.015770 0.145934 0.011738 -0.292049 -0.014912 Increased 0.574588 0.131598 Count **Price** 0.028635 -0.003667 0.000715 -0.031747 -0.037535 -0.075136 0.000942 0.022424 0.004554 Increased **Count Mm Price** -0.017513 0.057165 -0.075189 -0.121897 -0.054976 0.040652 Increased 0.077457 -0.005149 **Count Yy Price Reduced** 0.010009 0.111422 0.928124 0.036157 0.115221 0.152903 -0.067616 -0.276484 0.140865 -0.024554 Count **Price Reduced** -0.040304 -0.004978 -0.105423 -0.015387 -0.009680 0.437635 0.044552 0.024560 -0.054361 0.081486 **Count Mm Price Reduced** -0.008436 0.130472 0.113264 0.178245 -0.111669 **Count Yy** 0.045840 0.144675 0.147560 0.057001 **Quality Flag** 0.113044 -0.049454 0.098223 0.032747 -0.013107 -0.060395 quality 0.055840 -0.041734 0.055271 0.067298 -0.041791 -0.005640 0.154768 0.107987 0.068067 0.094711 -0.008980 flag_hotness **Supply Score** -0.109818 -0.078814 -0.012465 0.155106 0.002672 0.055309 0.275931 -0.784937 0.132956 **Total Listing** 0.006334 0.014002 -0.279728 0.081080 0.973679 0.059673 0.195915 -0.028395 0.165890 -0.019246 Count **Total Listing** 0.012469 -0.064242 0.014301 0.489837 0.056189 -0.010670 -0.017388 -0.060009 -0.019279 **Count Mm Total Listing** 0.109395 0.530479 0.081261 0.328752 0.059121 -0.054209 -0.164512 0.137017 -0.085880 **Count Yy** 0.040329 -0.366282 0.157289 -0.150586 0.119947 0.001844 0.000068 0.002249 year 0.291721 0.034892 -0.108106 0.001208 month -0.026378 -0.022161 0.055737 0.014983 -0.001172 -0.001076 -0.008937 -0.407659 0.010169 0.041034 0.192374 0.011082 -0.075443 -0.198316 0.029500 0.040731 Rank 0.007919 0.871316 0.008558 0.253533 -0.018532 2021_estimate 0.075072 0.038522 -0.230772 0.070427 -0.018271 0.039298 -0.229796 2020_Census 0.007894 0.075059 0.008757 0.254408 -0.018699 0.071103 -0.018364 0.018342 -0.003237 -0.052047 0.081583 -0.004444 0.003857 0.101061 0.042844 0.125086 -0.155897 0.021701 %_change In [84]: sns.boxplot(merged metro['Average Listing Price']) C:\Users\phill\anaconda3\lib\site-packages\seaborn\ decorators.py:36: FutureWarning: Pass the following variabl e 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(<AxesSubplot:xlabel='Average Listing Price'> Out[84]: Average Listing Price le6 In [85]: sns.boxplot(merged metro['Median Listing Price']) C:\Users\phill\anaconda3\lib\site-packages\seaborn_decorators.py:36: FutureWarning: Pass the following variabl e 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(<AxesSubplot:xlabel='Median Listing Price'> Out[85]: 2.0 0.0 0.5 1.0 1.5 Median Listing Price le6 There is a variable in both datasets called Quality Flag, which was created to indicate data values that are outside of the typical range ... they can be found with the "1" marker. I am going to remove any of these values from the dataset, since they can potentially represent outliers and skew the results of my reporting/modeling. In [86]: merged_metro = merged_metro[(merged_metro['Quality Flag']==0) | (merged_metro['quality flag_hotness']==0)] There are variables in both Realtor.com datasets which represent the same measures such as median days on market & median listing prices. I am going to remove one version of these variables, since they are highly correlated. The dataset also has median listing price and average listing price. In the case of different housing markets, there can be outliers among the metro areas since the listing prices depend on many factors such as location, square feet, year built, etc. Therefore, I am just going to keep the median listing price variables to avoid outlier skew. There is also Household Rank which is a rank indicator based on the household count for a metro area. This variable is highly correlated with the Rank from the Wikipedia population source; I am going to only keep the population one. I am also going to take out the Nielsen HH rank for the same. I am not very interested in the Pending Ratio variables, so I am going to remove them. In [87]: #dropping columns for reduction & simplicity merged metro = merged metro.drop(['metro area','combined statistical area','Hotness Score','Supply Score','Demo C:\Users\phill\AppData\Local\Temp/ipykernel 25828/1686688935.py:2: FutureWarning: In a future version of pandas all arguments of DataFrame.drop except for the argument 'labels' will be keyword-only merged metro = merged metro.drop(['metro_area','combined_statistical_area','Hotness Score','Supply Score','De mand Score', 'median days on market mm hotness', 'median days on market yy_hotness', 'median listing price mm_hotn ess','median listing price yy_hotness','Median Dom Mm Day','Median Dom Yy Day','Average Listing Price','Average Listing Price Mm', 'Average Listing Price Yy', 'Household Rank', 'Pending Ratio', 'Pending Ratio Mm', 'Pending Ratio' Yy','Nielsen Hh Rank','Quality Flag','quality flag_hotness'],1) In [88]: #renaming columns to replace spaces with underscores #updating column names with underscores merged metro.columns = [c.replace(' ', ' ') for c in merged metro.columns] In [89]: merged metro.head() Out[89]: Cbsa_Code Cbsa_Title Active_Listing_Count Active_Listing_Count_Mm Active_Listing_Count_Yy Hotness_Rank Hotness_Rank_Mm Hotness_Ra new yorknewark-0 35620 46763 -0.1382 -0.1824298 -1.0 jersey city, ny-nj-pa new yorknewark-1 35620 54263 -0.0524 -0.1738297 -6.0 jersey city, ny-nj-pa new yorknewark-2 57262 0.0060 -0.1657 291 -2.0 35620 jersey city, ny-nj-pa new yorknewark-3 35620 56921 -0.0022 -0.1521 289 6.0 jersey city, ny-nj-pa new yorknewark-35620 57049 0.0366 -0.1444295 3.0 iersev citv ny-nj-pa 5 rows × 46 columns In [90]: merged metro.dtypes #may need to convert month and year to date objects Out[90]: Cbsa_Code int64 Cbsa Title object Active Listing Count int64 Active Listing Count Mm float64 Active Listing Count Yy float64 Hotness Rank int64 Hotness Rank Mm float64 Hotness Rank Yy float64 Ldp Unique Viewers Per Property Mm float64 Ldp_Unique_Viewers_Per_Property_Vs_Us float64
Ldp_Unique_Viewers_Per_Property_Yy float64 Median Days On Market int64 Median Days On Market Mm float64 Median Days On Market Yy float64 Median Dom Vs Us float64 Median Listing Price int64 float64 Median Listing Price Mm Median_Listing_Price_Per_Square_Foot int64 Median Listing Price Per Square Foot Mm float64 Median Listing Price Per Square Foot Yy float64 Median Listing Price Vs Us float64 Median Listing Price Yy float64 Median Square Feet int64 Median Square Feet Mm float64 Median Square Feet Yy float64 New Listing Count int64 New Listing Count Mm float64 New Listing Count Yy float64 Pending Listing Count float64 Pending Listing Count Mm float64 Pending Listing Count Yy float64 Price Increased Count int64 Price Increased Count Mm float64 Price Increased Count Yy float64 Price Reduced Count int64 Price Reduced Count Mm float64 Price Reduced Count Yy float64 Total Listing Count int64 Total Listing Count Mm float64 Total Listing Count Yy float64 int64 month int64 Rank int64 2021 estimate int64 2020 Census int64 % change float64 dtype: object In [91]: merged metro['Hotness Rank'].describe() count 5317.000000 Out[91]: mean 154.785029 std 86.654318 1.000000 min 79.00000 25% 158.000000 50% 230.000000 75% max 300.000000 Name: Hotness Rank, dtype: float64 In [92]: merged metro.describe()



<pre># creating a confusion matrix Xmm predictions = Kmn.predict(X_test) Xmm_ms = confusion matrix(y_test, kmn_predictions) # fighting prediction and result xmm_mstations = prediction and result xmm_mstations = prediction and result xmm_mstations # confusion_mstations # confusion_mstations # confusion_mstations # confusion_mstation # training a Maive Payes classifier # training a Constantive Dayes import GaussianNB gph = CaussianNB().fit(X train, y_train) gph = creating a confusion matrix gph = creating a confusion matrix gph_accuracy = gph_score(X_test, y_test) print(gph_accuracy) # accuracy on X_test # accuracy on X_test # accuracy on A_test # accuracy</pre>		<pre># training a KNN classifier from sklearn.neighbors import KNeighborsClassifier knn = KNeighborsClassifier(n_neighbors = 7).fit(X_train, y_train) # accuracy on X_test knn_accuracy = knn.score(X_test, y_test) print(knn_accuracy)</pre>
# training a Naire Sayes classifier from sklearn.neive bayes import GaussianNB ghb = GaussianNB().fit(X_train, y_train) gnb_predictions = gnb.predict(X_test) # scoursey on X_test gnb_accuracy = gnb.acore(X_test, y_test) print(gnb_accuracy) # creating a confusion matrix gpb_on = confusion matrix(y_test, gnb_predictions) # getting precision and recall gnb_matrica = precision_recall_facore_support(y_test, gnb_predictions, average='macre',labels=mp.unique print(gnb_matrica) 0.00782122905027933 (0.00375458488041371, 0.03082725490196078, 0.007884959338371305, None) # Wultinomial Logistic Regression from sklearn.inear_model import logisticRegression from sklearn.model_selection import RepeatedStratifiedKFold # define the multimemial logistic regression moded		<pre>print(knn_accuracy) # creating a confusion matrix knn_predictions = knn.predict(X_test) knn_cm = confusion_matrix(y_test, knn_predictions) #getting precision and recall knn_metrics = precision_recall_fscore_support(y_test, knn_predictions, average='macro',labels=np.unique print(knn_metrics) 0.008938547486033519 (0.00953065134099617, 0.01958128078817734, 0.007683671692292383, None)</pre>
#getting precision and recall gnb_metrics = precision_recall_fscore_support(y_test, gnb_predictions, average='macro',labels=np.unique print(gnb_metrics) 0.00782122905027933 (0.005754585488041371, 0.03063725490196078, 0.007884959539371305, None) Multinomial Logistic Regression from sklearn.linear_model import LogisticRegression from sklearn.model_selection import RepeatedStratifiedKFold # define the multinomial logistic regression model model = LogisticRegression(multic_class='multinomial', solver='lbfgs') # define the model evaluation procedure cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1) # evaluate the model and collect the scores n_scores = cross_val_score(model, X_train, y_train, scoring='accuracy', cv=cv, n_jobs=-1) # report the model performance print('Mean Accuracy: %.3f (%.3f)' % (np.mean(n_scores), np.std(n_scores))) C:\Users\phill\anaconda3\lib\site-packages\sklearn\model_selection_split.py:676: UserWarning: The least ted class in y has only 5 members, which is less than n_splits=10. warnings.warn(C:\Users\phill\anaconda3\lib\site-packages\sklearn\model_selection_split.py:676: UserWarning: The least ted class in y has only 5 members, which is less than n_splits=10. warnings.warn(C:\Users\phill\anaconda3\lib\site-packages\sklearn\model_selection_split.py:676: UserWarning: The least ted class in y has only 5 members, which is less than n_splits=10. warnings.warn(Mean Accuracy: 0.014 (0.006) # fit the model on the whole dataset model.fit(X_train, y_train) # get predictions yhat = model.predict(X_test) # summarize the predicted class print('Predicted Class: %d' % yhat[0])	ð	<pre># training a Naive Bayes classifier from sklearn.naive_bayes import GaussianNB gnb = GaussianNB().fit(X_train, y_train) gnb_predictions = gnb.predict(X_test) # accuracy on X_test gnb_accuracy = gnb.score(X_test, y_test) print(gnb_accuracy) # creating a confusion matrix</pre>
<pre># define the multinomial logistic regression model model = LogisticRegression(multi_class='multinomial', solver='lbfgs') # define the model evaluation procedure cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1) # evaluate the model and collect the scores n_scores = cross_val_score(model, X_train, y_train, scoring='accuracy', cv=cv, n_jobs=-1) # report the model performance print('Mean Accuracy: %.3f (%.3f)' % (np.mean(n_scores), np.std(n_scores))) C:\Users\phill\anaconda3\lib\site-packages\sklearn\model_selection_split.py:676: UserWarning: The least ted class in y has only 5 members, which is less than n_splits=10. warnings.warn(C:\Users\phill\anaconda3\lib\site-packages\sklearn\model_selection_split.py:676: UserWarning: The least ted class in y has only 5 members, which is less than n_splits=10. warnings.warn(C:\Users\phill\anaconda3\lib\site-packages\sklearn\model_selection_split.py:676: UserWarning: The least ted class in y has only 5 members, which is less than n_splits=10. warnings.warn(Mean Accuracy: 0.014 (0.006) # fit the model on the whole dataset model.fit(X_train, y_train) # foot predictions yhat = model.predict(X_test) # summarize the predicted class print('Predicted Class: %d' % yhat[0])</pre>	, , , , , , , , , , , , , , , , , , ,	<pre>#getting precision and recall gnb_metrics = precision_recall_fscore_support(y_test, gnb_predictions, average='macro', labels=np.unique print(gnb_metrics) 0.00782122905027933 (0.005754585488041371, 0.03063725490196078, 0.007884959539371305, None) Multinomial Logistic Regression from sklearn.linear_model import LogisticRegression</pre>
<pre>warnings.warn(C:\Users\phill\anaconda3\lib\site-packages\sklearn\model_selection_split.py:676: UserWarning: The least ted class in y has only 5 members, which is less than n_splits=10. warnings.warn(C:\Users\phill\anaconda3\lib\site-packages\sklearn\model_selection_split.py:676: UserWarning: The least ted class in y has only 5 members, which is less than n_splits=10. warnings.warn(Mean Accuracy: 0.014 (0.006) # fit the model on the whole dataset model.fit(X_train, y_train) # get predictions yhat = model.predict(X_test) # summarize the predicted class print('Predicted Class: %d' % yhat[0])</pre>		<pre># define the multinomial logistic regression model model = LogisticRegression(multi_class='multinomial', solver='lbfgs') # define the model evaluation procedure cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1) # evaluate the model and collect the scores n_scores = cross_val_score(model, X_train, y_train, scoring='accuracy', cv=cv, n_jobs=-1) # report the model performance print('Mean Accuracy: %.3f (%.3f)' % (np.mean(n_scores), np.std(n_scores)))</pre> C:\Users\phill\anaconda3\lib\site-packages\sklearn\model_selection_split.py:676: UserWarning: The least
<pre>yhat = model.predict(X_test) # summarize the predicted class print('Predicted Class: %d' % yhat[0])</pre>		<pre>warnings.warn(C:\Users\phill\anaconda3\lib\site-packages\sklearn\model_selection_split.py:676: UserWarning: The least ted class in y has only 5 members, which is less than n_splits=10. warnings.warn(C:\Users\phill\anaconda3\lib\site-packages\sklearn\model_selection_split.py:676: UserWarning: The least ted class in y has only 5 members, which is less than n_splits=10. warnings.warn(Mean Accuracy: 0.014 (0.006) # fit the model on the whole dataset model.fit(X_train, y_train)</pre>
		<pre>print('Predicted Class: %d' % yhat[0])</pre>