

Background

As a Data Analyst I am always seeking opportunities to increase my skill set and better myself. I have completed a variety of courses surrounding the Data Science field including machine learning, neural networks, NLP, and web scraping. As a way to show case this I took it upon myself to create a WebScraper using the BeautifulSoup library in Python.

Overview

In this case study I set out to search Indeed for Zillow house postings and gather a few key features from the listings. I thought it most important to gather the price, address, number of beds, baths, and square feet if listed to compare listings in the current market. I did so by generating a WebScraper in Jupyter Notebook using the BeautifulSoup Library in Python. A spreadsheet was then generated and imported into Excel, where using the get data function, it was cleaned and organized for better interpretation. The cleaned data was then loaded into Tableau for visualization and analysis. I did this to quickly gain insight into markets for buying a future home or to potentially invest in.

Data Extracting

In Jupyter notebook I created a webscraper that will iterate through a selected number of pages on Zillow and pull all the data from the listings and format the price, address, beds, baths, and square feet separated into individual columns. I did so by using the BeautifulSoup library for html parsing as well as the Pandas library for creating and structing data frames. After populating my data frame with all the metrics I was looking for I exported this data frame into an excel file so I could do additional cleaning and visualization. On the next page the code and the associated commenting is shown.

Python Script

```
1 #Zillow Webscraper by Drew Fingerhut
2 #4/1/2022
3
4 #Importing the necessary Libraries
5 from bs4 import BeautifulSoup
6 import numpy as np
7 import pandas as pd
8 import requests
9
10 #Establishing the lists to be populated outside of iteration
11 price = []
12 bed = []
13 address = []
14
15 #Creating the loop for a specified range representing the number of pages to scrape
16 for i in range(30):
17
18     #Passing in a request header as to get the dynamic page rather than the static page that can not be scraped
19     req_headers = {
20         'accept': 'text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,image/apng,*/*;q=0.8',
21         'accept-encoding': 'gzip, deflate, br',
22         'accept-language': 'en-US,en;q=0.8',
23         'upgrade-insecure-requests': '1',
24         'user-agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/61.0.316
25     }
26
27     #Generating a new url for each page that is to be requested
28     with requests.Session() as s:
29         page_number = i+1
30         page = str(page_number)+'_p/'
31         city = 'salt-lake-city/' #*****change this city to what you want*****
32         url = 'https://www.zillow.com/homes/for_sale/'+city+page
33         r = s.get(url, headers=req_headers)
34
35     #Parse the requested html and look for each title card that has the listing info enclosed in it
36     soup = BeautifulSoup(r.content, 'html.parser')
37     house_elements = soup.find_all("div", class_="list-card-info")
38
39     #Search through each title card for price, beds, and address and populating to the preestablished lists only keeping the
40     for house_char in house_elements:
41         price_element = house_char.find("div", class_="list-card-price")
42         bed_element = house_char.find("ul", class_="list-card-details" )
43         address_element = house_char.find("address", class_="list-card-addr")
44         if price_element != None:
45             price.append(price_element.text)
46         elif price_element == None:
47             price.append('Blank')
48         if bed_element != None:
49             bed.append(bed_element.text)
50         elif bed_element == None:
51             bed.append('Blank')
52         if address_element != None:
53             address.append(address_element.text)
54         elif address_element == None:
55             address.append('Blank')
56
57     #Populate the data frame with the Listes that were just appended to
58     df = pd.DataFrame ({
59         "Price": price,
60         "Bed": bed,
61         "Address": address
62     })
63     #Create a save point of the data frame before altering the data to ensure there is not loss
64     df_new = df.copy()
```

Python Script (continued)

```
65
66 #Replacing unnecessary text in Bed column with spaces so that it can be seperated into individual columns
67 df_new['Beds'] = df_new['Bed'].str.replace(' bds', ' ')
68 df_new['Beds'] = df_new['Beds'].str.replace(' ba', ' ')
69 df_new['Beds'] = df_new['Beds'].str.replace(' sqft', ' ')
70 df_new['Beds'] = df_new['Beds'].str.replace(',', '')
71 df_new['Beds'] = df_new['Beds'].str.replace('--', '')
72 df_new['Beds'] = df_new['Beds'].str.replace('/', ' ')
73 df_new['Beds'] = df_new['Beds'].str.split('-')
74 df_new['Beds'] = df_new['Beds'].fillna(0)
75 df_new.drop(df_new.tail(1).index,inplace=True)
76
77 #Cleaning and dropping the wording such as Multi-family for sale or Townhouse for sale from data frame
78 split_df = pd.DataFrame(df_new['Beds'].tolist(),columns=['Info','Lose','Ex'])
79 split_df = split_df.drop(['Lose','Ex'],axis=1)
80
81 #Splititng Bed coloumn post cleaning into 3 seperate columns for bed, bath, and square feet
82 split_df[['Bedrooms','Bathrooms','Sq_Feet']] = split_df['Info'].str.split(' ', n=2,expand=True)
83 split_df = split_df.drop(['Info'],axis=1)
84
85 #Concatting our split data frame to the original and dropping now unnecessary columns
86 full_df = pd.concat([df_new, split_df],axis=1)
87 full_df = full_df.drop(['Bed','Beds'],axis=1)
88
89 #Displaying our full data frame after iterating through all the pages we defined it to
90 full_df
```

	Price	Address	Bedrooms	Bathrooms	Sq_feet
0	\$799,900	1768 S 1400 E, Salt Lake City, UT 84105	4	2	2222
1	\$569,000	32 W 1700 S #A13, Salt Lake City, UT 84115	3	3	1725
2	\$485,000	886 N Catherine St, Salt Lake City, UT 84116	3	2	1670
3	\$355,000	1162 S Emery St, Salt Lake City, UT 84104	2	1	772
4	\$525,000	1627 S Denver St, Salt Lake City, UT 84115	2	1	1683
...
244	\$815,000	166 S 1200 E, Salt Lake City, UT 84102	4	2	2096
245	\$679,900	1710 S West Temple St W UNIT 11, Salt Lake Cit...	3	4	2096
246	\$695,000	1202 E Princeton Ave, Salt Lake City, UT 84105	3	1	1730
247	\$525,000	1627 S Denver St, Salt Lake City, UT 84115	2	1	1683
248	\$589,900	753 E Loveland Ave, Salt Lake City, UT 84106	3	2	1702

249 rows × 5 columns

```
In [2]: 1 final_df = full_df.drop_duplicates()
```

```
In [4]: 1 final_df.to_csv('Salt_Lake_City_Zillow_Data.csv')
```

Transforming the Data

After running the script above for Seattle, Salt Lake City, Bozeman, and Saint Louis and saving each of them as separate csv files I used the get data function in Excel to gather them all into one sheet and do some transformation to make visualization easier. I went through and removed commas and dollar signs. I seperated the address into individual address, city, and state columns. I also replaced certain verbage that was not initially caught in the Python script from the bed column. After this I then loaded the data back into Excel.

Table.TransformColumnTypes(*Replaced Value6,({"Sq_Feet", Int64.Type}, {"Bathrooms", Int64.Type}, {"Bedrooms", Int64.Type}, {"Price", Currency.Type}))

	Price	Street Address	City	State	Zip Code	Bedrooms	Bathrooms	Sq_Feet
1	450,000.00	3324 W Babcock St	Bozeman	MT	59718	3	2	1
2	1,850,000.00	662 Coffee Creek Rd	Bozeman	MT	59713	4	5	5
3	599,000.00	2394 Renee Way	Bozeman	MT	59718	3	3	3
4	470,000.00	3650 Potosi St	Bozeman	MT	59718	3	3	3
5	789,000.00	1209 S Montana Ave	Bozeman	MT	59715	5	2	2
6	740,000.00	400 Buckhorn Trl	Bozeman	MT	59718	4	3	3
7	534,900.00	3758 Potosi St	Bozeman	MT	59718	3	3	3
8	699,000.00	325 Meriwether Ave	Bozeman	MT	59718	3	3	3
9	680,000.00	138 Albrecht Trl #B	Bozeman	MT	59718	3	2	2
10	Error	Blank		null	null	null	Error	null
11	4,950,000.00	4060 Johnson Rd	Bozeman	MT	59718	5	4	4
12	749,900.00	18 Rustler Trl	Bozeman	MT	59718	3	3	3
13	2,499,000.00	56 T Bone Way	Bozeman	MT	59718	4	4	4
14	899,000.00	921 Josephine Dr	Bozeman	MT	59715	2	2	2
15	1,895,000.00	2135 Barter Dr	Bozeman	MT	59715	4	4	4
16	3,800,000.00	168 Hyalite View Dr	Bozeman	MT	59718	5	4	4
17	400,000.00	848 Saint Andrews Dr	Bozeman	MT	59715	1	null	4
18	1,395,000.00	8121 Fowler Ln	Bozeman	MT	59718	3	4	4
19	16,495,000.00	13830 Cottonwood Rd	Bozeman	MT	59718	4	4	4
20	749,900.00	Willow Blvd	Bozeman	MT	59718	3	null	4
21	3,459,900.00	770 Mountain Moose Rd	Bozeman	MT	59715	1	null	4
22	375,000.00	LOT 1 Riparian Way	Bozeman	MT	59718	1	null	4
23	3,000,000.00	25 Leafmaster Trl	Bozeman	MT	59718	5	6	6
24	800,000.00	Nhn Hyalite Canyon Rd	Bozeman	MT	59715	1	null	4
25	2,499,000.00	333 Black Bull Trl	Bozeman	MT	59718	4	5	5
26	490,000.00	579 N Saint Andrews Dr	Bozeman	MT	59715	1	null	4
27	2,195,000.00	1 Jackson Creek Rd PARCEL 1A	Bozeman	MT	59715	40	40	40
28	2,588,000.00	13259 Astell Gateway Rd	Gallatin Gateway	MT	59730	13	13	13

AutoSave: ON | Zillow_data | Search (Alt+Q) | d.fingerhut38@gmail.com

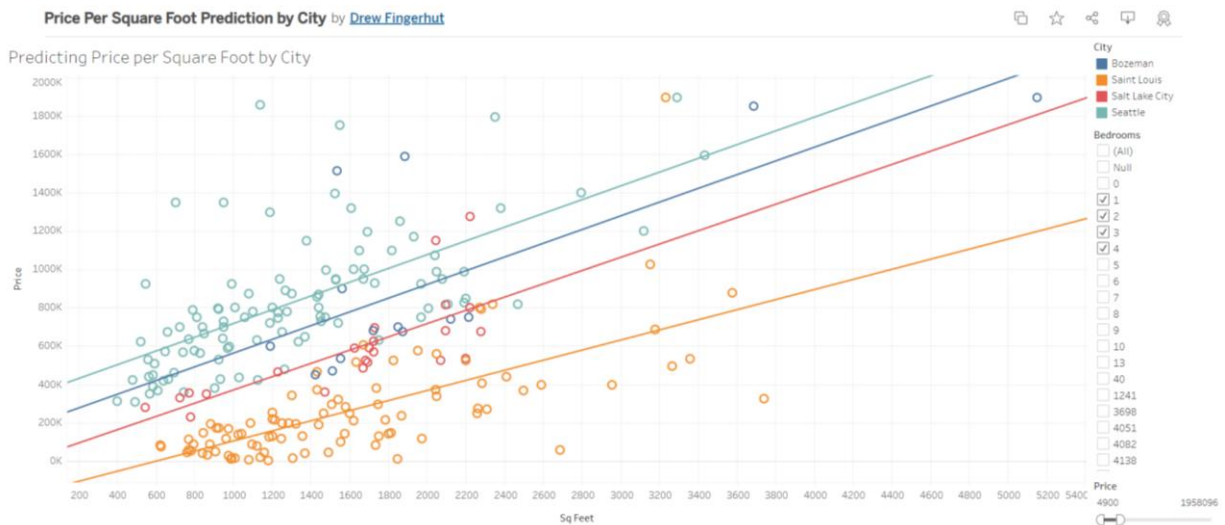
	Price	Street Address	City	State	Zip Code	Bedrooms	Bathrooms	Sq_Feet
30	1513000	315 N Tracy Ave #401	Bozeman	MT	59715	2	2	1536
31	337000	LOT 7 Riparian Way	Bozeman	MT	59718	1		
32	330000	LOT 6 Riparian Way	Bozeman	MT	59718	1		
33	1588000	315 N Tracy Ave #202	Bozeman	MT	59715	2	3	1886
34	445700	LOT 8 Audubon Way	Bozeman	MT	59715	8088		
35	850000	25 Green Mountain Way LOT 25	Bozeman	MT	59715	9		
36	2000000	Dovetail Ln	Bozeman	MT	59718	8		
37	2750000	N Wallace Ave #9	Bozeman	MT	59715	3	2	2240
38	340000	LOT 8 Riparian Way	Bozeman	MT	59718	1		
39	675000	1487 Bora Way	Bozeman	MT	59718	3	3	1875
40	815000	4488 Glenwood Dr	Bozeman	MT	59718	5	3	2382
41	325000	Leiden Ln	Bozeman	MT	59718	4051		
42	130000	5254 Maple Ave	Saint Louis	MO	63113	4	3	1750
43	125000	4642 Leona St	Saint Louis	MO	63116	3	1	1186
44	99900	4728 Hammett Pl	Saint Louis	MO	63113	3	2	1554
45	169900	615 Clara Ave APT 10	Saint Louis	MO	63112	2	2	975
46	140000	5408 Maple Ave	Saint Louis	MO	63112	5	2	2294
47	115000	4327 W Belle Pl	Saint Louis	MO	63108	8	10	3376
48	250000	3918 Kingsland Ct	Saint Louis	MO	63116	4	2	2160
49	175000	6919 Wise Ave	Saint Louis	MO	63139	2	1	924
50	339900	6211 Pernod Ave	Saint Louis	MO	63139	3	2	2047
51		Blank						
52	9900	4933 Plover Ave	Saint Louis	MO	63120	3	2	1080
53	144900	3722 Dunnica Ave	Saint Louis	MO	63116	3	3	1800
54	144900	1113 Washington Ave UNIT 509	Saint Louis	MO	63101	1		1043
55	75000	210 N 17th St UNIT 812	Saint Louis	MO	63103	1		624
56	535000	4467 McPherson Ave	Saint Louis	MO	63108	4	5	3358
57	235000	5114 Daggett Ave	Saint Louis	MO	63110	2	2	1870
58	79900	5037 Durant Ave	Saint Louis	MO	63115	2	1	1121

Data Visualization

With this new compounded Excel sheet I uploaded it into Tableau for visualization. I thought it best to break it down into a few separeate visuals. I first broke it into Average Price by City, Average Price by Zip Code, and Number of house available by Zip Code. I then combined these into a dashboard that can be filtered out by price. I did this for a large overarching visualization to show you a few key metrics about the different areas for further investigation depending on investing constraint or market availability.



I then created a visualization that can be filtered out by price and number of bedrooms. I then added coloring to the cities for easier interpretation and a trendline based on the price compared to the square footage. This is a key metric that when comparing the markets will show you how much house you are getting for the price.



Reviewing Success

At the end of this exercise, I view it as a massive success. I was able to write a universal script that can pull data for any city and number of pages from Zillow and save it to a csv file. I was then able to transform this data for further analysis and visualize it for key insights. I was able to identify that Seattle was the most expensive price by square foot and Saint Louis was the cheapest. I was able to create a visualization that based on the square feet and city could infer a price it would be listed for. I was also able to see how many houses were available per zip code and what the average price in that zip code was. This will give me the necessary key metrics I would need when starting to house hunt and give me the flexibility I need to sort by number of beds and price.

Future Use

Moving forward I will be adding more cities to this and narrow in on a particular region to invest in. With the prebuilt scripting and excel queries it will be a quick and painless process to analyze any area for future investing. Thanks for reading through this and feel free to reach out if there are any questions!