

Clustering of Houston City Neighborhood

For Coursera capstone project

Jason Chen, June 2020

Introduction

According to Wikipedia ^[1] Houston is the most populous city in TX with an estimated population in 2019 of 2.3 million ^[2]. **Figure 1** shows that from 1950 to 2010 the population of Houston city has a steady annual growth rate of 20,000 people/year. People move to Houston for its plentiful jobs, education, and favorable housing price ^[3]. There are approximately 90 so called super neighborhoods (**Figure 2**) in Houston city ^[4]. The neighborhoods have diverse attributes regarding population characteristics, age of population, ethnicity, income, education status, median household income, median housing value, and so on. For anyone who wants to move to Houston or anyone who wants to open a small business (for example, a day-care center, spider-smart training center, or a restaurant with a typical food type) it would be important to know the clustering and distribution of these diverse neighborhoods, so that an optimal location can be chosen for the small business. This study is about the clustering of Houston city neighborhoods using the various demographic and other data obtained online and the Foursquare location data.

Data

On its main website ^[4], the super neighborhood's demographics, land use map, and resource assessment are summarized in PDF files accessible through hyperlinks. **Figure 3** shows an example of the first page of the PDF file that contains the feature data regarding each neighborhood, including population characteristics (total population, and persons per sq. mile), age of population (percentage of under 5 years, 5-17 years, 18-64 years, 65 and over), ethnicity (non-Hispanic Whites, non-Hispanic Blacks, Hispanics, Non-Hispanic Asians, and non-Hispanic others), income categories (under 25,000, 25,001 – 50,000, 50,001 – 100,000, over 100,001), median household income, educational status (no diploma, high school diploma, some college, bachelor's or higher), and housing and households (total housing units, occupied, vacant, total households, family households, and median housing value). In total 76 PDF files are downloaded excluding those neighborhoods that do not have such online data.

The other major data type is the Foursquare search results of the various vendors close by each neighborhood. The longitude and latitude of the neighborhood for the Foursquare search is obtained through Google search and summarized in an Excel file. Three neighborhoods are excluded in this process as they are in the airport or lake area.

Data preparation

Data preparation is performed with the following steps:

- Latitude and longitude data for each neighborhood is obtained through Google search
- 73 PDF files are parsed using PDFMINER to convert PDF into text
- The feature data of each neighborhood is obtained through string search of the above text file
- The various data from above is combined into a data frame.

The resulting data frame is shown in **Figure 4**. This data frame along with the Foursquare search results will be used to cluster the neighborhood using different approaches.

Methodology

The methodology applied in the data analysis include poly-nomial regression and the k-mean clustering. Poly-nomial regression is used to explore the relationship between median household income and median housing value. There are 73 data points in total. The data points are split into training and test data set (15% of data is used for test). The test set is used to determine the degree of the poly-nomial fitting. Maximum R^2 is reached when the degree = 2 (**Figure 5**). Therefore, 2nd order poly-nomial regression is used to fit the data between median household income and median housing value. R^2 of the regression is 0.723 (**Figure 5**), which suggests that there is a good correlation between median household income and median housing value. Not surprising as people with higher income tend to buy more expensive houses.

k-mean clustering is used to segment the neighborhoods into different groups. The segmentation is done using 5 types of different data:

- 1.) median household income,
- 2.) median housing value,
- 3.) age of population,
- 4.) Ethnicity,
- 5.) Vendor data from Foursquare searching.

The focus of the segmentation or clustering of the neighborhoods is to provide answer to the optimum location for a potential small business, for example, new day-care center for Hispanics or new senior center for Asian group. 7 clusters are chosen as it provides good solutions. 5 clusters are used for the vendor data from Foursquare searching for simplicity purpose.

Results

The clustering results are shown in the following Figures.

Figure 6: clustering by median household income.

Figure 7: clustering by median housing value.

Figure 8: clustering by age of population.

Figure 9: clustering by age of population: potential location for day-care center and senior center.

Figure 10: clustering by Ethnicity.

Figure 11: clustering by Ethnicity: location for Hispanic community and non-Hispanic Asian community.

Figure 12: clustering by vendors from Foursquare searching results: Cluster 1 ● : Asian restaurant, Cluster 2 ● : Mexican restaurant, Cluster 3 ● : Coffee/Hotel, Cluster 4 ● : Mexican / ice-cream, Cluster 5 ● : Discount store.

Figure 13: Potential location for a new day – care center for the Hispanic community.

Figure 14: Potential location for a new senior care center for the non-Hispanic Asian community.

Conclusions

The following conclusion can be drawn from this study:

- 1.) Clustering method can be used to identify potential business locations. For example, this study identifies potential new day – care center for Hispanic community and new senior care center for non-Hispanic Asian community.
- 2.) Median housing value has a linear relationship with median household income, the clustering results with respect to median housing value and median household income are identical (shown in **Figure 15**).

Discussion

The foursquare searching results tend to be concentrated with vendors such as restaurants, coffee shops, and stores.

Figures

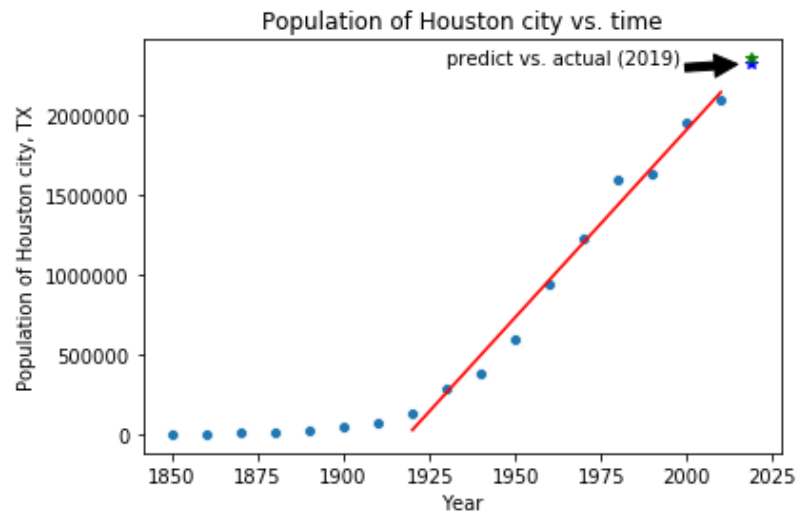


Figure 1. population of Houston from 1850 to 2010. A linear regression was performed using the historical data. The prediction of the population in 2019 matches well with the actual data.

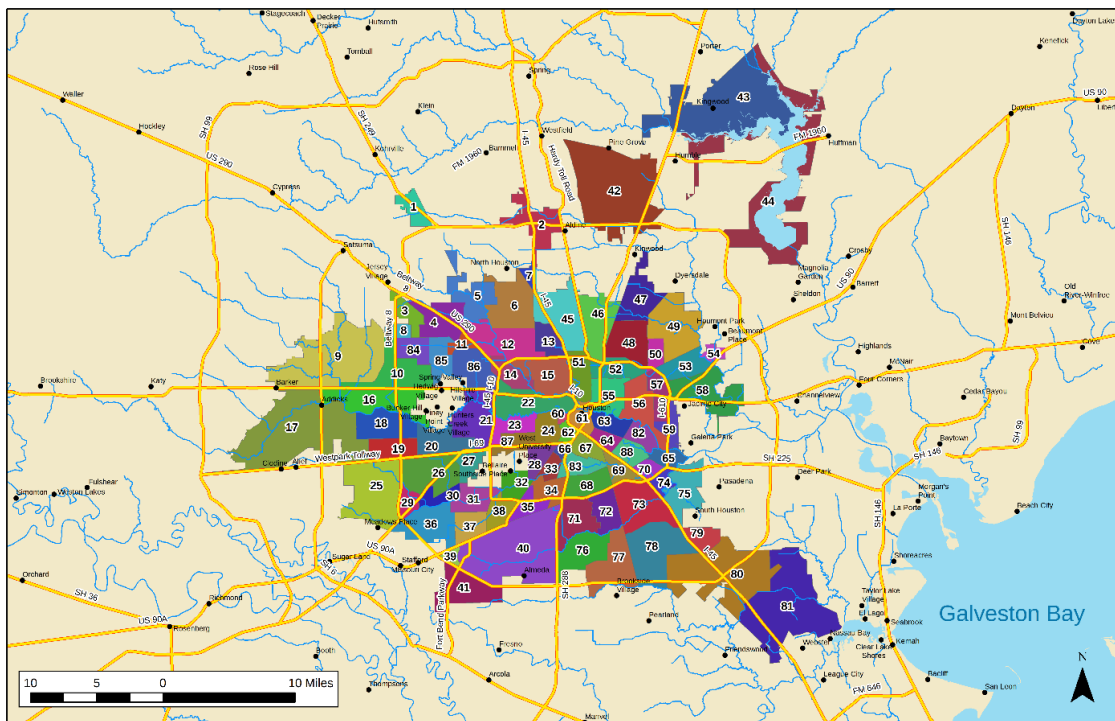


Figure 2. Neighborhoods of Houston city, TX ^[1].

SUPER

Neighborhood

RESOURCE ASSESSMENT



	Super Neighborhood		Houston	
Pop. characteristics	2000	2015	2000	2015
Total population	2,741	8,509	1,953,631	2,217,706
Persons per sq. mile	908	2,818	3,166	3,314
Age of Population				
Under 5 years	8%	9%	8%	8%
5- 17 years	13%	18%	19%	17%
18- 64 years	60%	66%	64%	65%
65 and over	19%	7%	9%	10%
Ethnicity				
Non Hispanic Whites	56%	27%	31%	26%
Non Hispanic Blacks	19%	26%	25%	22%
Hispanics	15%	27%	37%	44%
Non Hispanic Asians	8%	17%	6%	7%
Non Hispanic Others	2%	3%	1%	1%
Income				
Under \$25,000	37%	23%	33%	27%
\$25,001 to \$50,000	32%	20%	31%	25%
\$50,001 to \$100,000	26%	37%	24%	26%
Over \$100,001	5%	19%	12%	22%
Median Household Income	\$32,366	\$58,713	\$36,616	\$46,187
Educational Status				
No Diploma	13%	9%	30%	23%
High School Diploma	28%	19%	20%	23%
Some College	28%	36%	23%	24%
Bachelor's or Higher	31%	36%	27%	31%
Housing and Households				
Total housing units	1,660	3,660	782,378	909,336
Occupied	85%	91%	92%	88%
Vacant	15%	9%	8%	12%
Total households	1,416	3,320	717,945	799,714
Family households	546	1,107	457,548	491,778
Median Housing Value	\$64,976	\$153,557	\$79,300	\$131,700

Description

Willowbrook is in northwest Harris County. The area generally surrounds Willowbrook Mall and is primarily commercial, with about 63.9% undeveloped land. The area includes apartment complexes, office buildings, a major Houston Lighting and Power electric generating station and several retail shopping centers in addition to the regional mall. The area was annexed by the City in 1993.

Highlights

- ✓ Houston City Council District A
- ✓ Cypress-Fairbanks Independent School District
- ✓ 1 Police beat (includes bordering beats)
- ✓ 1,933 acres (3.02 sq. miles)



PLANNING &
DEVELOPMENT
DEPARTMENT

Figure 3. Example of PDF file of super neighborhood data sheet ^[4].

```
df_data = pd.DataFrame.from_dict(data, orient='index',
                                columns=['Income_Median', 'Housing_Median', 'Population_Total',
                                         'Psmile_Person', 'Age%_Under5', 'Age%_5_17', 'Age%_18_64',
                                         'Age%_65up', 'Non_Hispanic%_W', 'Non_Hispanic%_B',
                                         'Hispanic%', 'Non_Hispanic%_A', 'Nin_Hispanic%_O',
                                         'Income%_U25k', 'Income%_U50k', 'Income%_U100k',
                                         'Income%_100kup', 'No_Diploma%', 'High_School%',
                                         'College%', 'Bachelor_Or_Higher%'])

df_data['Index'] = df_n['Index'].to_numpy()
dft = pd.merge(df_n, df_data, on='Index')
dft.head()
```

	HoustonNeighborhoods	Index	Latitude	Longitude	Income_Median	Housing_Median	Population_Total	Psmile_Person	Age%_Under5	Age%_5_17	Age%_18_64	Age%_65up	Non_Hispanic%_W	Non_Hispanic%_B	Hispanic%	Non_Hispanic%_A	Nin_Hispanic%_O	Income%_U25k	Income%_U50k	Income%_U100k	Income%_100kup	No_Diploma%	High_School%	College%	Bachelor_Or_Higher%
0	9_Addicks_Park_Ten	9	29.813300	-95.645500	80584.0	168155.0	19683.0	840.0	5.0	12.0	23.0	60.0	78.0	18.0	12.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0
1	23_AftonOaks_RiverOaks	23	29.749994	-95.433234	95682.0	499169.0	14518.0	4021.0	5.0	12.0	23.0	60.0	78.0	18.0	12.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0
2	25_Alief	25	29.682700	-95.593200	41833.0	90655.0	106657.0	7544.0	7.0	12.0	23.0	60.0	78.0	18.0	12.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0
3	34_Astrodome_Area	34	29.685045	-95.409813	46284.0	102268.0	18223.0	4846.0	5.0	12.0	23.0	60.0	78.0	18.0	12.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0
4	30_Braeburn	30	29.682779	-95.534980	42958.0	116547.0	18843.0	4711.0	8.0	12.0	23.0	60.0	78.0	18.0	12.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0

5 rows × 25 columns

Figure 4. The head information for the combined data frame.

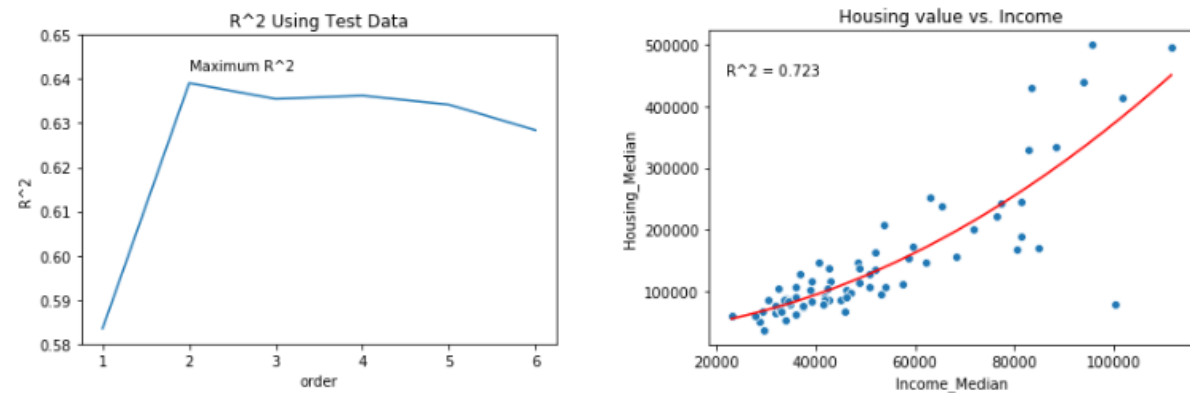


Figure 5. 2nd order poly-nomial regression on the relationship between median housing value and median household income ($R^2 = 0.723$).

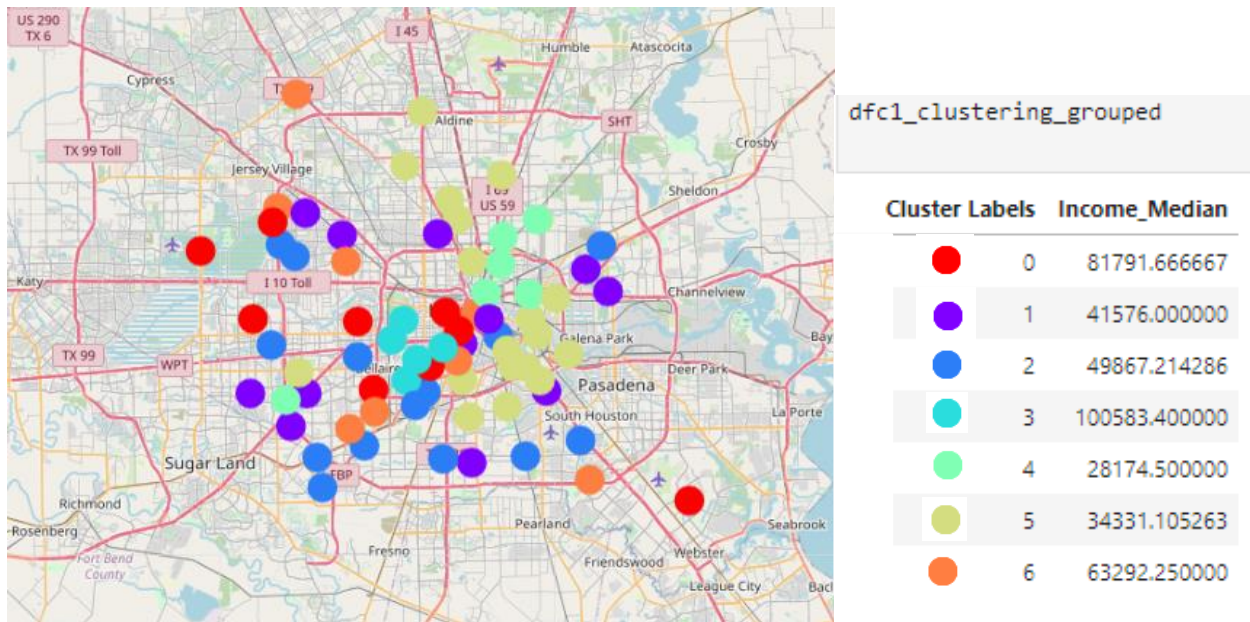


Figure 6: clustering by median household income.

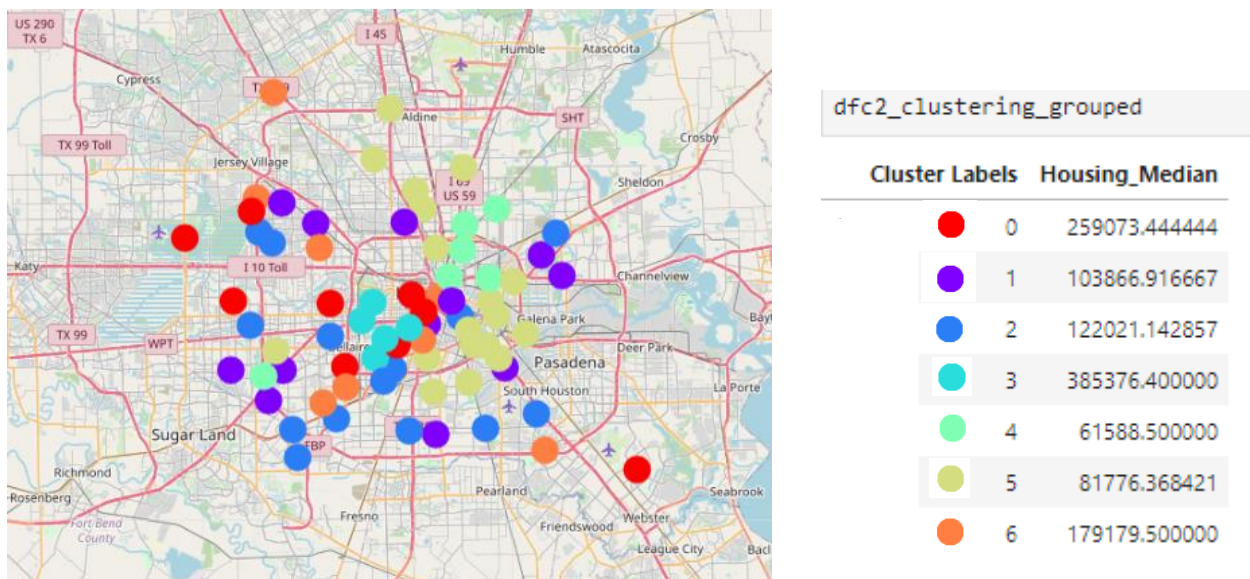


Figure 7: clustering by median housing value.

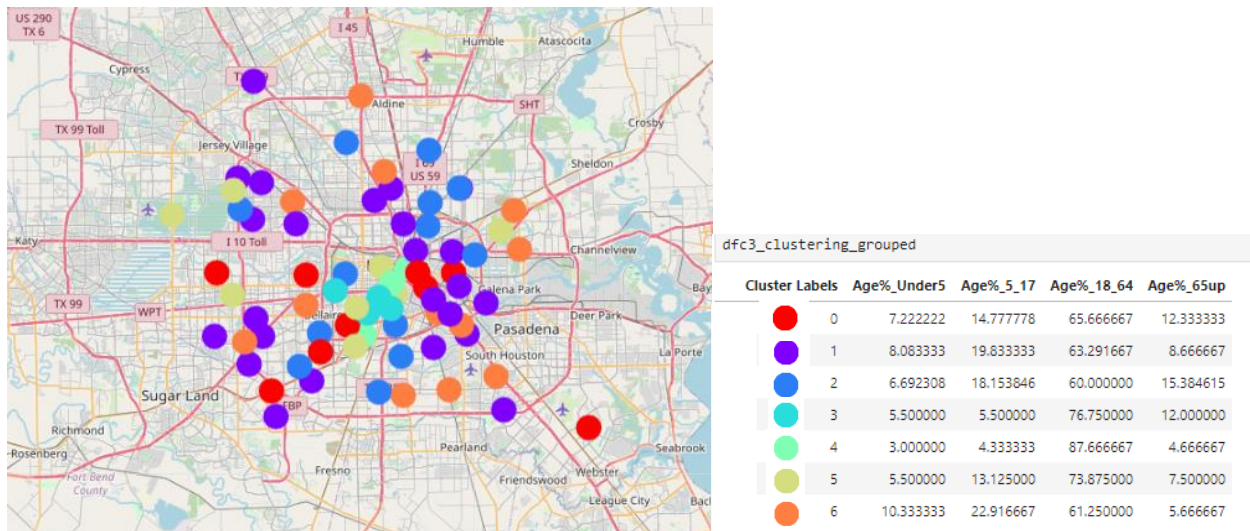


Figure 8: clustering by age of population.

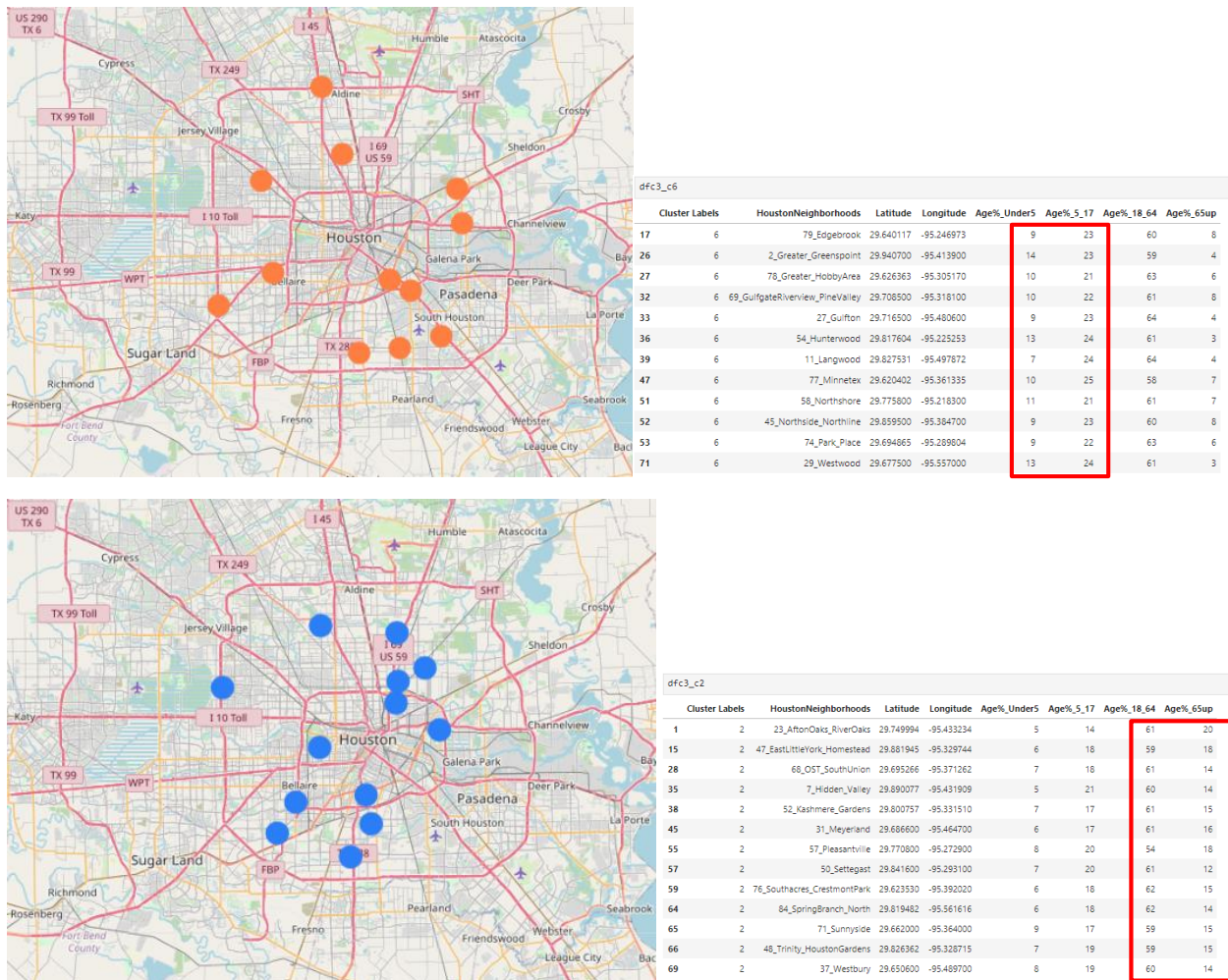


Figure 9: clustering by age of population: potential location for day-care center (top figure) and senior center (bottom figure).

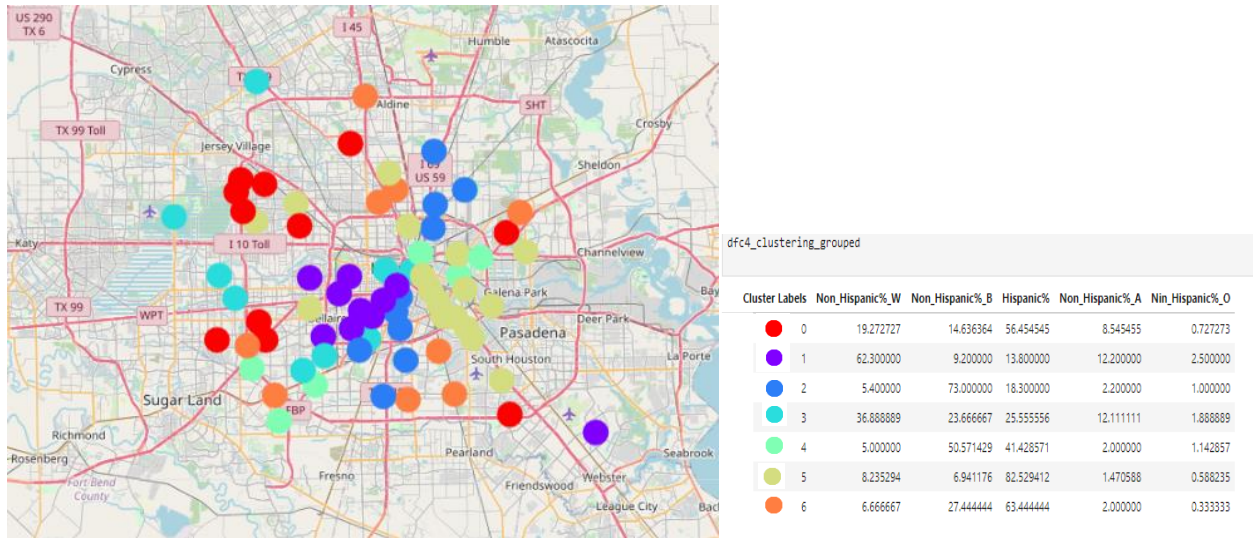


Figure 10: clustering by Ethnicity.

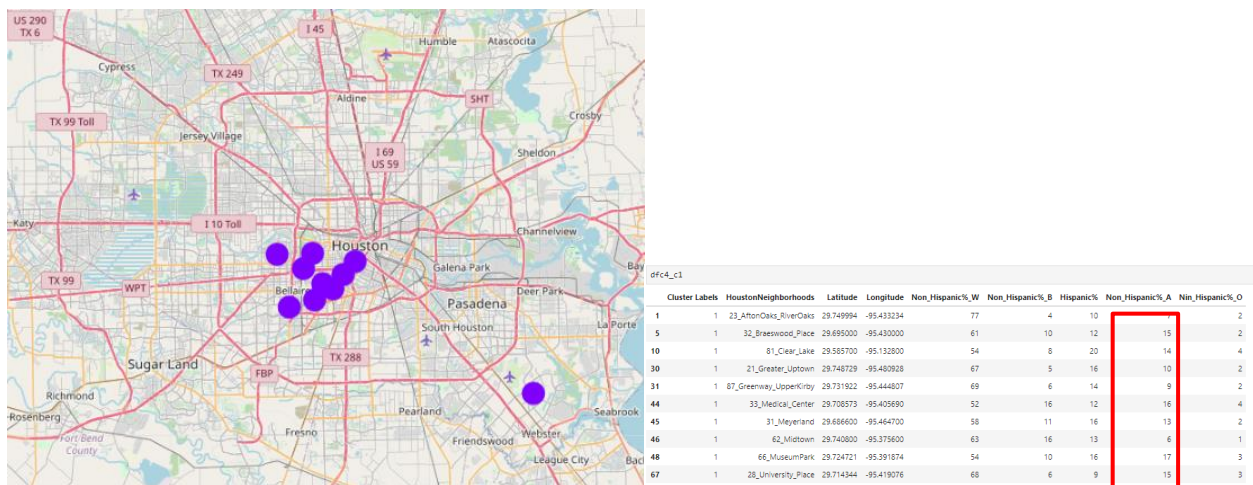
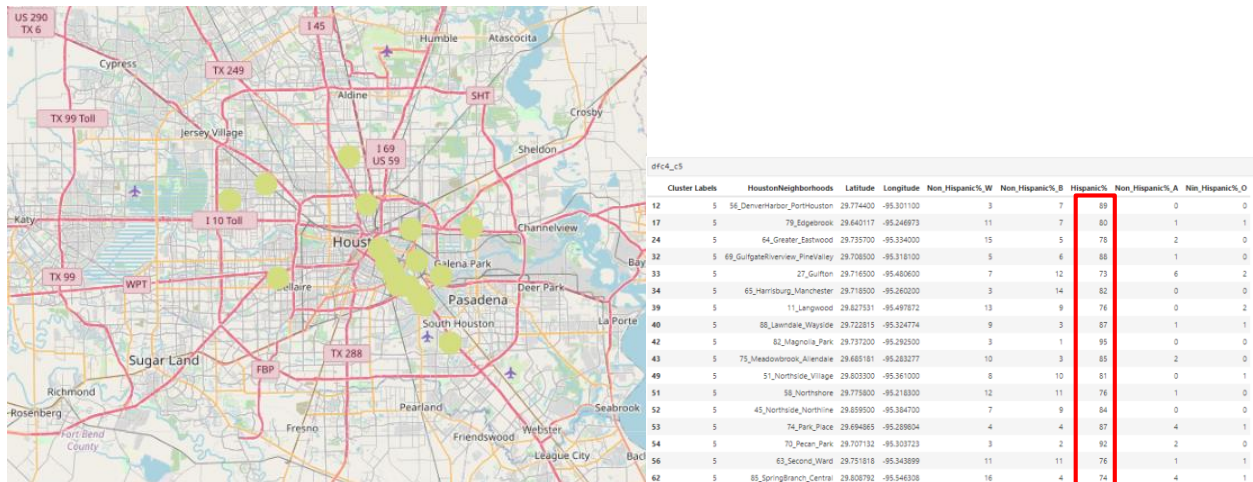


Figure 11: clustering by Ethnicity: location for Hispanic community (top figure) and non-Hispanic Asian community (bottom figure).

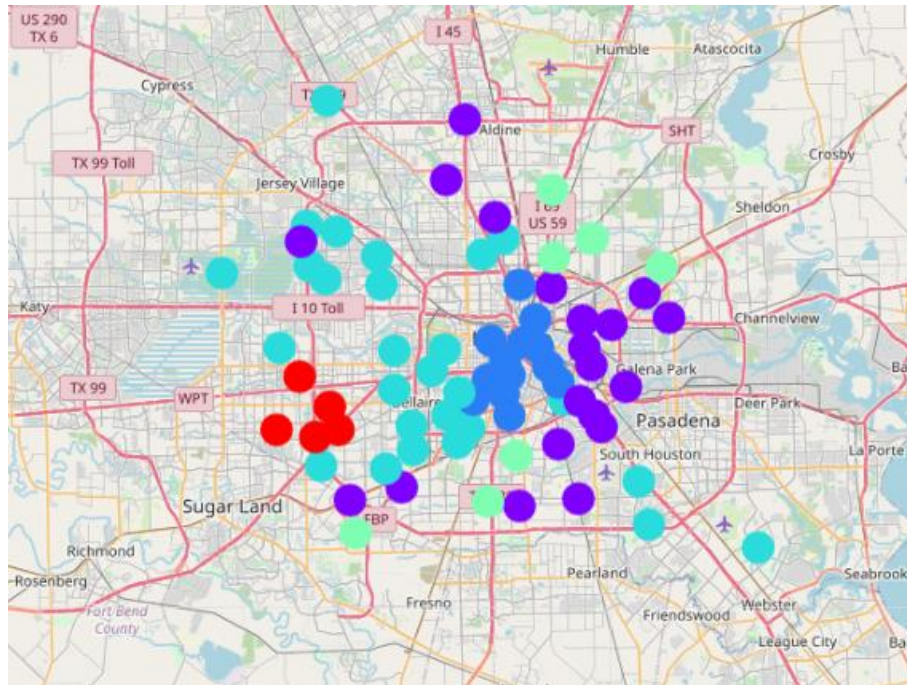
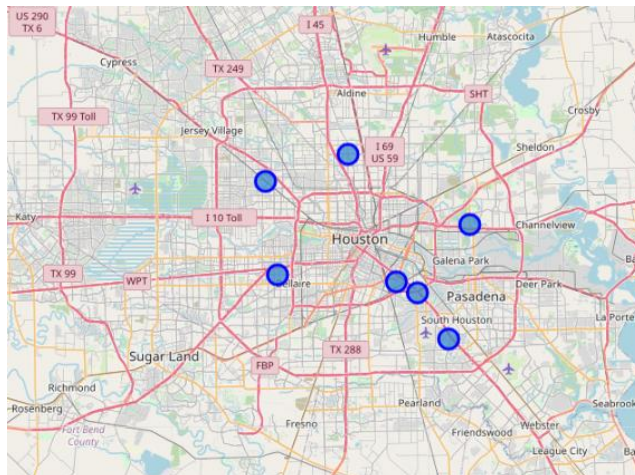


Figure 12: clustering by vendors from Foursquare searching results: Cluster 1 ● : Asian restaurant, Cluster 2 ● : Mexican restaurant, Cluster 3 ● : Coffee/Hotel, Cluster 4 ● : Mexican / ice-cream, Cluster 5 ● : Discount store.



	HoustonNeighborhoods	Latitude	Longitude	Age%_Under5	Age%_5_17	Hispanic%
0	79_Edgebrook	29.640117	-95.246973	9	23	80
1	69_GulfgateRiverview_PineValley	29.708500	-95.318100	10	22	88
2	27_Gulfton	29.716500	-95.480600	9	23	73
3	11_Langwood	29.827531	-95.497872	7	24	76
4	58_Northshore	29.775800	-95.218300	11	21	76
5	45_Northside_Northline	29.859500	-95.384700	9	23	84
6	74_Park_Place	29.694865	-95.289804	9	22	87

Figure 13: Potential location for a new day – care center for the Hispanic community.

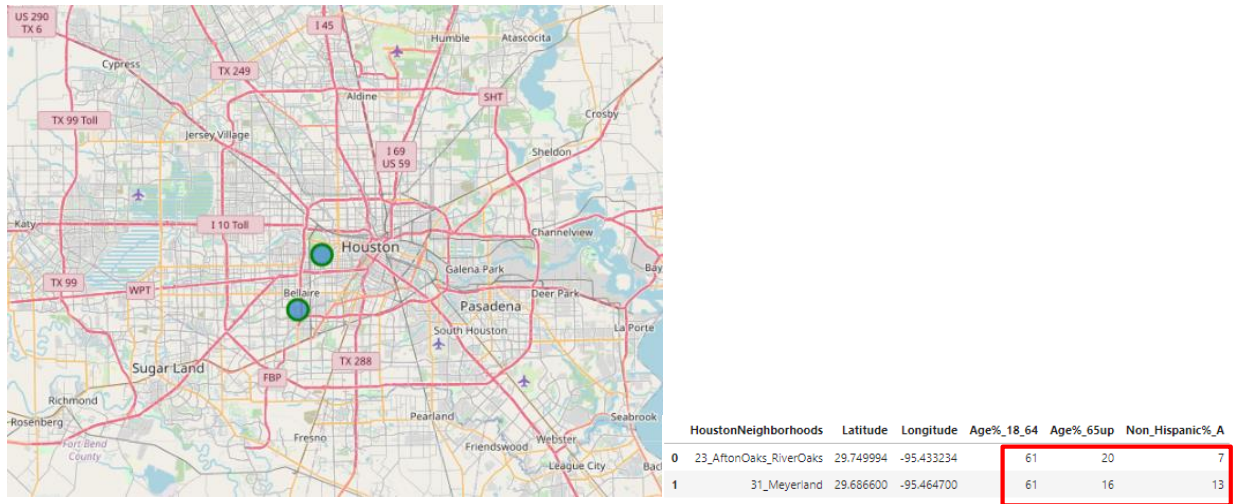


Figure 14: Potential location for a new senior care center for the non-Hispanic Asian community.

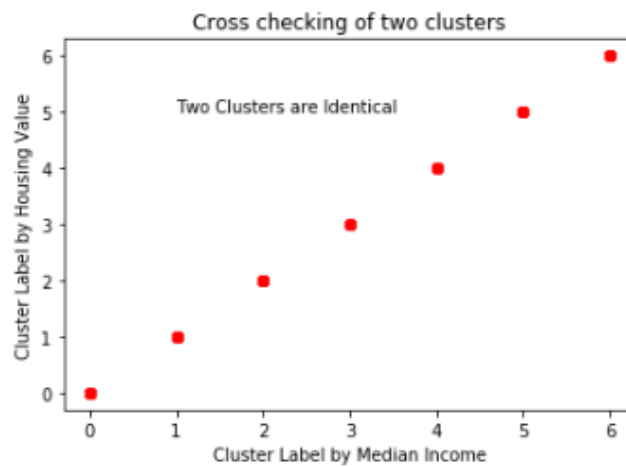


Figure 15: Identical clustering from “median income” and “housing value” when “median income” and “housing value” has a linear relationship ($R^2 = 0.723$).

References

- [1]. <https://en.wikipedia.org/wiki/Houston>, accessed on June 2020.
- [2]. Annual Estimates of the Resident Population for Incorporated Places of 50,000 or More, Ranked by July 1, 2019 Population: April 1, 2010 to July 1, 2019". United States Census Bureau, Population Division. Retrieved May 21, 2020.
- [3]. <https://www.houstonchronicle.com/news/houston-texas/houston/article/Houston-growth-shows-no-sign-of-waning-4540887.php>, accessed on June 14, 2020.
- [4]. Super neighborhood, <https://www.houstontx.gov/superneighborhoods/recognized.html>, accessed on June 14, 2020.