

Correlation Analysis Between Airbnb Prices & Subway Distance

Gangyi Li
12/04/2025

Issues and Objectives

- **Core issue:** Are Airbnb prices really higher when you are closer to the subway?
- **Code:** Multi-library collaboration (pandas/geopandas/sklearn/folium)
- **Process:** data → modeling → visualization
- **Outcome:** Python implementation for spatial data processing and regression analysis.

Data Source

Airbnb Listing Data:

Content: 20,000+ NYC records in CSV/GZ format

Python Implementation: `pd.read_csv()` with native compression support

Subway Station Data:

Content: Latitude, Longitude, Station Name

Python Implementation: Custom `read_subway()` function for automatic column name detection

Neighborhood Income Data:

Content: Sourced from Census API (Population, Median Income)

Python Implementation: Data fetched via `requests` and parsed into a `pd.DataFrame()`

Analysis and Models

Data Cleaning:

- *Remove Missing Values*
- *Standardize Parameters*
- *Spatial Format Conversion*

Spatial Analysis:

Use KMeans to cluster properties with similar geographic locations into 5 categories

Regression Modeling:

- *Why $\text{np.log1p}(y)$? Reduce impact of \$1000+ luxury listings outliers*
- *Outputs: Coefficients (impact strength) + R^2 (model fit)*

Output files

airbnb_processed.csv:

Content: Cleaned data (including distance, clustering, and community attributes)

Python Implementation: Save using `df.to_csv()` with new added features.

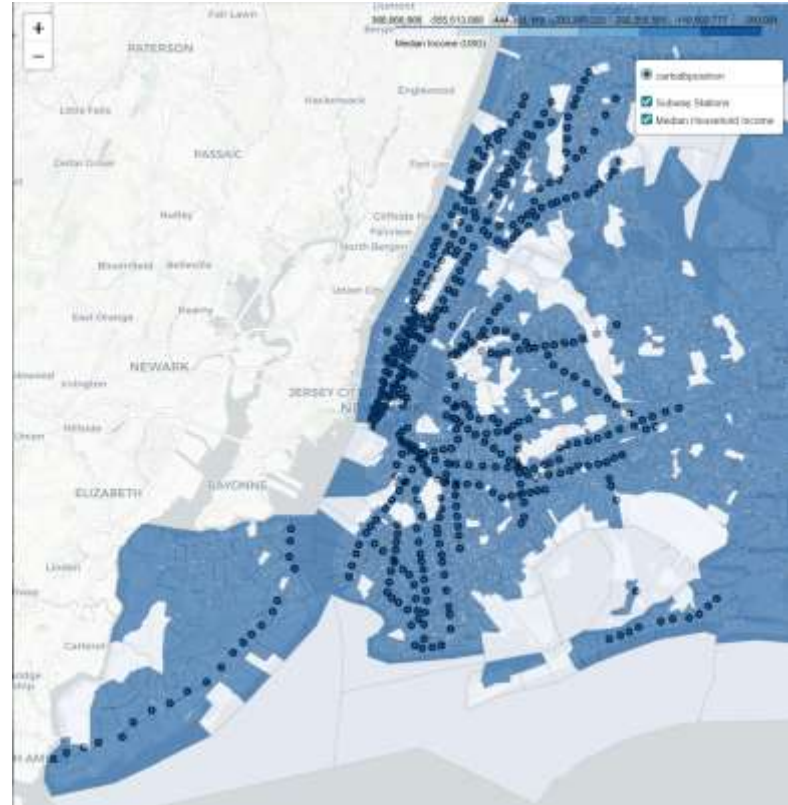
#	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	id	name	host_id	neighbour	latitude	longitude	price	room_type	number_of_index	right	SECD	median_hc	population	dist_to_star	nearest_street	cluster
2	40824219	Room ckt	3.18E+08	Neighbour	40.74698	-73.9176	66	Private room	16	274	3.61E+10	93366	5387	572.4562	46 St-Bless	0
3	40843980	Cozy 2 Br	2.95E+08	Neighbour	40.6823	-73.8455	97	Entire home	93	11	3.61E+10	84630	2152	327.1845	Rodaway	9
4	40824301	Cozy room	1.4890430	Neighbour	40.71316	-73.9421	60	Private room	26	2594	3.6E+10	96875	2939	232.7909	Graham Av	1
5	40825740	House of i	7728764	Neighbour	40.67412	-73.9412	425	Entire home	1	2387	3.6E+10	50417	4489	700.6149	Kingston P	1
6	2595	Skyllt Stud	2645	Neighbour	40.75356	-73.9856	240	Entire home	47	1160	3.61E+10	91250	122	149.6056	42 St-Brya	0
7	6848	Only 2 stc	15091		40.70935	-73.9634	96	Entire home	195	2624	3.6E+10	63608	5921	364.1034	Hewes St	1
8	6672	Uptown S.	1.6104	Neighbour	40.80107	-73.9426	59	Private room	1	4473	3.61E+10	39395	7810	373.0827	116 St	2
9	6090	LES Beaut	16800	Neighbour	40.78778	-73.9478	73	Private room	249	4460	3.61E+10	60778	6011	414.8134	103 St	2
10	7097	Perfect fol	17571	Neighbour	40.69194	-73.9739	218	Private room	423	1694	3.6E+10	-5.7E+08	0	726.8216	Fulton St	1
11	60846106	Comfortal	3.02E+08	Neighbour	40.70652	-73.9029	93	Private room	25	384	3.61E+10	67001	5136	308.9911	Forest Av	1
12	40861640	Essex Hou	36703116		40.76561	-73.9773	420	Entire home	12	4414	3.61E+10	165600	6352	241.4655	57 St	0
13	40860871	Sonder at	2.2E+08	Neighbour	40.70797	-74.0068	448	Entire home	34	1059	3.61E+10	182348	8485	214.2973	Fulton St	0
14	40873866	Home awi	1.43E+08	Neighbour	40.65953	-73.8937	99	Private room	64	3370	3.6E+10	58547	4303	630.1081	New Lots /	1
15	40874667	Full size b	1.43E+08	Neighbour	40.65797	-73.8929	90	Private room	76	4951	3.6E+10	75833	4571	709.3146	New Lots /	1
16	8490	Maison ok	25183		40.68455	-73.9296	170	Entire home	189	2319	3.6E+10	104537	5303	693.875	Kingston-T	1
17	9357	Midtown l	30193	Neighbour	40.76724	-73.9896	176	Entire home	68	1000	3.61E+10	106915	9316	544.9131	99 St-Cok	0
18	10452	Radiant O	95035	Neighbour	40.68294	-73.9568	90	Private room	82	2247	3.6E+10	121250	4425	229.0201	Franklin Av	1
19	12937	1 Stop to	50124	Neighbour	40.74757	-73.9457	232	Private room	456	574	3.61E+10	122125	4131	51.83518	Court Sq-C	0
20	12940	Charming	50148	Neighbour	40.67946	-73.9542	151	Entire home	80	3478	3.6E+10	61250	2499	248.8577	Franklin Av	1
21	14314	Greenport	58246		40.73535	-73.9558	115	Entire home	176	1608	3.6E+10	146458	4778	607.4033	Greenport	0

Output files

nyc_airbnb_map.html:

Content: Interactive Maps (Property Listings, Subway Lines, Income Heat Maps)

Python Implementation: Built with `folium.Map()`, supporting layer control.



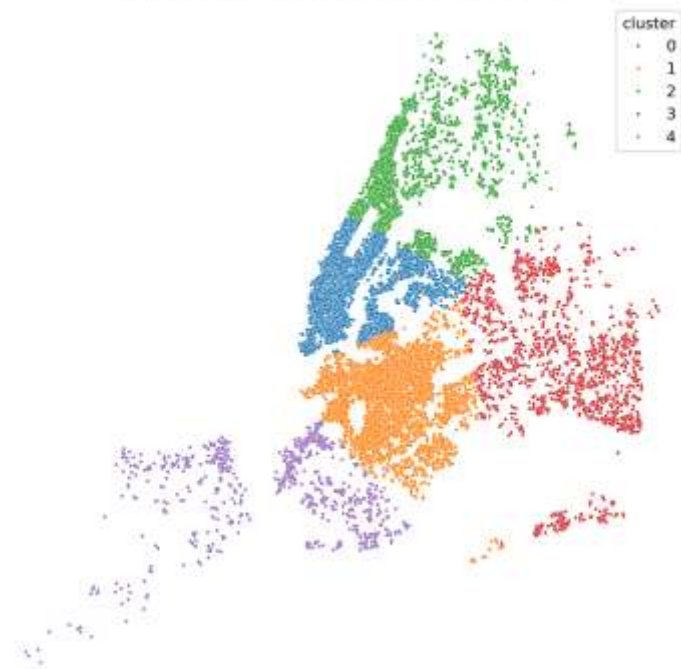
Output files

airbnb_clusters.png:

Content: Spatial cluster results (5 groups)

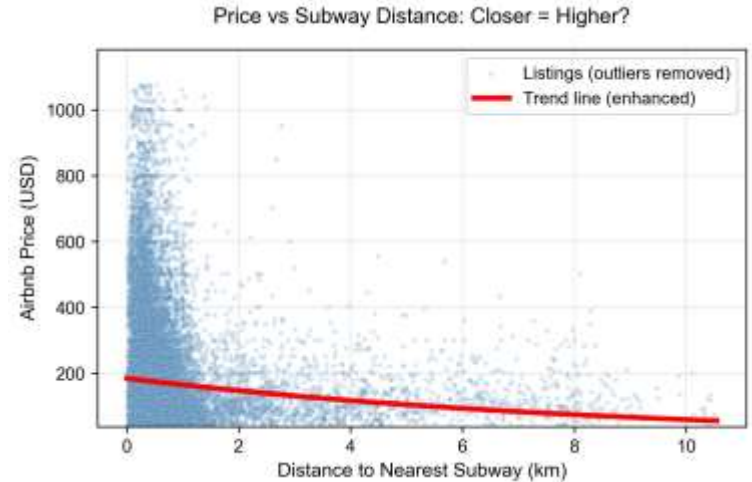
Python Implementation: Clustered via **sklearn.KMeans**, visualized with seaborn.

Airbnb Spatial Clusters (Projected Coordinates)



Result

- The spatial autocorrelation index Moran's I is 0.035 ($p=0.026$) : Weak but significant price clustering.
- In the regression model, the coefficient for subway distance is -0.05: 5% price drop per 1km.
- Model fit $R^2 = 0.45$: 45% price variation explained (distance & income)



- The red trendline illustrates the pattern that 'the farther from the subway, the lower the housing prices,' controlling for confounding factors such as community income.

Challenge

Data Format Inconsistency:

- The data contains numerous mixed string formats (e.g., price information: \$100, \$200, “350”, etc.).
- `def parse_price(x): s = str(x).replace('$', '').replace(',', '').strip()`

Outlier Impact on Regression:

- High-priced properties distort the results.
- `y = np.log1p(df_model['price'])`