# Week 7 Report

# A Clustering Model Based on Japan's housing transaction data

37237256 薛丁銘

## 1. Background

Similar to last week, I decide to do a cluster model by myself. This time I have got a dataset of Japan's housing transaction. Then I clustered the transaction into 10 clusters, showing the 10 most frequent housing transaction areas in Japan.

## 2. Dataset

Two datasets are involved in this task. The first one is the detailed housing transaction data in Japan for the last four month with the machi code of the houses as shown in Figure 2.1. In order to cluster the locations, another dataset with latitudes/longitudes of Japan machi code was prepared as shown in Figure 2.2.



Figure 2.1 The Housing transaction data



Figure 2.2 The latitude/longitude data

## 3. Clustering Model

In the first part I conducted the data preprocess including loading the data and merge the transaction data with the location data. There was something wrong with the length of the machi code so I had to convert the formation. The preprocess code is as shown in Figure 3.1.



```python
# Reading house price data and latitude_longitude data
house_pricing_data = pd.read_csv("school\Introduction to Machine Learning\Week07_exercise\data\house_pricing_data.csv")
latitude_longitude = pd.read_csv("school\Introduction to Machine Learning\Week07_exercise\data\latitude_longitude.csv")
# print(house_pricing_data.columns)
# print(latitude_longitude.columns)


house_machi_code = pd.DataFrame(house_pricing_data['市区町村コード'])
latitude_longitude_reference = latitude_longitude[['コード', '緯度', '経度']]

house_machi_code.dropna(inplace = True)
latitude_longitude_reference.dropna(inplace = True)

datacheck_house_machi_code = pd.DataFrame(house_machi_code).map(lambda x : len(str(x)))
datacheck_latitude_longitude_reference = pd.DataFrame(latitude_longitude_reference['コード']).map(lambda x : len(str(x)))

print(datacheck_house_machi_code.value_counts())
print(datacheck_latitude_longitude_reference.value_counts())

code2lng = {}
code2lat = {}
for i in range(len(latitude_longitude_reference)):
    code = str(latitude_longitude_reference['コード'].iloc[i])
    if (len(code) == 5):
        code = code[0:4]
    else:
        code = code[0:5]
    code2lng.update({int(code) : latitude_longitude_reference["経度"].iloc[i]})
    code2lat.update({int(code) : latitude_longitude_reference["緯度"].iloc[i]})

house_machi_code['latitude'] = house_machi_code['市区町村コード'].map(code2lat)
house_machi_code['longitude'] = house_machi_code['市区町村コード'].map(code2lng)

house_location = []
for i in range(len(house_machi_code)):
    lat = house_machi_code.iloc[i, 1]
    long = house_machi_code.iloc[i, 2]
    house_location.append((lat, long))

print(len(house_location))

# # plot the house positions
# house_machi_code.plot(x = 'latitude', y = 'longitude', kind = 'scatter', label = 'House location', marker = '+', color = 'black')
# plt.title("Scatterplot of house locations")
# plt.xlabel("Latitude")
# plt.ylabel("Longitude")
# plt.show()

location_array = house_machi_code[['latitude', 'longitude']].to_numpy()
```

Figure 3.1 Preprocess of the data

After the preprocessing, I drew I map of the transaction data. The distribution of the transaction data looks like the map of Japan as shown in Figure 3.2. Then I used kmeans to cluster the housing transaction data into 10 clusters. I also drew another map of the clustering result with different colors. The result map is shown in Figure 3.3 and the clustering code is shown in Figure 3.4.
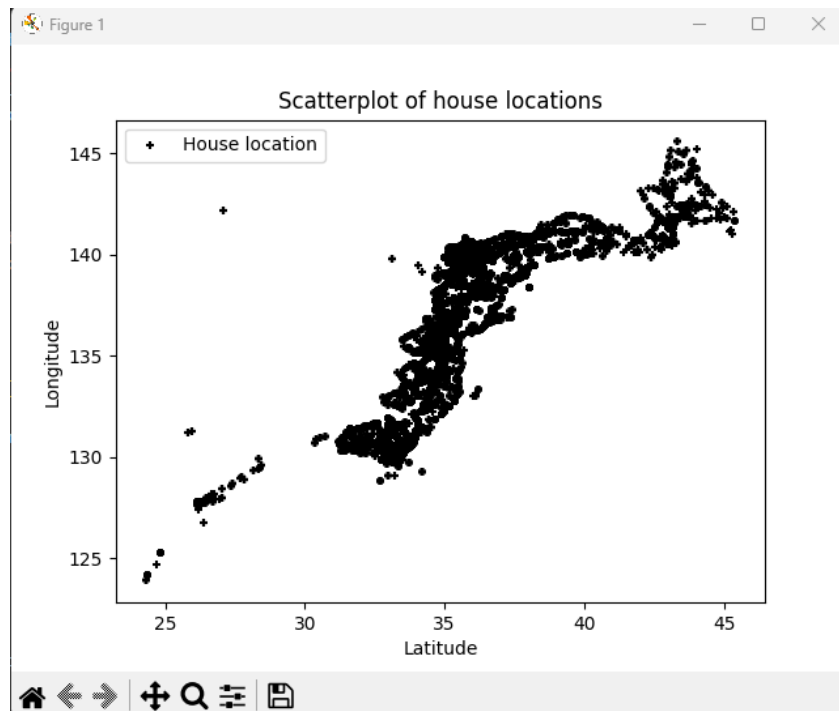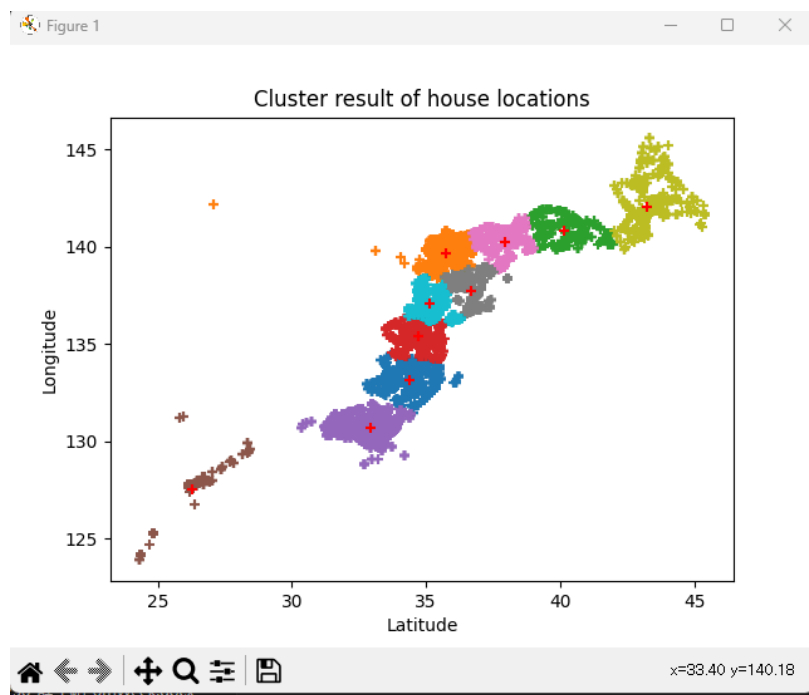
Figure 3.2 Shape of the transaction data



Figure 3.3 Shape of the transaction data

```
kmeans = KMeans(n_clusters=10, random_state = 42)
kmeans.fit(location_array)

centers = kmeans.cluster_centers_
print(centers)

predicted_labels = kmeans.labels_
print(len(predicted_labels))

merged_array = np.concatenate((location_array, predicted_labels.reshape(-1, 1)), axis =1)

print(merged_array)

labels = np.unique(merged_array[:, 2])
split_house_array = [merged_array[merged_array[:,2] == label] for label in labels]
label_count = 0
for array in split_house_array:
    plt.scatter(x = array[:,0], y = array[:,1], label = 'house cluster ' + str(label_count), marker = '+')
    label_count += 1
plt.scatter(x = centers[:,0], y = centers[:,1], label = 'house location', marker = '+', color = 'red')
plt.title("Cluster result of house locations")
plt.xlabel("Latitude")
plt.ylabel("Longitude")
plt.show()
```

Figure 3.4 The clustering code

Finaly I ran the evaluation code as shown in Figure 3.5 and the results is shown in Figure 3.6. There are no other cluster methods for comparing so the clustering evaluation is just for reference. The evaluation took longer time than I expected. I guess it was because the dataset was too large.

```
## Evaluation
silhousette = metrics.silhouette_score(location_array, predicted_labels)
print("Silhousette Score: ", silhousette)

calinski = metrics.calinski_harabasz_score(location_array, predicted_labels)
print("Calinski-Harabasz: ", calinski)

Davies = metrics.davies_bouldin_score(location_array, predicted_labels)
print("Davies-Bouldin: ", Davies)
```

Figure 3.5 The evaluation code

```
Silhousette Score:  0.59788599918214449
Calinski-Harabasz:  362853.02294252854
Davies-Bouldin:  0.6126386250450064
PS D:\code> 
```

Figure 3.6 The evaluation result