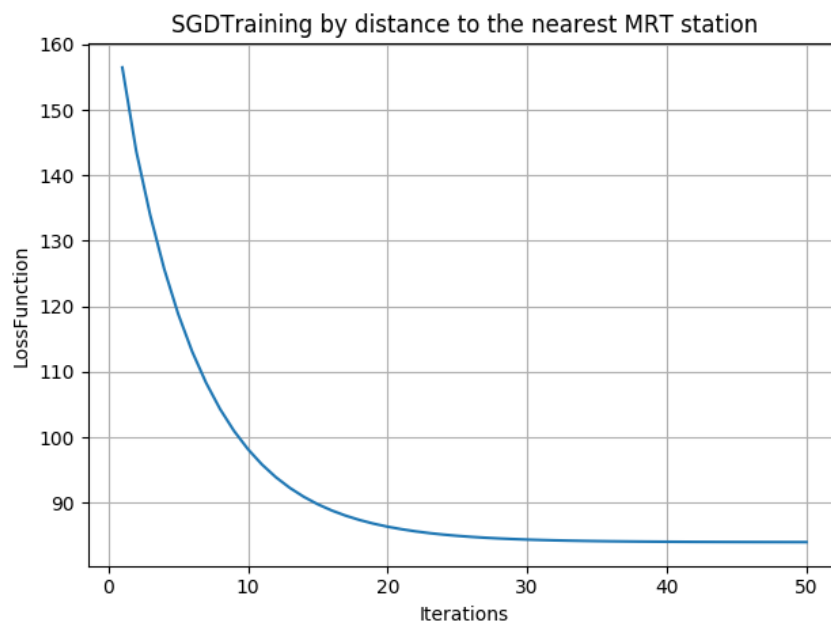
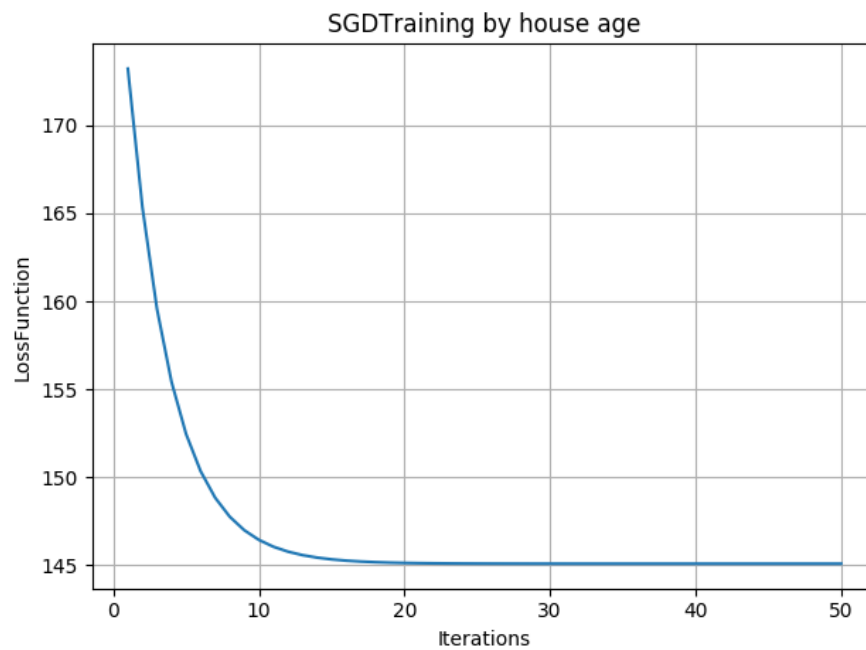
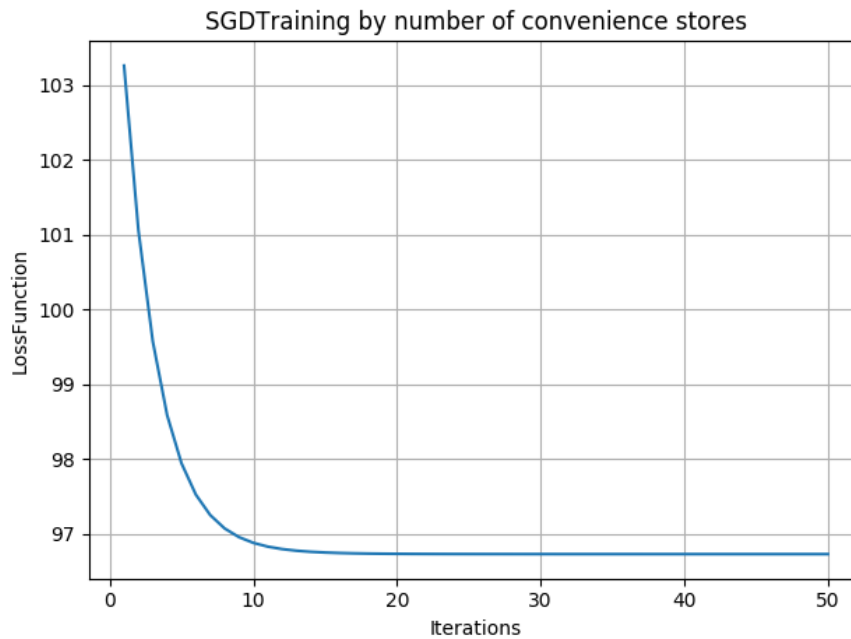


Question 1: The  $\theta$  parameters ( $\theta_0, \theta_1$ ) from step 3 when you are using house age feature. (2 marks)

$$\theta_0 = 42.54078538, \theta_1 = -10.31939902$$

Question 2: A plot, which visualises the change in cost function  $J(\theta)$  at each iteration. (1 mark)





Question 3: RMSE for your training set when you use house age feature. (0.5 mark)

$$RMSE_{\text{training set/house age feature}} = 12.04551030591235$$

Question 4: RMSE for test set, when you use house age feature. (0.5 mark)

$$RMSE_{\text{test set/house age feature}} = 16.58731450340051$$

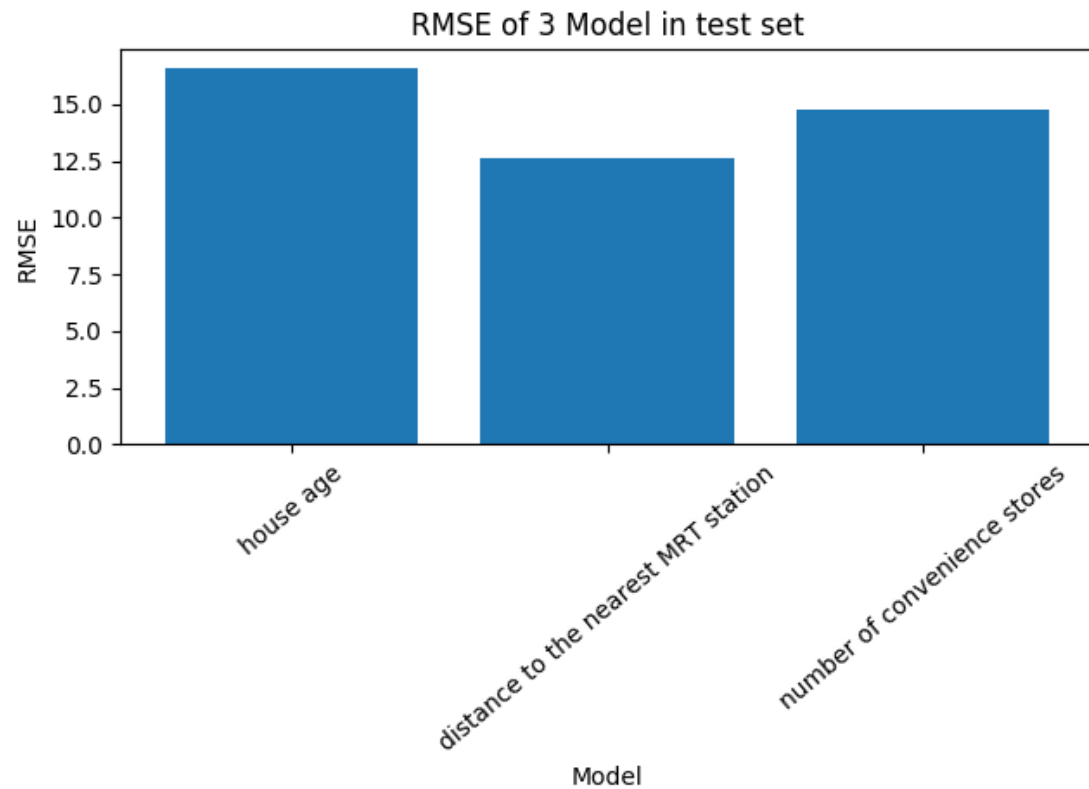
Question 5: RMSE for test set, when you use distance to the station feature. (0.25 mark)

$$RMSE_{\text{test set/distance to the station feature}} = 12.652088009723935$$

Question 6: RMSE for test set, when you use number of stores feature. (0.25 mark)

$$RMSE_{\text{test set/number of stores feature}} = 14.731993508206783$$

Question 7: Compare the performance of your three models and rank them accordingly. (0.5 mark)



The 1<sup>st</sup> is the model training by distance to the station feature.

The 2<sup>nd</sup> is the model training by number of stores feature.

The 3<sup>rd</sup> is the model training by house age feature.

## The code

```
1. import numpy as np
2. import pandas as pd
3. import matplotlib.pyplot as plt
4.
5.
6. def normalize(data):
7.     return (data - data.min()) / (data.max() - data.min())
8.
9.
10. def costFunction(x, y, theta):
11.     return np.dot((np.dot(x, theta) - y).T, (np.dot(x, theta) - y))
12.
13.
14. def targetFunction(x, theta):
15.     return theta[0] + x * theta[1]
16.
17.
18. def SGDTraining(X, Y, parameters, draw = 0):
19.     iteration_times = 0
20.     loss_list = []
21.     theta = np.array(parameters['theta'])
22.     alpha = parameters['alpha']
23.     iterations_max = parameters['iterations_max']
24.     while iteration_times < iterations_max:
25.         for i in range(X.shape[0]):
26.             grad = alpha * (Y[i, :] - X[i, :].dot(theta)) * X[i, :]
27.             theta = theta + np.mat(grad).T
28.             loss = (np.square(Y - X.dot(theta))).sum() / Y.shape[0]
29.             loss_list.append(loss)
30.             iteration_times += 1
31.         if draw == 1:
32.             plt.title("SGDTraining by " + parameters['name'])
33.             plt.grid()
34.             plt.plot(range(1, iterations_max + 1), loss_list)
35.             plt.xlabel("Iterations")
36.             plt.ylabel("LossFunction")
37.             plt.show()
38.             SGDPlot(X, Y, theta, parameters['name'])
39.         return theta
40.
41.
42. def RMSE(RMSE_Data, name, theta):
43.     X = RMSE_Data[['One', name]].values
```

```
44.     Y = RMSE_Data[['house price of unit area']].values
45.     return np.sqrt((np.square(Y - X.dot(theta))).sum() / Y.shape[0])
46.
47.
48. def SGDPlot(X, Y, theta, name):
49.     plt.title("SGDTraining by " + name)
50.     plt.xlabel(name)
51.     plt.ylabel('house price of unit area')
52.     theta = theta.tolist()
53.     for i in range(300):
54.         plt.scatter(X[i][1], Y[i], s = 5, c = 'c')
55.     plt.plot((0, 1), (theta[0][0], theta[0][0] + theta[1][0]))
56.     plt.show()
57.
58.
59. def splitTrain(training_Data, parameters):
60.     X = training_Data[['One', parameters['name']]].values
61.     Y = training_Data[['house price of unit area']].values
62.     theta = SGDTraining(X, Y, parameters, draw = 1)
63.     return theta
64.
65.
66. def preProcessData(data):
67.     data_Norm = normalize(data.iloc[:, 1:4])
68.     data_Norm.insert(0, 'One', 1)
69.     data_Norm.insert(4, 'house price of unit area', data.iloc[:, 4:])
70.     return data_Norm
71.
72. def main():
73.     data = pd.read_csv("house_prices.csv")
74.     data_Norm = preProcessData(data)
75.     training_Data = data_Norm[0:300]
76.     test_Data = data_Norm[300:]
77.
78.
79.     training_Parameters = [{"name": "house age",
80.                             "theta": [[-1], [-0.5]],
81.                             "alpha": 0.01,
82.                             "iterations_max" : 50},
83.                             {"name": "distance to the nearest MRT station",
84.                             "theta": [[-1], [-0.5]],
85.                             "alpha": 0.01,
86.                             "iterations_max" : 50},
87.                             {"name": "number of convenience stores",
```

```
88.         "theta": [[-1], [-0.5]],
89.         "alpha": 0.01,
90.         "iterations_max" : 50}]
91.
92.
93.     theta = {}
94.     RMSE_data = {}
95.     for data in training_Parameters:
96.         theta[data['name']] = splitTrain(training_Data, data)
97.     print(theta)
98.     for name in theta.keys():
99.         RMSE_data[name] = RMSE(test_Data, name, theta[name])
100.    print(RMSE_data)
101.    plt.bar(*zip(*RMSE_data.items()))
102.    plt.xlabel("Model")
103.    plt.ylabel('RMSE')
104.    plt.title("RMSE of 3 Model in test set")
105.    plt.xticks(rotation = 40)
106.    plt.show()
107.
108.
109.
110. if __name__ == '__main__':
111.     main()
```