



实 验 报 告

(2017 / 2018 学年 第 2 学期)

课程名称	机器学习导论
实验名称	Location Assignment
实验时间	2018 年 6 月 20 日
指导教师	王邦

姓名	游浩然	学号	U201515429
----	-----	----	------------

1 问题重述

- 如下图 1 所示，为南一楼二楼局部区域室内布局图。绿色区域为数据采集区域。共包括四个教室和室外的走廊。如下所示建立坐标轴，规定左上角坐标为 $(0,0)$ 。现收集各坐标的不同 BSSID 的信号强度值，拟据此推断其位置坐标。

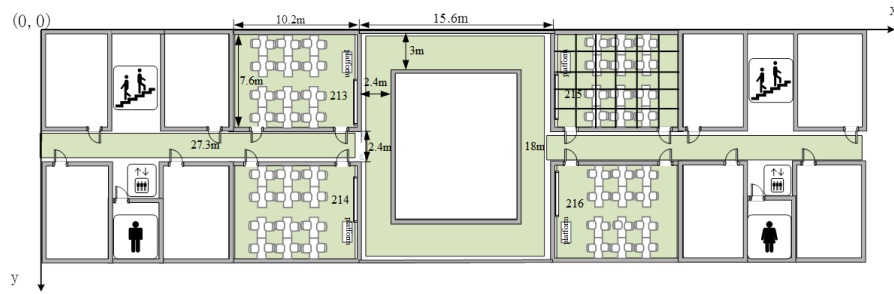


Figure 1: 南一楼局部区域室内布局图

- 训练集一共有 14798 个样本，为了模拟实际场景中收集的样本标注位置存在误差，因此给每个样本的实际坐标添加一定的扰动。具体地， x , y 轴分别添加的扰动为：均值为 0，标准差为 1.2m 的一组随机生成的值。训练集中共接收到 577 个 BSSID，文件中‘负值’表示相应样本接收到相应 BSSID 的信号强度值，‘0 值’表示相应样本未接收到该 BSSID 的值。文件中 x, y 分别表示样本实际采集位置添加扰动后的坐标。
- 测试集一共包括 652 个测试指纹，每个测试指纹是 10 个测试样本的均值，测试集共接收到 458 个 BSSID。

2 问题分析

- 首先观察训练集 S 的特征向量，发现其特征向量并不一致，故取其交集（相较于取并集预测精度有显著提高）。
- 对数据集中特殊点的处理：由于数据集中存在 0 值，可将其变为极弱信号 (-100) 以便于算法预测精度的提升。
- 算法选择方面：考虑该具体问题的性质，虽然是回归问题，但不适合使用 svm、neural network 等强拟合算法。外加噪声会极大地影响其效果，平均可在一定程度上弱化噪声的影响，但综合考虑使用 K 个最近邻样本输出的平均值作为回归预测值 (KNN) 即简单高效，避免复杂计算，又可以得到不错的结果，为上策。

3 算法设计

- main-idea: 用 K 个最近邻样本输出的平均值作为回归预测值。

4 Python Code for Location Assignment

4.1 data

```
1 # -*- coding: utf-8 -*-
2 """
3 @author : Haoran You
4
5 """
6 import csv
7 import random
8 import numpy as np
9
```

```

10 def getFeat():
11     # train dataset
12     trainDT = csv.reader(open('train.csv', 'r'))
13     dataset = []
14     for line in trainDT:
15         dataset.append(line)
16     train_feat_index = dataset[0][1:-2]
17     del (dataset[0])
18     # test dataset
19     testDT = csv.reader(open('test.csv', 'r'))
20     test_feature = []
21     for line in testDT:
22         test_feature.append(line)
23     test_feat_index = test_feature[0][1:]
24     del (test_feature[0])
25     # all feature
26     all_feat = list(set(train_feat_index).intersection(set(test_feat_index)))
27     train_index, test_index = [], []
28     for item in train_feat_index:
29         if item in all_feat:
30             train_index.append(all_feat.index(item))
31         else:
32             train_index.append(-1)
33     for item in test_feat_index:
34         if item in all_feat:
35             test_index.append(all_feat.index(item))
36         else:
37             test_index.append(-1)
38     num_feat = len(all_feat)
39     return num_feat, train_index, test_index, dataset, test_feature
40
41 def traincsv(dataset, train_index, num_feat):
42     feature = []; x = []; y = []
43     for item in dataset:
44         feature.append(item[1:-2])
45         x.append(item[-2])
46         y.append(item[-1])
47     feature = [[float(j) for j in i] for i in feature]
48     train_feature = []
49     for i in range(len(feature)):
50         _feature = list(np.zeros(num_feat))
51         for j in range(len(feature[i])):
52             if feature[i][j] < 0.0:
53                 if train_index[j] == -1:
54                     pass
55                 else:
56                     # _feature[train_index[j]] = 1.0
57                     _feature[train_index[j]] = feature[i][j]
58             else:
59                 if train_index[j] == -1:
60                     pass
61                 else:
62                     _feature[train_index[j]] = -100
63         train_feature.append(_feature)
64     x = [float(i) for i in x]
65     y = [float(i) for i in y]
66     return train_feature, x, y
67
68 def testcsv(feature, test_index, num_feat):
69     for i in range(len(feature)):
70         feature[i] = feature[i][1:]
71     feature = [[float(j) for j in i] for i in feature]
72     test_feature = []
73     for i in range(len(feature)):
74         _feature = list(np.zeros(num_feat))
75         for j in range(len(feature[i])):
76             if feature[i][j] < 0.0:
77                 if test_index[j] == -1:
78                     pass
79                 else:
80                     # _feature[test_index[j]] = 1.0
81                     _feature[test_index[j]] = feature[i][j]
82             else:
83                 if test_index[j] == -1:
84                     pass
85                 else:
86                     _feature[test_index[j]] = -100
87         test_feature.append(_feature)
88     return test_feature
89
90 def divideTrainVal(feature, x, y, ratio):
91     num_dataset = len(x)
92     index_train = random.sample(range(num_dataset), int(num_dataset*ratio))
93     train_feature, train_x, train_y = [], [], []
94     val_feature, val_x, val_y = [], [], []
95     for i in range(num_dataset):
96         if i in index_train:
97             train_feature.append(feature[i])

```

```

98         train_x.append(x[i])
99         train_y.append(y[i])
100     else:
101         val_feature.append(feature[i])
102         val_x.append(x[i])
103         val_y.append(y[i])
104     return train_feature, train_x, train_y, val_feature, val_x, val_y
105
106 def dataset():
107     num_feat, train_index, test_index, dataset, test_feature = getFeat()
108     feature, x, y = traincsv(dataset, train_index, num_feat)
109     train_feature, train_x, train_y, val_feature, val_x, val_y = \
110     divideTrainVal(feature[:, x[:, y[:, ratio=0.9)
111     test_feature = testcsv(test_feature, test_index, num_feat)
112     print('num of trainset : ', len(train_feature))
113     print('num of valset : ', len(val_feature))
114     print('num of testset : ', len(test_feature))
115     print('total features : ', num_feat)
116     return train_feature, train_x, train_y, val_feature, val_x, val_y, test_feature

```

4.2 Main

```

1  # -*- coding: utf-8 -*-
2  """
3  @author : Haoran You
4
5  """
6  import os
7  import csv
8  from data import *
9  import matplotlib.pyplot as plt
10
11 # load data
12 train_feat, train_x, train_y, val_feat, val_x, val_y, test_feat = dataset()
13
14 # train
15 def calDistance(x, y):
16     return np.sqrt(np.square(x) + np.square(y))
17 def method(model_x, model_y):
18     model_x.fit(train_feat, train_x)
19     score = model_x.score(val_feat, val_x)
20     result_val_x = model_x.predict(val_feat)
21     result_test_x = model_x.predict(test_feat)
22     print('score of val_x : ', score)
23     model_y.fit(train_feat, train_y)
24     score = model_y.score(val_feat, val_y)
25     result_val_y = model_y.predict(val_feat)
26     result_test_y = model_y.predict(test_feat)
27     print('score of val_y : ', score)
28     print('average deviation of val_x : ', np.average(abs(val_x - result_val_x)))
29     print('average deviation of val_y : ', np.average(abs(val_y - result_val_y)))
30     print('average deviation of val distance : ', np.average(calDistance(val_x - result_val_x, val_y - result_val_y)))
31     if os.path.exists('result.csv'):
32         os.remove('result.csv')
33     f = open('result.csv', 'a', newline='')
34     csv_write = csv.writer(f, dialect='excel')
35     for i in range(len(result_test_x)):
36         result = []
37         result.append(i)
38         result.append(result_test_x[i])
39         result.append(result_test_y[i])
40         csv_write.writerow(result)
41
42 # plot val result figure
43 plt.figure(1)
44 plt.subplot(131)
45 plt.plot(train_x, train_y, 'ro', label='real')
46 plt.title('trainset distribution')
47 plt.xlabel('x'); plt.ylabel('y')
48 plt.legend(loc='upper right', ncol=1)
49 plt.subplot(132)
50 plt.plot(val_x, val_y, 'ro', label='real')
51 plt.plot(result_val_x, result_val_y, 'bo', label='predict')
52 plt.title('valset distribution')
53 plt.xlabel('x'); plt.ylabel('y')
54 plt.legend(loc='upper right', ncol=1)
55 plt.subplot(133)
56 plt.plot(result_test_x, result_test_y, 'bo', label='predict')
57 plt.title('testset distribution')
58 plt.xlabel('x'); plt.ylabel('y')
59 plt.legend(loc='upper right', ncol=1)
60 plt.show()
61
62 def run(type):

```

```

63 if type == 'decision_tree':
64     from sklearn import tree
65     model = tree.DecisionTreeRegressor()
66 elif type == 'linear':
67     from sklearn import linear_model
68     model = linear_model.LinearRegression()
69 elif type == 'svm':
70     from sklearn import svm
71     model = svm.SVR()
72 elif type == 'KNN':
73     from sklearn import neighbors
74     model = neighbors.KNeighborsRegressor()
75 elif type == 'random_forest':
76     from sklearn import ensemble
77     model = ensemble.RandomForestRegressor(n_estimators=20)
78 elif type == 'adaboost':
79     from sklearn import ensemble
80     model = ensemble.AdaBoostRegressor(n_estimators=50)
81 elif type == 'extra_tree':
82     from sklearn.tree import ExtraTreeRegressor
83     model = ExtraTreeRegressor()
84     method(model, model)
85
86 run('KNN')

```

5 结果讨论

5.1 Location Distribution

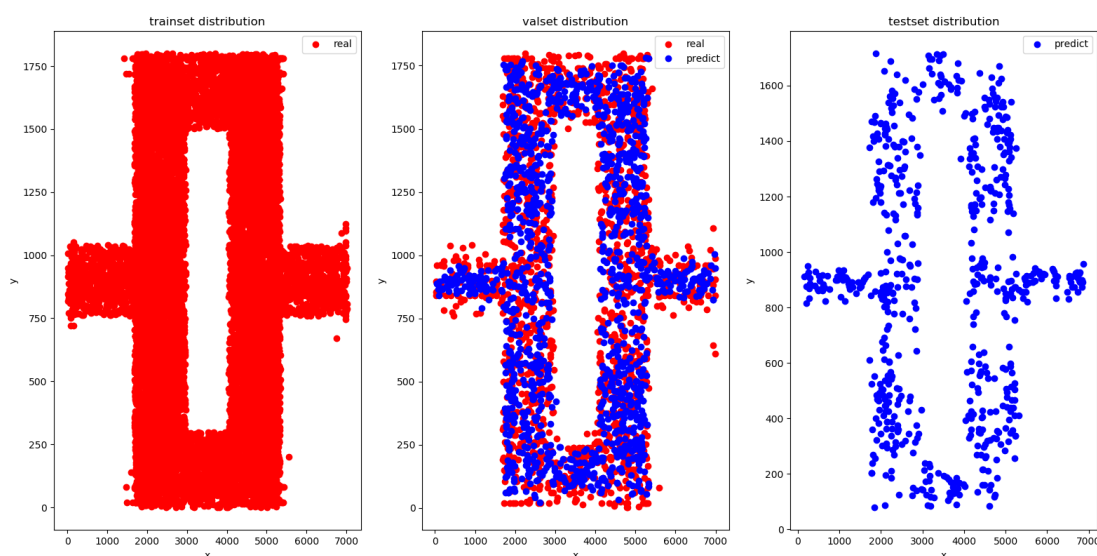


Figure 2: Location Distribution

5.2 CDF

CDF 为累计分布函数，结果的平均定位误差为：

$$err = \frac{\sum_{n=1}^N \sqrt{(x_n^e - x_n)^2 + (y_n^e - y_n)^2}}{N} \quad (5.1)$$

结果如图3所示。

6 课程评价

- 老师准备得很细，看得出画了很大的功夫，但是讲课风格不合我的口味，我比较偏向于原理部分而不是繁琐的计算，老师在讲原理的时候没有讲得很深入或者很形象，可以考虑将课堂的计算换

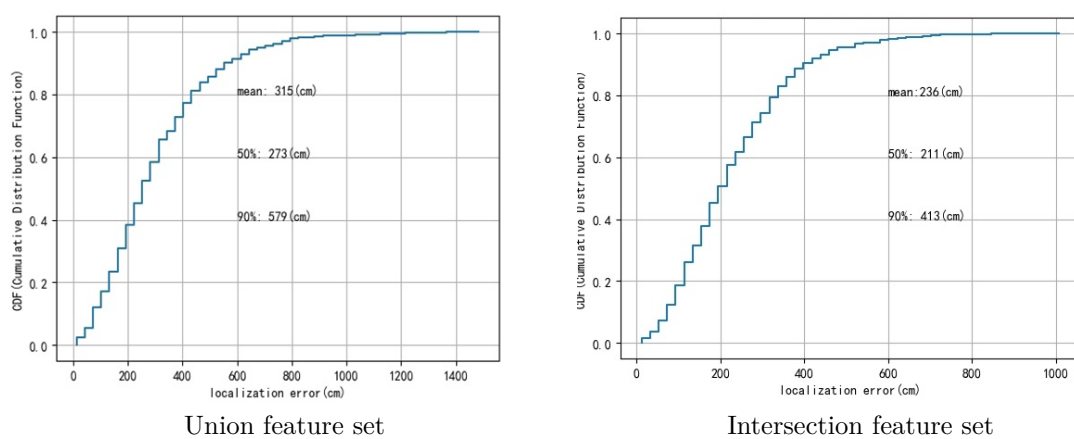


Figure 3: Cumulative Distribution Function

成更简单但也能体现原理的算例。

- 课后作业还可以，编程作业做了还是有收获的。但是 Neural Network 作业貌似没布置。
- 作为学院的第一次尝试来讲可以了。