

实验报告

1 基本思路

代码指路链接:

[林胜峰的Git Hub](#)

先看 `main()` 函数

```
1 def main(k):
2     global accuracy
3     # step 1. train()
4     train()
5     # step 2. validate(i)
6     for i in range(1, k + 1):
7         accuracy.append(validate(i))
8     # step 3. draw()
9     draw(k, accuracy)
10    # step 4. 获取验证效果最好的k, 即 accuracy最大值的下标 + 1
11    index_max = accuracy.index(max(accuracy)) + 1
12    # step 5. 预测, 并计算预测精度
13    do_pre(index_max)
14
15
16 if __name__ == "__main__":
17     main(5) # 超参数从1到5试验一遍
18
```

1.1 读取训练数据集并划分训练集和验证集

`train()`

```
1 # 这个函数的作用有:
2 # 1. 读取 train.csv 中的数据
3 # 2. 划分数数据集(3:1 = 训练数据集:验证数据集)
4 def train():
5     global training_labels, training_vectors, validate_labels,
6     validate_vectors
7     with open('train.csv') as training_csv:
8         training_reader = csv.reader(training_csv, delimiter=' ',
9         quotechar='|')
10        # get rid of header of csv
11        next(training_reader, None)
```

```

10         for row in training_reader:
11             row = [int(i) for i in row[0].split(',')]
12             training_labels.append(row[0]) # 把标签加入训练标签列标中
13             training_vectors.append(row[1:]) # 把数据加入训练数据集中
14             # 划分数据集, 训练:验证 = 3:1, 完全随机划分
15             training_vectors, validate_vectors, training_labels,
validate_labels = \
16                 train_test_split(training_vectors, training_labels,
test_size=0.25, random_state=0)
17             print("Done training!")

```

1.2 超参数设置为k,并用验证集验证 `validate(k)`

```

1  # 验证, k表示knn的超参数k
2  def validate(k):
3      global validate_vectors, validate_labels, validate_temp_labels
4      i = 0
5      validate_temp_labels = [-1] * len(validate_vectors)
6      for row in validate_vectors:
7          min_dist = [sys.maxsize] * k # 长度为k的最小距离列表, 最小距离计算方法
在dist()方法中
8          min_dist_label = [-1] * k # 长度为k的标签列表, 存放可能的标签
9          for count, (vec, num) in enumerate(zip(training_vectors,
training_labels)): # 遍历训练集
10             if count % 2000 == 0: # 其实只有 count 为 0 的时候才生效
11                 print("Validating on No." + str(i))
12                 d = dist(row, vec) # 计算 row(即当前数据) 与 vec(当前测试数据) 的
距离
13                 if d < max(min_dist): # 如果距离比 min_dist 中的最大值小
14                     min_dist_label[min_dist.index(max(min_dist))] = num # 替换
对应最大距离的标签为vec的标签num
15                     min_dist[min_dist.index(max(min_dist))] = d # 替换最大距离为
当前距离d
16                 validate_temp_labels[i] = get_most_label(min_dist, min_dist_label)
# 检查最有可能的标签
17                 i += 1
18             print("Validate done by KNN, where k = " + str(k))
19             # 返回当前k获取的精度
20             return metrics.accuracy_score(validate_labels, validate_temp_labels)

```

1.2.1 计算两个向量距离的函数 `dist(v1, v2)`

```

1 # 计算两个向量距离的函数，其实，可以理解为  $sum = \sum |v1[i] - v2[i]|$ , for i in
  range(len(v1)), 注意 v1与v2长度必相同
2 def dist(vector_one, vector_two):
3     d = 0
4     for i, j in zip(vector_one, vector_two):
5         d += abs(i - j)
6     return d

```

1.3 绘制评价图 `draw()`

```

1 # 绘制评价图，可以是精度也可以是f1，两个参数中必定有一个为True，另一个为 False
2 def draw(k, score, accuracy_score=True, f1_score=False):
3     x = [i for i in range(1, k + 1)]
4     name = ''
5     if accuracy_score and not f1_score:
6         name = 'Accuracy'
7     if f1_score and not accuracy_score:
8         name = 'F1'
9     plt.plot(x, score, label=name + ' Score', linewidth=3, color='r',
10             marker='o')
11     plt.xlabel('number of hyper parameter k')
12     plt.ylabel(name + ' Score')
13     plt.title(name + ' Score of k from 1 to 5')
14     plt.show()

```

1.4 测试 `do_pre(k)`

```

1 # 设置超参数为k，然后进行predict(k)，再根据真实的结果(保存在truth.csv)计算超参数为k
  得到的精度
2 def do_pre(k):
3     global truth_labels, learning_labels
4     predict(k)
5     truth_labels = get_labels('truth.csv', truth_labels)
6     result = metrics.accuracy_score(truth_labels, learning_labels)
7     print(result)

```

1.4.1 测试的核心函数 `predict(k)`

```

1 # 预测，逻辑大致与validate()相同
2 def predict(k):
3     global training_labels, training_vectors, learning_labels,
4     learning_vectors
5     with open('test.csv') as test_csv: # 测试数据在 test.csv

```

```

5         learning_reader = csv.reader(test_csv, delimiter=' ',
quotechar='|')
6         next(learning_reader, None)
7         learning_vectors = list(learning_reader)
8         i = 0
9         learning_labels = [-1] * len(learning_vectors)
10        for row in learning_vectors:
11            row = [int(j) for j in row[0].split(',')]
12            min_dist = [sys.maxsize] * k
13            min_dist_label = [-1] * k
14            closest_num = -1
15            for count, (vec, num) in enumerate(zip(training_vectors,
training_labels)):
16                if count % 10000 == 0:
17                    print("Working on NO." + str(i))
18                d = dist(row, vec)
19                if d < max(min_dist):
20                    min_dist_label[min_dist.index(max(min_dist))] = num
21                    min_dist[min_dist.index(max(min_dist))] = d
22                    closest_num = get_most_label(min_dist, min_dist_label)
23            learning_labels[i] = closest_num
24            i += 1
25        print("KNN Done!")

```

1.4.2 获取真实标签的函数 `get_labels(f, labels)`

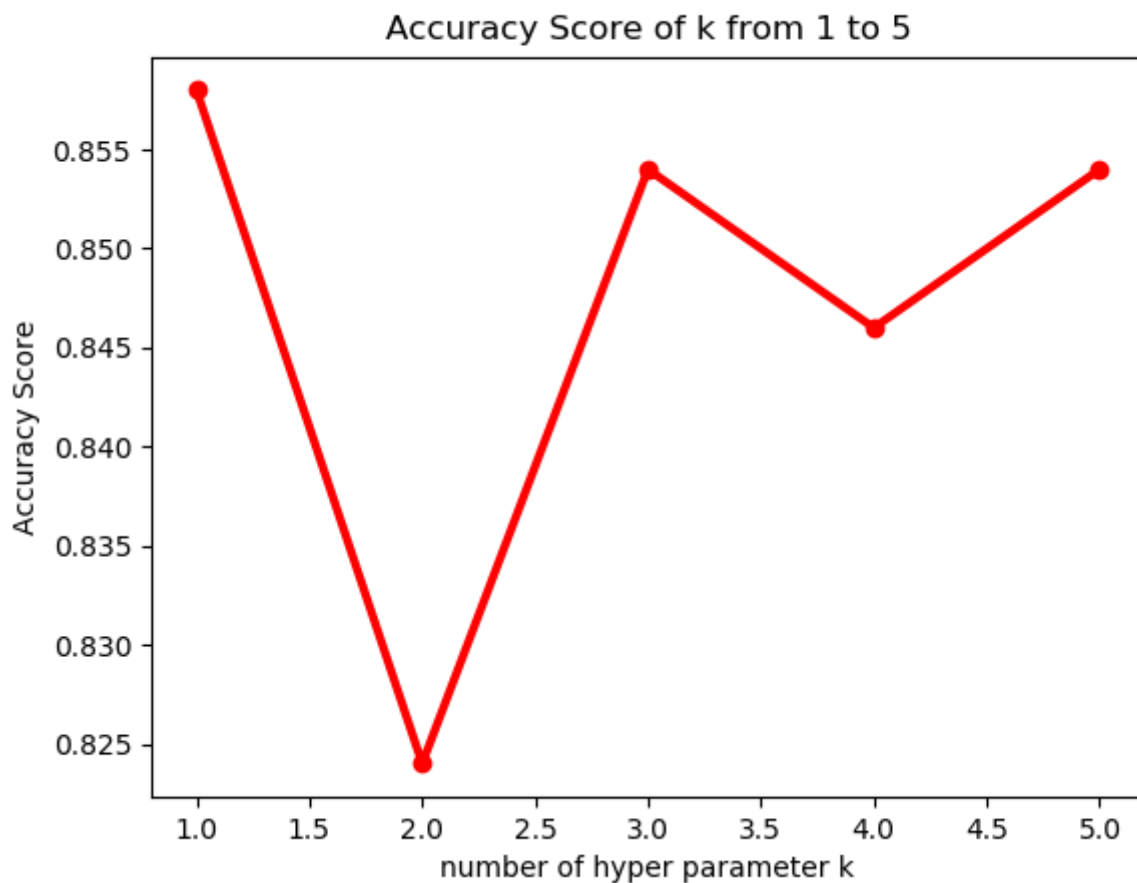
```

1  # 从csv文件获取标签的函数
2  def get_labels(csv_file, labels):
3      with open(csv_file) as temp_csv:
4          temp_reader = csv.reader(temp_csv, delimiter=' ', quotechar='|')
5          next(temp_reader, None)
6          for row in temp_reader:
7              row = [int(i) for i in row[0].split(',')]
8              labels.append(row[0])
9          print("Get " + csv_file + " labels done!")
10         return labels

```

2 实验结果

针对这个数据集, 我将 `train` 数据随机 `3:1` 划分为 `train` 和 `validate`, 并将 `k` 从 `1`到`5` 测试了一遍, 每个`k`得到的精度(`accuracy`)图如下:



由图可知, `k=1` 时, 精度值最高, 所以设置超参数 `k=1` 进行测试。

测试结果如下:

```
1 k: 1
2 accuracy: 0.8866171003717472
```

3 全部代码

```
1 import csv
2 import sys
3 import math
4 from sklearn.model_selection import train_test_split
5 from sklearn import metrics
6 import random
7 import matplotlib.pyplot as plt
8
9 # 训练标签
10 training_labels = []
11 # 训练数据
12 training_vectors = []
13
14 # 学习(预测)标签
15 learning_labels = []
```

```

16 # 学习数据
17 learning_vectors = []
18
19 # 验证标签
20 validate_labels = []
21 # 验证数据
22 validate_vectors = []
23 # 验证预测标签: 即 由knn算法为验证集预测的标签, 它与真实的验证集标签可以用来判断超
    参数k的好坏
24 validate_temp_labels = []
25 # 精度: 存放每个k进行实验后的精度
26 accuracy = []
27
28 # 学习数据的真实标签
29 truth_labels = []
30
31
32 # 这个函数用于获取一个数据可能的标签 这个列表中占比最多的标签, 如果有一样多的标签,
    则随机取
33 def get_most_label(dists, labels):
34     # 0-9所有数的出现次数, 以字典的形式保存
35     times = {0: 0, 1: 0, 2: 0, 3: 0, 4: 0, 5: 0, 6: 0, 7: 0, 8: 0, 9: 0}
36     for order in range(len(dists)):
37         if labels[order] != -1: # -1是无效标签, 要排除
38             times[labels[order]] += 1
39     tem = sorted(times.items(), key=lambda x: x[1], reverse=True) # 倒序排
    序, 也就是从大到小排序
40     most = [tem[0][0]]
41     for i in range(1, len(tem)):
42         if tem[i][1] == tem[0][1]:
43             most.append(tem[i][0])
44         else:
45             break
46     if len(most) == 1:
47         return most[0]
48     else:
49         return most[random.randint(0, len(most) - 1)]
50
51
52 # 这个函数的作用有:
53 # 1. 读取 train.csv 中的数据
54 # 2. 划分数据集(3:1 = 训练数据集:验证数据集)
55 def train():
56     global training_labels, training_vectors, validate_labels,
    validate_vectors
57     with open('train.csv') as training_csv:
58         training_reader = csv.reader(training_csv, delimiter=' ',
    quotechar='|')
59         # get rid of header of csv
60         next(training_reader, None)
61         for row in training_reader:
62             row = [int(i) for i in row[0].split(',')]

```

```

63         training_labels.append(row[0]) # 把标签加入训练标签列标中
64         training_vectors.append(row[1:]) # 把数据加入训练数据集中
65         # 划分数数据集, 训练:验证 = 3:1, 完全随机划分
66         training_vectors, validate_vectors, training_labels,
validate_labels = \
67             train_test_split(training_vectors, training_labels,
test_size=0.25, random_state=0)
68         print("Done training!")
69
70
71 # 验证, k表示knn的超参数k
72 def validate(k):
73     global validate_vectors, validate_labels, validate_temp_labels
74     i = 0
75     validate_temp_labels = [-1] * len(validate_vectors)
76     for row in validate_vectors:
77         min_dist = [sys.maxsize] * k # 长度为k的最小距离列表, 最小距离计算方法
在dist()方法中
78         min_dist_label = [-1] * k # 长度为k的标签列表, 存放可能的标签
79         for count, (vec, num) in enumerate(zip(training_vectors,
training_labels)): # 遍历训练集
80             if count % 2000 == 0: # 其实只有 count 为 0 的时候才生效
81                 print("Validating on No." + str(i))
82                 d = dist(row, vec) # 计算 row(即当前数据) 与 vec(当前测试数据) 的
距离
83                 if d < max(min_dist): # 如果距离比 min_dist 中的最大值小
84                     min_dist_label[min_dist.index(max(min_dist))] = num # 替换
对应最大距离的标签为vec的标签num
85                     min_dist[min_dist.index(max(min_dist))] = d # 替换最大距离
为当前距离d
86                 validate_temp_labels[i] = get_most_label(min_dist, min_dist_label)
# 检查最有可能的标签
87                 i += 1
88                 print("Validate done by KNN, where k = " + str(k))
89                 # 返回当前k获取的精度
90                 return metrics.accuracy_score(validate_labels, validate_temp_labels)
91
92 # 预测, 逻辑大致与validate()相同
93 def predict(k):
94     global training_labels, training_vectors, learning_labels,
learning_vectors
95     with open('test.csv') as test_csv: # 测试数据在 test.csv
96         learning_reader = csv.reader(test_csv, delimiter=' ',
quotechar='|')
97         next(learning_reader, None)
98         learning_vectors = list(learning_reader)
99         i = 0
100         learning_labels = [-1] * len(learning_vectors)
101         for row in learning_vectors:
102             row = [int(j) for j in row[0].split(',')]
103             min_dist = [sys.maxsize] * k
104             min_dist_label = [-1] * k

```

```

105         closest_num = -1
106         for count, (vec, num) in enumerate(zip(training_vectors,
training_labels)):
107             if count % 10000 == 0:
108                 print("Working on NO." + str(i))
109                 d = dist(row, vec)
110                 if d < max(min_dist):
111                     min_dist_label[min_dist.index(max(min_dist))] = num
112                     min_dist[min_dist.index(max(min_dist))] = d
113                     closest_num = get_most_label(min_dist, min_dist_label)
114                 learning_labels[i] = closest_num
115                 i += 1
116         print("KNN Done!")
117
118
119 # 把超参数k获取的结果写入一个csv中
120 def write(k):
121     global learning_labels
122     file_k = 'answers' + str(k) + '.csv'
123     with open(file_k, 'w', newline='') as csvfile:
124         answer_writer = csv.writer(csvfile, delimiter=' ', quotechar='|',
quoting=csv.QUOTE_MINIMAL)
125         answer_writer.writerow(['Answers'])
126         for label in learning_labels:
127             answer_writer.writerow(str(label))
128         print("Done writing CSV!")
129
130
131 # 计算两个向量距离的函数，其实，可以理解为 sum = abs(v1[i] - v2[i]), for i in
range(len(v1)), 注意 v1与v2长度必相同
132 def dist(vector_one, vector_two):
133     d = 0
134     for i, j in zip(vector_one, vector_two):
135         d += abs(i - j)
136     return d
137
138
139 # 从csv文件获取标签的函数
140 def get_labels(csv_file, labels):
141     with open(csv_file) as temp_csv:
142         temp_reader = csv.reader(temp_csv, delimiter=' ', quotechar='|')
143         next(temp_reader, None)
144         for row in temp_reader:
145             row = [int(i) for i in row[0].split(',')]
146             labels.append(row[0])
147         print("Get " + csv_file + " labels done!")
148     return labels
149
150
151 # 绘制评价图，可以是精度也可以是f1，两个参数中必定有一个为True，另一个为 False
152 def draw(k, score, accuracy_score=True, f1_score=False):
153     x = [i for i in range(1, k + 1)]

```



```

154     name = ''
155     if accuracy_score and not f1_score:
156         name = 'Accuracy'
157     if f1_score and not accuracy_score:
158         name = 'F1'
159     plt.plot(x, score, label=name + ' Score', linewidth=3, color='r',
160             marker='o')
161     plt.xlabel('number of hyper parameter k')
162     plt.ylabel(name + ' Score')
163     plt.title(name + ' Score of k from 1 to 5')
164     plt.show()
165
166 # 设置超参数为k, 然后进行predict(k), 再根据真实的结果(保存在truth.csv)计算超参数
167 # 为k得到的精度
168 def do_pre(k):
169     global truth_labels, learning_labels
170     predict(k)
171     truth_labels = get_labels('truth.csv', truth_labels)
172     result = metrics.accuracy_score(truth_labels, learning_labels)
173     print(result)
174
175 def main(k):
176     global accuracy
177     # step 1. train()
178     train()
179     # step 2. validate(i)
180     for i in range(1, k + 1):
181         accuracy.append(validate(i))
182     # step 3. draw()
183     draw(k, accuracy)
184     # step 4. 获取验证效果最好的k, 即 accuracy最大值的下标 + 1
185     index_max = accuracy.index(max(accuracy)) + 1
186     # step 5. 预测, 并计算预测精度
187     do_pre(index_max)
188
189
190 if __name__ == "__main__":
191     main(5)

```