

20220426-C++

[1.过程描述](#)

[2.结果输出](#)

1.过程描述

▼ knn.h C++ 复制代码

```
1  #ifndef __KNN_H
2  #define __KNN_H
3
4  ▼ #include "Common.h"
5
6  // O(k*n) where k is the number of neighbors and N is the size of
   training data
7  // O(n) + O(k*n) + k
8
9  class KNN : public CommonData
10 ▼ {
11     int k;
12     std::vector<Data *> * neighbors;
13
14     public:
15     KNN(int);
16     KNN();
17     ~KNN();
18
19     void findKnearest(Data *queryPoint);
20     void setK(int val);
21     int findMostFrequentClass();
22     double calculateDistance(Data* queryPoint, Data* input);
23     double validatePerformance();
24     double testPerformance();
25 };
26
27 #endif
```

```
1  ▼ #include "knn.h"
2  #include <cmath>
3  #include <limits>
4  #include <map>
5  #include "stdint.h"
6  #include "DataHandler.h"
7
8
9  KNN::KNN(int val)
10 ▼ {
11     k = val;
12 }
13
14 KNN::KNN()
15 ▼ {
16 }
17
18
19 KNN::~~KNN()
20 ▼ {
21     // NOTHING TO DO
22 }
23
24 void KNN::findKnearest(Data* queryPoint)
25 ▼ {
26     neighbors = new std::vector<Data*>;
27     double min = std::numeric_limits<double>::max();
28     double previousMin = min;
29     int index;
30     for (int i = 0; i < k; i++)
31 ▼     {
32         if (i == 0)
33 ▼         {
34             for (int j = 0; j < trainingData->size(); j++)
35 ▼             {
36                 double dist = calculateDistance(queryPoint,
trainingData->at(j));
37                 trainingData->at(j)->setDistance(dist);
38                 if (dist < min)
39 ▼                 {
40                     min = dist;
41                     index = j;
42                 }
43             }
44             neighbors->push_back(trainingData->at(index));
```

```

45         previousMin = min;
46         min = std::numeric_limits<double>::max();
47     }
48     else
49     {
50         for (int j = 0; j < trainingData->size(); j++)
51         {
52             double dist = trainingData->at(j)->getDistance();
53             if (dist > previousMin && dist < min)
54             {
55                 min = dist;
56                 index = j;
57             }
58         }
59         neighbors->push_back(trainingData->at(index));
60         previousMin = min;
61         min = std::numeric_limits<double>::max();
62     }
63 }
64 }
65 void KNN::setK(int val)
66 {
67     k = val;
68 }
69
70 int KNN::findMostFrequentClass()
71 {
72     std::map<uint8_t, int> frequencyMap;
73     for (int i = 0; i < neighbors->size(); i++)
74     {
75         if (frequencyMap.find(neighbors->at(i)->getLabel()) ==
frequencyMap.end())
76         {
77             frequencyMap[neighbors->at(i)->getLabel()] = 1;
78         }
79         else
80         {
81             frequencyMap[neighbors->at(i)->getLabel()]++;
82         }
83     }
84
85     int best = 0;
86     int max = 0;
87
88     for (auto kv : frequencyMap)
89     {
90         if (kv.second > max)
91         {

```

```

92         max = kv.second;
93         best = kv.first;
94     }
95 }
96 delete neighbors;
97 return best;
98
99 }
100
101 double KNN::calculateDistance(Data* queryPoint, Data* input)
102 {
103     double value = 0;
104     if (queryPoint->getNormalizedFeatureVector()->size() != input-
105 >getNormalizedFeatureVector()->size())
106     {
107         printf("Vector size mismatch.\n");
108         exit(1);
109     }
110 #ifdef EUCLID
111     for (unsigned i = 0; i < queryPoint->getNormalizedFeatureVector()-
112 >size(); i++)
113     {
114         value += pow(queryPoint->getNormalizedFeatureVector()->at(i) -
115 input->getNormalizedFeatureVector()->at(i), 2);
116     }
117     return sqrt(value);
118 #elif defined MANHATTAN
119     //do some stuff
120 #endif
121 }
122
123 double KNN::validatePerformance()
124 {
125     double current_performance = 0;
126     int count = 0;
127     int data_index = 0;
128     for (Data* queryPoint : *validationData)
129     {
130         findKnearest(queryPoint);
131         int prediction = findMostFrequentClass();
132         data_index++;
133         if (prediction == queryPoint->getLabel())
134         {
135             count++;
136         }
137         printf("Current Performance: %.3f %%\n", ((double)count) * 100.0
138 / ((double)data_index));
139     }

```

```

136     current_performance = ((double)count) * 100.0 /
((double)validationData->size());
137     printf("Validation Performance for K = %d: %.3f\n", k,
current_performance);
138     return current_performance;
139 }
140 double KNN::testPerformance()
141 {
142     double current_performance = 0;
143     int count = 0;
144     for (Data* queryPoint : *testData)
145     {
146         findKnearest(queryPoint);
147         int prediction = findMostFrequentClass();
148         if (prediction == queryPoint->getLabel())
149         {
150             count++;
151         }
152     }
153     current_performance = ((double)count) * 100.0 / ((double)testData-
>size());
154     printf("Validation Performance for K = %d: %.3f\n", k,
current_performance);
155     return current_performance;
156 }
157
158 int
159 main()
160 {
161     DataHandler* dh = new DataHandler();
162     //dh->read_csv("/home/gerardta/iris.data",",");
163     dh->readInputData("train-images.idx3-ubyte");
164     dh->readLabelData("train-labels.idx1-ubyte");
165     dh->countClasses();
166     dh->splitData();
167     KNN* nearest = new KNN();
168     nearest->setK(3);
169     nearest->setTrainingData(dh->getTrainingData());
170     nearest->setTestData(dh->getTestData());
171     nearest->setValidationData(dh->getValidationData());
172     double performance = 0;
173     double best_performance = 0;
174     int best_k = 1;
175     for (int k = 1; k <= 3; k++)
176     {
177         if (k == 1)
178         {
179             performance = nearest->validatePerformance();

```

```

180         best_performance = performance;
181     }
182     else
183     {
184         nearest->setK(k);
185         performance = nearest->validatePerformance();
186         if (performance > best_performance)
187         {
188             best_performance = performance;
189             best_k = k;
190         }
191     }
192 }
193 nearest->setK(best_k);
194 nearest->testPerformance();
195 }

```

▼ common.h

C++ | 复制代码

```

1  #ifndef __COMMON_HPP
2  #define __COMMON_HPP
3  ▼ #include "data.h"
4  #include <vector>
5  class CommonData
6  ▼ {
7  protected:
8      std::vector<Data*>> trainingData;
9      std::vector<Data*>> testData;
10     std::vector<Data*>> validationData;
11 public:
12     void setTrainingData(std::vector<Data*>> vect);
13     void setTestData(std::vector<Data*>> vect);
14     void setValidationData(std::vector<Data*>> vect);
15 };
16 #endif
17

```

```
1 ▼ #include "common.h"
2
3 void CommonData::setTrainingData(std::vector<Data*>* vect)
4 ▼ {
5     trainingData = vect;
6 }
7 void CommonData::setTestData(std::vector<Data*>* vect)
8 ▼ {
9     testData = vect;
10 }
11 void CommonData::setValidationData(std::vector<Data*>* vect)
12 ▼ {
13     validationData = vect;
14 }
```

```
1  #ifndef __DATA_H
2  #define __DATA_H
3
4  ▼ #include <vector>
5  #include "stdint.h" // uint8_t
6  #include "stdio.h"
7  class Data
8  ▼ {
9      std::vector<uint8_t>* featureVector;
10     std::vector<double>* normalizedFeatureVector;
11     std::vector<int>* classVector;
12     uint8_t label;
13     uint8_t enumeratedLabel; // A -> 1
14     double distance;
15
16     public:
17         void setDistance(double);
18         void setFeatureVector(std::vector<uint8_t>*);
19         void setNormalizedFeatureVector(std::vector<double>*);
20         void setClassVector(int counts);
21         void appendToFeatureVector(uint8_t);
22         void appendToFeatureVector(double);
23         void setLabel(uint8_t);
24         void setEnumeratedLabel(uint8_t);
25         void printVector();
26         void printNormalizedVector();
27
28         double getDistance();
29         int getFeatureVectorSize();
30         uint8_t getLabel();
31         uint8_t getEnumeratedLabel();
32
33         std::vector<uint8_t>* getFeatureVector();
34         std::vector<double>* getNormalizedFeatureVector();
35         std::vector<int> getClassVector();
36
37     };
38
39     #endif
40
41
```



```
1  ▾ #include "data.h"
2
3  void Data::setDistance(double dist)
4  ▾ {
5      distance = dist;
6  }
7  void Data::setFeatureVector(std::vector<uint8_t>* vect)
8  ▾ {
9      featureVector = vect;
10 }
11
12
13 void Data::setNormalizedFeatureVector(std::vector<double>* vect)
14 ▾ {
15     normalizedFeatureVector = vect;
16 }
17 void Data::appendToFeatureVector(uint8_t val)
18 ▾ {
19     featureVector->push_back(val);
20 }
21 void Data::appendToFeatureVector(double val)
22 ▾ {
23     normalizedFeatureVector->push_back(val);
24 }
25 void Data::setLabel(uint8_t val)
26 ▾ {
27     label = val;
28 }
29 void Data::setEnumeratedLabel(uint8_t val)
30 ▾ {
31     enumeratedLabel = val;
32 }
33
34 void Data::setClassVector(int classCounts)
35 ▾ {
36     classVector = new std::vector<int>();
37     for (int i = 0; i < classCounts; i++)
38 ▾ {
39         if (i == label)
40             classVector->push_back(1);
41         else
42             classVector->push_back(0);
43     }
44 }
45
```

```

46 void Data::printVector()
47 {
48     printf("[ ");
49     for (uint8_t val : *featureVector)
50     {
51         printf("%u ", val);
52     }
53     printf("]\n");
54 }
55
56 void Data::printNormalizedVector()
57 {
58     printf("[ ");
59     for (auto val : *normalizedFeatureVector)
60     {
61         printf("%.2f ", val);
62     }
63     printf("]\n");
64 }
65
66
67 double Data::getDistance()
68 {
69     return distance;
70 }
71
72 int Data::getFeatureVectorSize()
73 {
74     return featureVector->size();
75 }
76 uint8_t Data::getLabel()
77 {
78     return label;
79 }
80 uint8_t Data::getEnumeratedLabel()
81 {
82     return enumeratedLabel;
83 }
84
85 std::vector<uint8_t>* Data::getFeatureVector()
86 {
87     return featureVector;
88 }
89 std::vector<double>* Data::getNormalizedFeatureVector()
90 {
91     return normalizedFeatureVector;
92 }
93

```

```
94     std::vector<int> Data::getClassVector()  
95     {  
96         return *classVector;  
97     }
```

```
1  #ifndef __DATA_HANDLER_H
2  #define __DATA_HANDLER_H
3
4  ▼ #include "fstream"
5  #include "stdint.h"
6  #include "data.h"
7  #include <vector>
8  #include <string>
9  #include <map>
10 #include <unordered_set>
11 #include <math.h>
12
13 class DataHandler
14 ▼ {
15     std::vector<Data*> dataArray; // all of the data
16     std::vector<Data*> trainingData;
17     std::vector<Data*> testData;
18     std::vector<Data*> validationData;
19     int class_counts;
20     int featureVectorSize;
21     std::map<uint8_t, int> classFromInt;
22     std::map<std::string, int> classFromString; //string key
23
24     public:
25         const double TRAIN_SET_PERCENT = .1;
26         const double TEST_SET_PERCENT = .075;
27         const double VALID_SET_PERCENT = 0.005;
28
29         DataHandler();
30         ~DataHandler();
31
32         void readCsv(std::string, std::string);
33         void readInputData(std::string path);
34         void readLabelData(std::string path);
35         void splitData();
36         void countClasses();
37         void normalize();
38         void print();
39
40         int getClassCounts();
41         int getDataArraySize();
42         int getTrainingDataSize();
43         int getTestDataSize();
44         int getValidationSize();
45 }
```

```
46         uint32_t format(const unsigned char* bytes);
47
48         std::vector<Data*>* getTrainingData();
49         std::vector<Data*>* getTestData();
50         std::vector<Data*>* getValidationData();
51         std::map<uint8_t, int> getClassMap();
52
53     };
54
55 #endif
```

```
1  ▼ #include "DataHandler.h"
2  #include <algorithm>
3  #include <random>
4  #pragma warning(disable:4996)
5
6  DataHandler::DataHandler()
7  ▼ {
8      dataArray = new std::vector<Data*>;
9      trainingData = new std::vector<Data*>;
10     testData = new std::vector<Data*>;
11     validationData = new std::vector<Data*>;
12 }
13
14 DataHandler::~DataHandler()
15 ▼ {
16     // FIX ME
17 }
18
19 void DataHandler::readCsv(std::string path, std::string delimiter)
20 ▼ {
21     class_counts = 0;
22     std::ifstream data_file;
23     data_file.open(path.c_str());
24     std::string line;
25
26     while (std::getline(data_file, line))
27 ▼     {
28         if (line.length() == 0) continue;
29         Data* d = new Data();
30         d->setNormalizedFeatureVector(new std::vector<double>());
31         size_t position = 0;
32         std::string token;
33         while ((position = line.find(delimiter)) != std::string::npos)
34 ▼         {
35             token = line.substr(0, position);
36             d->appendToFeatureVector(std::stod(token));
37             line.erase(0, position + delimiter.length());
38         }
39
40         if (classFromString.find(line) != classFromString.end())
41 ▼         {
42             d->setLabel(classFromString[line]);
43         }
44         else
45 ▼         {
```

```

46         classFromString[line] = class_counts;
47         d->setLabel(classFromString[token]);
48         class_counts++;
49     }
50     dataArray->push_back(d);
51 }
52 for (Data* data : *dataArray)
53     data->setClassVector(class_counts);;
54 //normalize();
55 featureVectorSize = dataArray->at(0)->getNormalizedFeatureVector()-
>size();
56 }
57
58 void DataHandler::readInputData(std::string path)
59 {
60     uint32_t magic = 0;
61     uint32_t num_images = 0;
62     uint32_t num_rows = 0;
63     uint32_t num_cols = 0;
64
65     unsigned char bytes[4];
66     FILE* f = fopen(path.c_str(), "r");
67     if (f)
68     {
69         int i = 0;
70         while (i < 4)
71         {
72             if (fread(bytes, sizeof(bytes), 1, f))
73             {
74                 switch (i)
75                 {
76                     case 0:
77                         magic = format(bytes);
78                         i++;
79                         break;
80                     case 1:
81                         num_images = format(bytes);
82                         i++;
83                         break;
84                     case 2:
85                         num_rows = format(bytes);
86                         i++;
87                         break;
88                     case 3:
89                         num_cols = format(bytes);
90                         i++;
91                         break;
92                 }

```

```

93         }
94     }
95     printf("Done getting file header.\n");
96     uint32_t image_size = num_rows * num_cols;
97     for (i = 0; i < num_images; i++)
98     {
99         Data* d = new Data();
100        d->setFeatureVector(new std::vector<uint8_t>());
101        uint8_t element[1];
102        for (int j = 0; j < image_size; j++)
103        {
104            if (fread(element, sizeof(element), 1, f))
105            {
106                d->appendToFeatureVector(element[0]);
107            }
108        }
109        dataArray->push_back(d);
110        dataArray->back()->setClassVector(class_counts);
111    }
112    normalize();
113    featureVectorSize = dataArray->at(0)->getFeatureVector()-
>size();
114    printf("Successfully read %lu data entries.\n", dataArray-
>size());
115    printf("The Feature Vector Size is: %d\n", featureVectorSize);
116    }
117    else
118    {
119        printf("Invalid Input File Path\n");
120        exit(1);
121    }
122 }
123 void DataHandler::readLabelData(std::string path)
124 {
125     uint32_t magic = 0;
126     uint32_t num_images = 0;
127     unsigned char bytes[4];
128     FILE* f = fopen(path.c_str(), "r");
129     if (f)
130     {
131         int i = 0;
132         while (i < 2)
133         {
134             if (fread(bytes, sizeof(bytes), 1, f))
135             {
136                 switch (i)
137                 {
138                     case 0:

```



```

139             magic = format(bytes);
140             i++;
141             break;
142         case 1:
143             num_images = format(bytes);
144             i++;
145             break;
146     }
147 }
148 }
149
150 for (unsigned j = 0; j < num_images; j++)
151 {
152     uint8_t element[1];
153     if (fread(element, sizeof(element), 1, f))
154     {
155         dataArray->at(j)->setLabel(element[0]);
156     }
157 }
158
159 printf("Done getting Label header.\n");
160 }
161 else
162 {
163     printf("Invalid Label File Path\n");
164     exit(1);
165 }
166 }
167 void DataHandler::splitData()
168 {
169     std::unordered_set<int> used_indexes;
170     int train_size = dataArray->size() * TRAIN_SET_PERCENT;
171     int test_size = dataArray->size() * TEST_SET_PERCENT;
172     int valid_size = dataArray->size() * VALID_SET_PERCENT;
173
174     std::random_shuffle(dataArray->begin(), dataArray->end());
175
176     // Training Data
177
178     int count = 0;
179     int index = 0;
180     while (count < train_size)
181     {
182         trainingData->push_back(dataArray->at(index++));
183         count++;
184     }
185
186     // Test Data

```

```

187     count = 0;
188     while (count < test_size)
189     {
190         testData->push_back(dataArray->at(index++));
191         count++;
192     }
193
194     // Test Data
195
196     count = 0;
197     while (count < valid_size)
198     {
199         validationData->push_back(dataArray->at(index++));
200         count++;
201     }
202
203     printf("Training Data Size: %lu.\n", trainingData->size());
204     printf("Test Data Size: %lu.\n", testData->size());
205     printf("Validation Data Size: %lu.\n", validationData->size());
206 }
207
208 void DataHandler::countClasses()
209 {
210     int count = 0;
211     for (unsigned i = 0; i < dataArray->size(); i++)
212     {
213         if (classFromInt.find(dataArray->at(i)->getLabel()) ==
classFromInt.end())
214         {
215             classFromInt[dataArray->at(i)->getLabel()] = count;
216             dataArray->at(i)->setEnumeratedLabel(count);
217             count++;
218         }
219         else
220         {
221             dataArray->at(i)->setEnumeratedLabel(classFromInt[dataArray-
>at(i)->getLabel()]);
222         }
223     }
224
225     class_counts = count;
226     for (Data* data : *dataArray)
227         data->setClassVector(class_counts);
228     printf("Successfully Extraced %d Unique Classes.\n", class_counts);
229 }
230
231 void DataHandler::normalize()
232 {

```

```

233     std::vector<double> mins, maxs;
234     // fill min and max lists
235
236     Data* d = dataArray->at(0);
237     for (auto val : *d->getFeatureVector())
238     {
239         mins.push_back(val);
240         maxs.push_back(val);
241     }
242
243     for (int i = 1; i < dataArray->size(); i++)
244     {
245         d = dataArray->at(i);
246         for (int j = 0; j < d->getFeatureVectorSize(); j++)
247         {
248             double value = (double)d->getFeatureVector()->at(j);
249             if (value < mins.at(j)) mins[j] = value;
250             if (value > maxs.at(j)) maxs[j] = value;
251         }
252     }
253     // normalize data array
254
255     for (int i = 0; i < dataArray->size(); i++)
256     {
257         dataArray->at(i)->setNormalizedFeatureVector(new
std::vector<double>());
258         dataArray->at(i)->setClassVector(class_counts);
259         for (int j = 0; j < dataArray->at(i)->getFeatureVectorSize();
j++)
260         {
261             if (maxs[j] - mins[j] == 0) dataArray->at(i)-
>appendToFeatureVector(0.0);
262             else
263                 dataArray->at(i)->appendToFeatureVector(
264                     (double)(dataArray->at(i)->getFeatureVector()->at(j)
- mins[j]) / (maxs[j] - mins[j]));
265         }
266     }
267 }
268
269 int DataHandler::getClassCounts()
270 {
271     return class_counts;
272 }
273
274 int DataHandler::getDataArraySize()
275 {
276     return dataArray->size();

```

```

277     }
278     int DataHandler::getTrainingDataSize()
279     {
280         return trainingData->size();
281     }
282     int DataHandler::getTestDataSize()
283     {
284         return testData->size();
285     }
286     int DataHandler::getValidationSize()
287     {
288         return validationData->size();
289     }
290
291     uint32_t DataHandler::format(const unsigned char* bytes)
292     {
293         return (uint32_t)((bytes[0] << 24) |
294             (bytes[1] << 16) |
295             (bytes[2] << 8) |
296             (bytes[3]));
297     }
298
299     std::vector<Data*> DataHandler::getTrainingData()
300     {
301         return trainingData;
302     }
303     std::vector<Data*> DataHandler::getTestData()
304     {
305         return testData;
306     }
307     std::vector<Data*> DataHandler::getValidationData()
308     {
309         return validationData;
310     }
311
312     std::map<uint8_t, int> DataHandler::getClassMap()
313     {
314         return classFromInt;
315     }
316
317     void DataHandler::print()
318     {
319         printf("Training Data:\n");
320         for (auto data : *trainingData)
321         {
322             for (auto value : *data->getNormalizedFeatureVector())
323             {
324                 printf("%.3f,", value);

```

```

325         }
326         printf(" ->  %d\n", data->getLabel());
327     }
328     return;
329
330     printf("Test Data:\n");
331     for (auto data : *testData)
332     {
333         for (auto value : *data->getNormalizedFeatureVector())
334         {
335             printf("%.3f,", value);
336         }
337         printf(" ->  %d\n", data->getLabel());
338     }
339
340     printf("Validation Data:\n");
341     for (auto data : *validationData)
342     {
343         for (auto value : *data->getNormalizedFeatureVector())
344         {
345             printf("%.3f,", value);
346         }
347         printf(" ->  %d\n", data->getLabel());
348     }
349
350 }
351

```

2.结果输出

今天只看了一个C++ ML的视频，效率低下，赶紧找回状态！