

# 一、环境配置

本书中所有代码讲解与运行均以百度AI Studio作为开发平台。因此，下面我们以运行LSTM算法为例，为读者讲解如何使用百度Studio，创建机器学习项目并运行。

## 1.1 百度AI Studio简介

百度AI Studio是针对AI学习者的在线一体化学习与实训社区。平台集合了AI教程，深度学习样例工程，各领域的经典数据集，云端的超强运算及存储资源，以及比赛平台和社区。我们将依托于这一平台讲解与运行机器学习算法。

## 1.2 百度AI Studio配置

### 账号登录

我们在打开浏览器，输入网址：<https://aistudio.baidu.com/aistudio/index>，进入百度Studio，点击右上角【登录】，读者选择合适的登陆方式即可。



### 创建项目

登陆后，点击右上角的个人头像，点击【个人中心】进入个人主页，点击【项目】右侧的下拉三角，点击【创建和Fork的项目】，点击【创建项目】。



我们选择类型为Notebook，点击【下一步】，

创建项目

×

1 选择类型

2 配置环境

3 项目描述



Notebook  
在线编程、优越算力  
所见即所得



脚本任务  
高速多卡、性能强大  
运行时间更长



图形化任务  
图形拖拽、快速部署  
简单易用

下一步

配置环境选择默认的 `PaddlePaddle 2.1.0` 和 `python3.7` 即可，点击【下一步】，根据实际情况填写项目描述，点击【创建】，完成项目的创建。

< 返回

创建项目

×

1 选择类型

2 配置环境

3 项目描述

\* 项目名称

test

\* 项目标签

已选中1个标签

▼

\* 项目描述

测试

数据集

+ 添加数据集

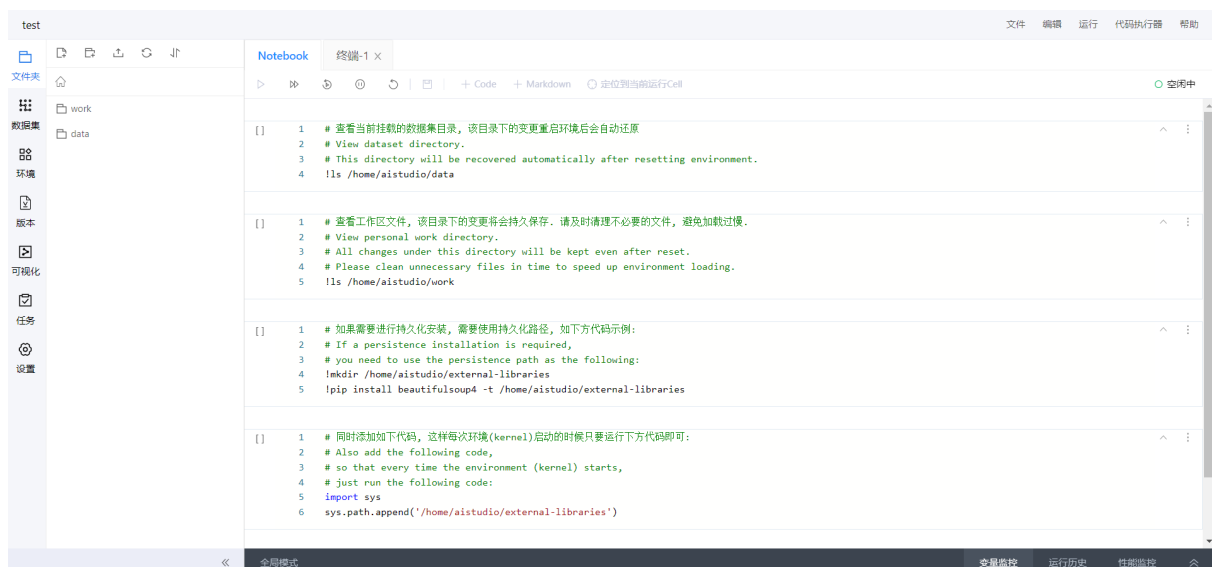
创建数据集

创建

项目创建完成后，点击【启动环境】，我们在本例中使用免费的CPU基本版，点击【确定】，



等待片刻后，我们进入如下图所示的界面：



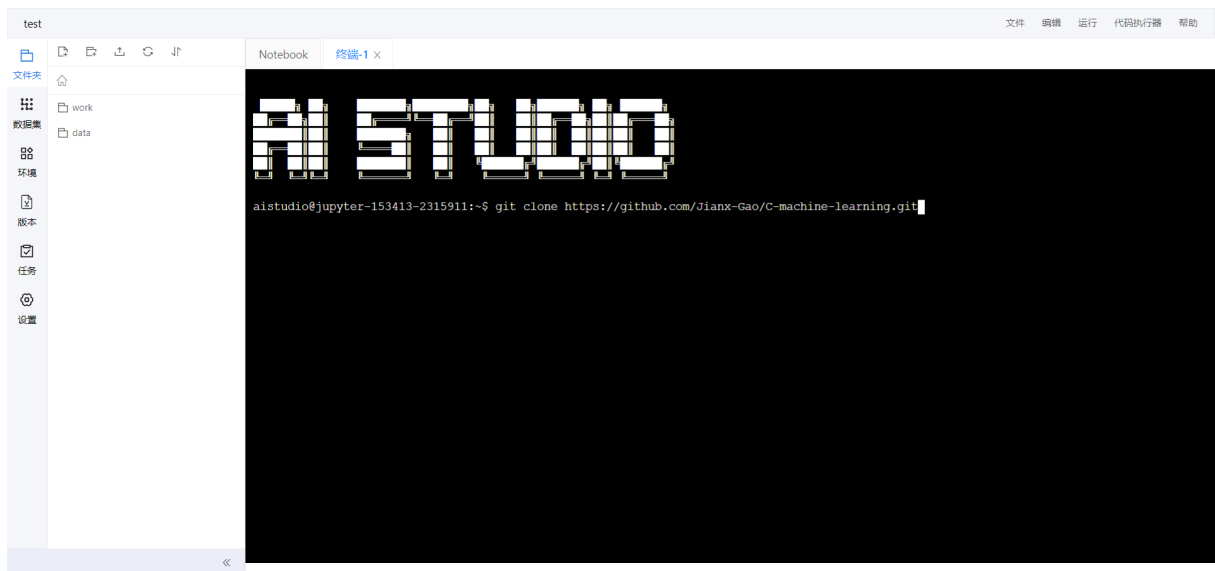
## 1.3 百度AI Studio案例运行

在该部分中，我们从github中将所有代码下载至AI Studio，并安装所有依赖，再通过随机森林算法演示如何在AI Studio运行机器学习代码。

### 下载代码

我们点击【终端】，输入指令：

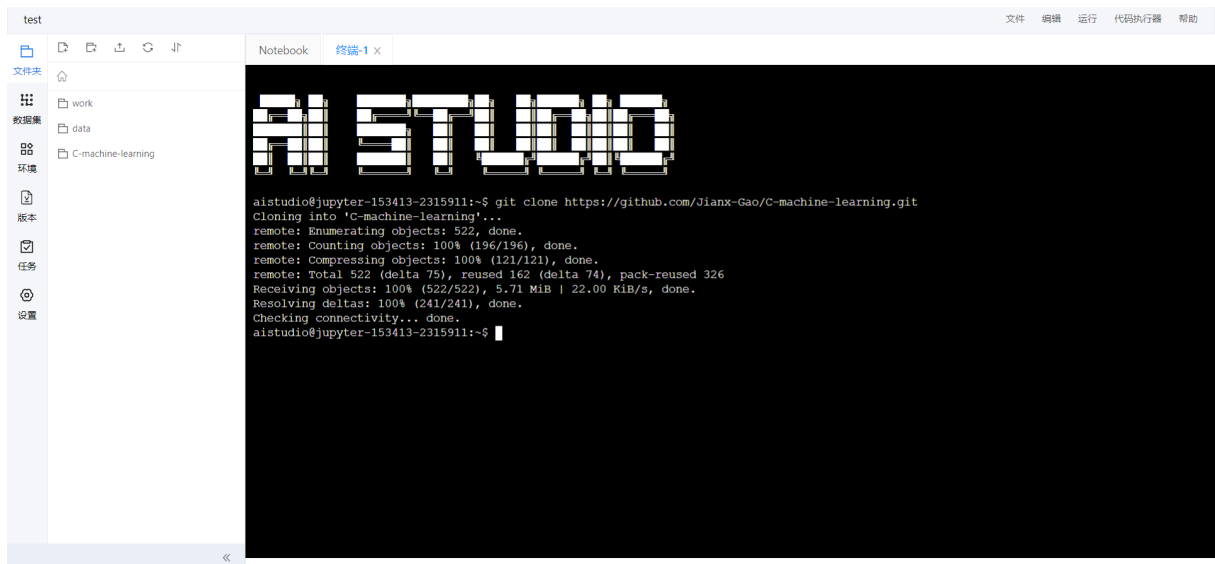
```
1 | git clone https://github.com/Jianx-Gao/C-machine-learning.git
```



如果下载失败，我们可以从github的镜像中下载代码：

```
1 | git clone https://github.com.cnpmjs.org/Jianx-Gao/C-machine-learning.git
```

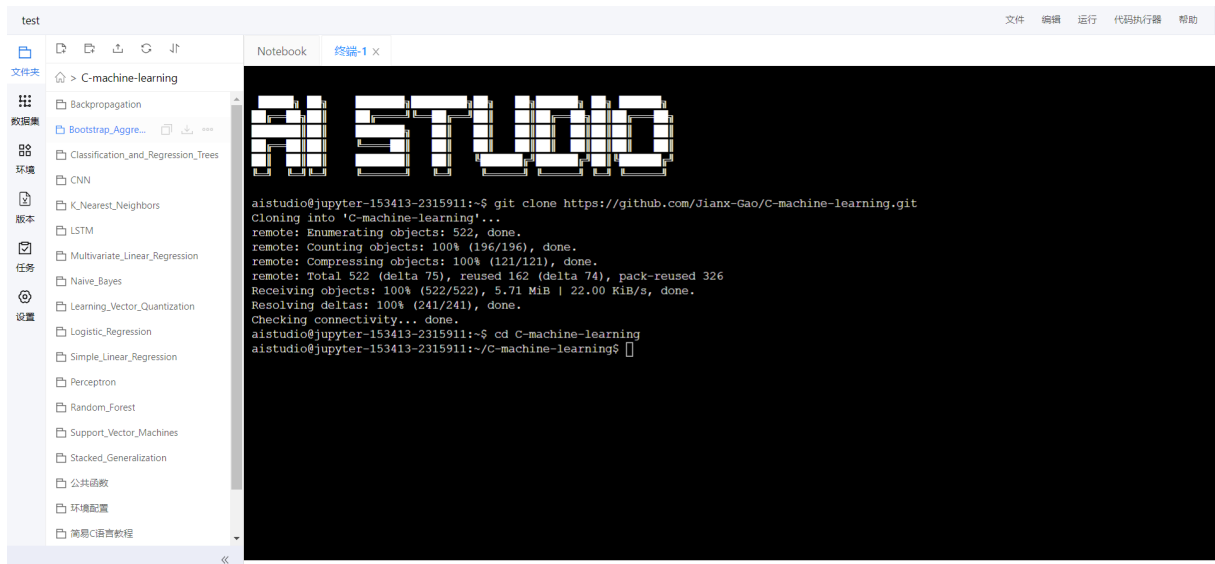
下载完成后，如下图所示：



## 环境配置

我们切换目录，进入 `cd C-machine-learning` 文件夹，同时点击左侧【C-machine-learning】文件夹

```
1 | cd C-machine-learning
```



我们安装 `requirements.txt` 中写的所有算法的依赖

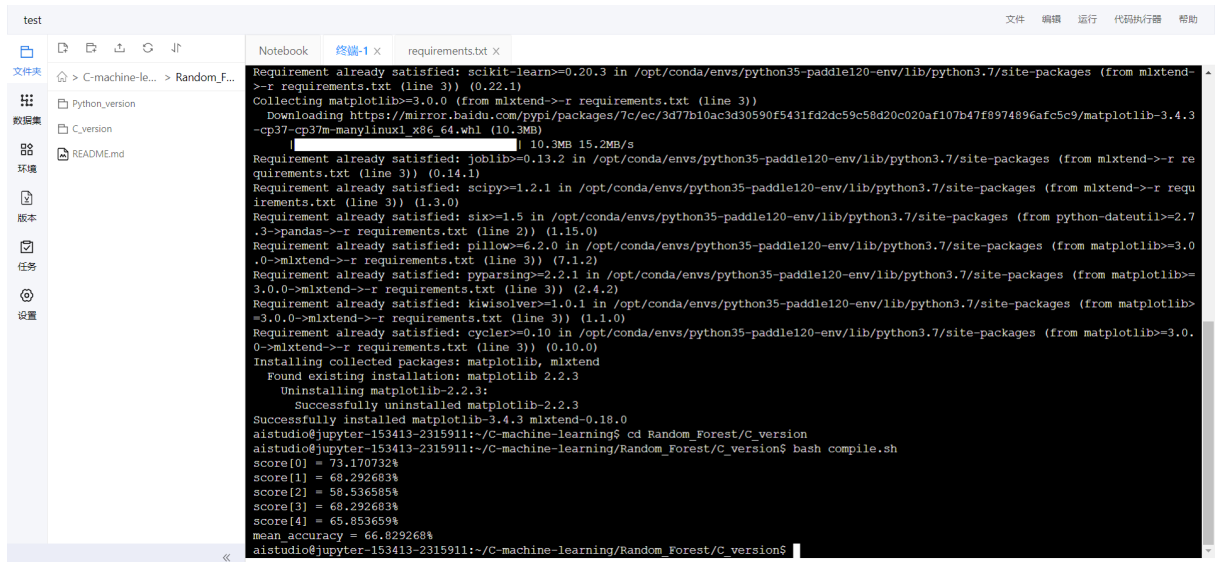
```
1 | pip install -r requirements.txt
```

## 案例运行

下面，我们就可以运行我们的案例来测试代码了，我们分别测试C语言和Python版本的随机森林算法，检验环境是否安装成功

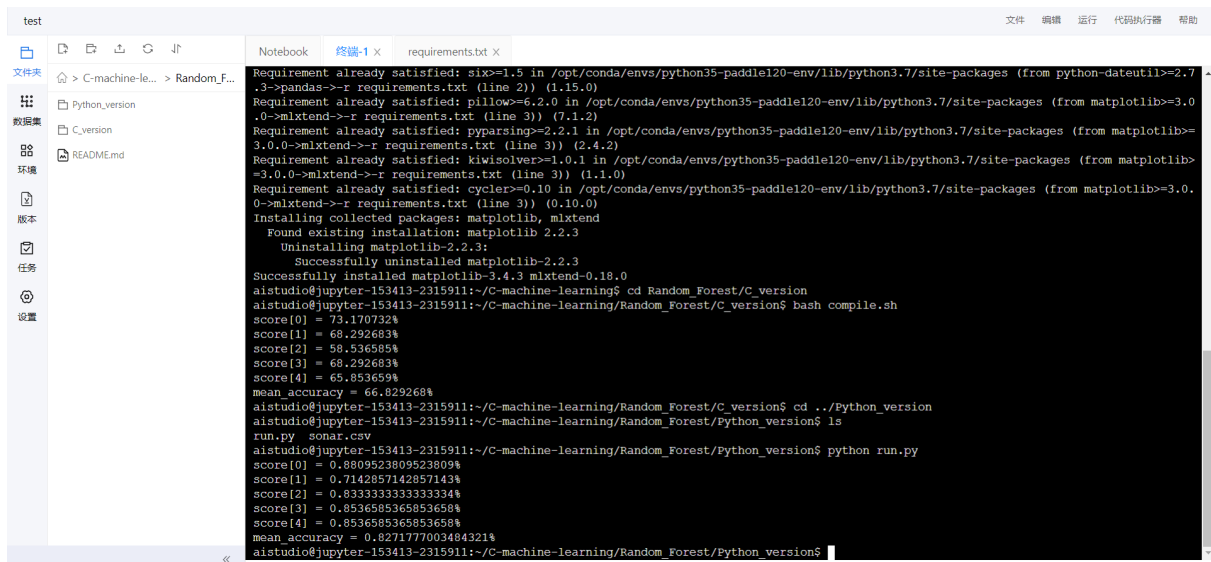
### C语言版本

```
1 | cd Random_Forest/C_version
2 | bash compile.sh
```



### Python语言版本

```
1 | cd ../Python_version
2 | python run.py
```



## 二、公共函数

在本书中，我们把多数机器-深度学习算法需要使用的函数提取出来，作为公共函数，在后续讲解中，直接调用，或简单修改即可。本章节中，我们将详细讲解这些函数的原理与实现，我们将按照如下顺序进行讲解：读取csv文件数据、数据K折交叉验证、数据标准化、计算算法结果、评价验证算法结果。

### 2.1 读取csv文件数据

在训练模型之前，应该先获得要用的数据集，数据集通常需要保存在文本文件中，这就要求我们要对数据集进行读取，以下 `get_row()`、`get_col()` 函数可以实现对csv文件的行列数的读取，它们需要字符串类型的文件名作为参数输入，而 `get_two_dimension()` 可以对csv文件的内容进行读取，它需要字符串类型的文件名作为输入参数。

- 输入：字符串（文件名）
- 输出：csv文件的行数、列数及数据的二维数组
- 功能——读取csv文件

```
1  int get_row(char *filename) //获取行数
2  {
3      char line[1024];
4      int i = 0;
5      FILE *stream = fopen(filename, "r");
6      while (fgets(line, 1024, stream))
7      {
8          i++;
9      }
10     fclose(stream);
11     return i;
12 }
13
14 int get_col(char *filename) //获取列数
15 {
16     char line[1024];
17     int i = 0;
18     FILE *stream = fopen(filename, "r");
19     fgets(line, 1024, stream);
20     char *token = strtok(line, ",");
21     while (token)
22     {
23         token = strtok(NULL, ",");
24         i++;
25     }
26     fclose(stream);
```

```

27     return i;
28 }
29
30 void get_two_dimension(char *line, double **data, char *filename)
31 {
32     FILE *stream = fopen(filename, "r");
33     int i = 0;
34     while (fgets(line, 1024, stream)) //逐行读取
35     {
36         int j = 0;
37         char *tok;
38         char *tmp = strdup(line);
39         for (tok = strtok(line, ","); tok && *tok; j++, tok = strtok(NULL, ",\n"))
40         {
41             data[i][j] = atof(tok); //转换成浮点数
42         } //字符串拆分操作
43         i++;
44         free(tmp);
45     }
46     fclose(stream); //文件打开后要进行关闭操作
47 }

```

我们创建 `data.csv` 文件，添加如下数据集：

1	1	12.2	12.5	11.1
2	2.5	555.2	121.4	2.1

调用函数：

```

1  #include <stdio.h>
2  #include <string.h>
3  #include <stdlib.h>
4  #include <math.h>
5
6  int main()
7  {
8      char filename[] = "data.csv";
9      char line[1024];
10     double **data;
11     int row, col;
12     row = get_row(filename);
13     col = get_col(filename);
14     data = (double **)malloc(row * sizeof(double *));
15     for (int i = 0; i < row; ++i){
16         data[i] = (double *)malloc(col * sizeof(double));
17     }
18     get_two_dimension(line, data, filename);
19     printf("row = %d\n", row);
20     printf("col = %d\n", col);
21
22     int i, j;
23     for(i=0; i<row; i++){
24         for(j=0; j<col; j++){
25             printf("%f\t", data[i][j]);
26         }
27         printf("\n");
28     }
29 }

```

得到结果如下：

```

1 row = 2
2 col = 4
3 1.000 12.2.000 12.5.000 11.1.000
4 2.5.000 555.2.000 121.4.000 2.1.000

```

## 2.2 数据K折交叉验证

K折交叉验证,将数据集等比例划分成K份,以其中的一份作为测试数据,其他的K-1份数据作为训练数据。然后,这样算是一次实验,而K折交叉验证只有实验K次才算完成完整的一次,也就是说交叉验证实际是把实验重复做了K次,每次实验都是从K个部分选取一份不同的数据部分作为测试数据(保证K个部分的数据都分别做过测试数据),剩下的K-1个当作训练数据,最后把得到的K个实验结果进行平分。

此函数则用于将原始数据划分为k等份,以用于k折交叉验证。

- 输入: 二维数组数据集, 数据集行数, 交叉验证折数, 交叉验证数组长度
- 输出: 划分后的三维数组
- 功能——划分数据为k折

```

1 double*** cross_validation_split(double **dataset, int row, int n_folds, int
  fold_size)
2 {
3     srand(10); // 随机种子
4     double ***split;
5     int i = 0, j = 0, k = 0;
6     int index;
7     int num;
8     num = row / n_folds;
9     double **fold;
10    split = (double***)malloc(n_folds*sizeof(double**));
11    for(i = 0; i < n_folds; i++)
12    {
13        fold = (double**)malloc(num * sizeof(double *));
14        while(j<num)
15        {
16            fold[j] = (double*)malloc(fold_size * sizeof(double));
17            index = rand() % row;
18            fold[j] = dataset[index];
19            for(k = index; k < row - 1; k++)//for循环删除数组中被rand取到的元素
20            {
21                dataset[k] = dataset[k + 1];
22            }
23            row--; //每次随机取出一个后总行数-1, 保证不会重复取某一行
24            j++;
25        }
26        j = 0;//清零j
27        split[i] = fold;
28    }
29    return split;
30 }

```

我们运行如下代码, 测试函数:



```

1 double data[6][2];
2 double *data_ptr[6];
3 double *** split;
4 for(int i=0;i<6;i++)
5 {
6     for(int j=0;j<2;j++)
7     {
8         data[i][j]=i+j;
9     }
10    data_ptr[i] = data[i];
11 };
12 split = cross_validation_split(data_ptr,6,3,2);
13 printf("%f",split[0][0][1]);

```

结果如下：

```
1 | 6.0000
```

## 2.3 数据标准化

数据标准化（归一化）处理是数据挖掘的一项基础工作，不同评价指标往往具有不同的量纲和量纲单位，这样的情况会影响到数据分析的结果，为了消除指标之间的量纲影响，需要进行数据标准化处理，以解决数据指标之间的可比性。原始数据经过数据标准化处理后，各指标处于同一数量级，适合进行综合对比评价。

### min-max Normalization

也称为离差标准化，是对原始数据的线性变换，使结果值映射到[0 - 1]之间。转换函数如下：

$$x^* = \frac{x - \min(x)}{\max(x) - \min(x)}$$

其中max为样本数据的最大值，min为样本数据的最小值。这种方法有个缺陷就是当有新数据加入时，可能导致max和min的变化，需要重新定义。

- 输入：二维数组数据集，数据集行数，数据集列数
- 输出：标准化过后的二维数组数据集
- 功能——数据集标准化

```

1 void normalize_dataset(double **dataset,int row, int col)
2 {
3     // 先 对列循环
4     double maximum, minimum;
5     for (int i = 0; i < col; i++)
6     {
7         // 第一行为标题，值为0，不能参与计算最大最小值
8         maximum = dataset[0][i];
9         minimum = dataset[0][i];
10        //再 对行循环
11        for (int j = 0; j < row; j++)
12        {
13            maximum = (dataset[j][i]>maximum)?dataset[j][i]:maximum;
14            minimum = (dataset[j][i]<minimum)?dataset[j][i]:minimum;
15        }
16        // 归一化处理
17        for (int j = 0; j < row; j++)
18        {
19            dataset[j][i] = (dataset[j][i] - minimum) / (maximum - minimum);
20        }
21    }
22 }

```

我们使用如下数据调用函数

1	0.000000	1.000000	2.000000	3.000000	4.000000	5.000000
2	10.000000	11.000000	12.000000	13.000000	14.000000	15.000000
3	20.000000	21.000000	22.000000	23.000000	24.000000	25.000000
4	30.000000	31.000000	32.000000	33.000000	34.000000	35.000000
5	40.000000	41.000000	42.000000	43.000000	44.000000	45.000000
6	50.000000	51.000000	52.000000	53.000000	54.000000	55.000000
7	60.000000	61.000000	62.000000	63.000000	64.000000	65.000000
8	70.000000	71.000000	72.000000	73.000000	74.000000	75.000000
9	80.000000	81.000000	82.000000	83.000000	84.000000	85.000000
10	90.000000	91.000000	92.000000	93.000000	94.000000	95.000000

归一化:

1	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
2	0.111111	0.111111	0.111111	0.111111	0.111111	0.111111
3	0.222222	0.222222	0.222222	0.222222	0.222222	0.222222
4	0.333333	0.333333	0.333333	0.333333	0.333333	0.333333
5	0.444444	0.444444	0.444444	0.444444	0.444444	0.444444
6	0.555556	0.555556	0.555556	0.555556	0.555556	0.555556
7	0.666667	0.666667	0.666667	0.666667	0.666667	0.666667
8	0.777778	0.777778	0.777778	0.777778	0.777778	0.777778
9	0.888889	0.888889	0.888889	0.888889	0.888889	0.888889
10	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000

## 2.4 计算算法结果

我们将训练集、测试集、学习率，epoch数，交叉验证fold的长度输入函数，调用算法，利用函数中的模型框架对测试集进行分类，或者回归，并返回预测结果。

- 输入：训练集、测试集、学习率，epoch数，数组长度
- 输出：预测结果
- 功能——计算算法结果

函数如下：

```
1 double get_test_prediction(double **train, double **test, double l_rate, int n_epoch,
2 int fold_size)
3 {
4     double *weights = (double*)malloc(col * sizeof(double));
5     // weights数组的长度就是列数（少一个结果位，多一个bias）
6     double *predictions = (double*)malloc(fold_size * sizeof(double));
7     // 预测集的行数就是数组prediction的长度
8     weights = train_weights(train, l_rate, n_epoch);
9     int i;
10    for(i = 0; i < fold_size; i++)
11    {
12        predictions[i] = predict(test[i], weights);
13    }
14    return predictions; // 返回对test的预测数组
15 }
```

## 2.5 评价验证算法结果

## 2.5.1 计算RMSE

衡量预测值与真实值之间的偏差。常用来作为机器学习模型预测结果衡量的标准。可以由以下公式计算：

$$RMSE(X, h) = \sqrt{\frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i)^2} \quad (1.8)$$

以下 `rmse_metric()` 实现了对测试集RMSE的计算，它需要真实值数组、预测值数组及交叉验证fold的长度作为输入参数。

- 输入：真实值数组，预测值数组，数组长度
- 输出：RMSE值
- 功能——计算算法在测试集的RMSE

```
1 double rmse_metric(double *actual, double *predicted, int fold_size)
2 {
3     double sum_err = 0.0;
4     int i;
5     int len = sizeof(actual)/sizeof(double);
6     for (i = 0; i < fold_size; i++)
7     {
8         double err = predicted[i] - actual[i];
9         sum_err += err * err;
10    }
11    double mean_err = sum_err / len;
12    return sqrt(mean_err);
13 }
```

我们利用如下数据，调用函数：

```
1 double act[] = {1, 2, 3, 4};
2 double pre[] = {4, 3, 2, 1};
3 printf("%f", rmse_metric(act, pre, 3));
```

计算得到得到RMSE值：

```
1 3.316625
```

## 2.5.2 计算准确率

该函数用于计算预测所得到的结果的准确率，其基本原理为：将预测正确的结果记为1，错误为0，最终求和得到正确结果个数，利用此个数除以总个数，从而得到正确率。

- 输入：真实值数组，预测值数组，数组长度
- 输出：准确率
- 功能——计算准确率

```
1 double accuracy_metric(double *actual, double *predicted, int fold_size)
2 {
3     int correct = 0;
4     int i;
5     int len = sizeof(actual);
6     for (i = 0; i < fold_size; i++)
7     {
8         if (actual[i] == predicted[i])
9             correct += 1;
10    }
11    return (correct / (double)len)*100.0;
12 }
```

我们利用如下数据，调用函数：

```
1 double act[] = {1.0, 2.0, 3.0, 4.0};
2 double pre[] = {1.0, 2.0, 3.0, 3.0};
3 printf("%f", (accuracy_metric(act, pre, 4)));
```

计算得到得到准确率：

```
1 75.0000
```

## 2.5.3 整体算法框架

- 输入：训练集、测试集、学习率，epoch数，交叉验证折数，数组长度
- 输出：准确率
- 功能——调用完整算法

```
1 #include<stdlib.h>
2 #include<stdio.h>
3
4 extern double* get_test_prediction(double **train, int train_size, double **test, int
test_size, int col, int num_neighbors, double l_rate, int n_epoch);
5 extern double accuracy_metric(double *actual, double *predicted, int fold_size);
6 extern double*** cross_validation_split(double **dataset, int row, int col, int
n_folds, int fold_size);
7
8 void evaluate_algorithm(double **dataset, int row, int col, int n_folds, int
num_neighbors, double l_rate, int n_epoch) {
9     int fold_size = (int)row / n_folds;
10    double ***split = cross_validation_split(dataset, row, n_folds, fold_size, col);
11    int i, j, k, l;
12    int test_size = fold_size;
13    int train_size = fold_size * (n_folds - 1);
14    double* score = (double*)malloc(n_folds * sizeof(double));
15
16    for (i = 0; i < n_folds; i++) {
17        double*** split_copy = (double***)malloc(n_folds * sizeof(double**));
18        for (j = 0; j < n_folds; j++) {
19            split_copy[j] = (double**)malloc(fold_size * sizeof(double*));
20            for (k = 0; k < fold_size; k++) {
21                split_copy[j][k] = (double*)malloc(col * sizeof(double));
22            }
23        }
24        for (j = 0; j < n_folds; j++)
25        {
26            for (k = 0; k < fold_size; k++)
27            {
28                for (l = 0; l < col; l++)
29                {
30                    split_copy[j][k][l] = split[j][k][l];
31                }
32            }
33        }
34        double** test_set = (double**)malloc(test_size * sizeof(double*));
35        for (j = 0; j < test_size; j++) {
36            test_set[j] = (double*)malloc(col * sizeof(double));
37            for (k = 0; k < col; k++) {
38                test_set[j][k] = split_copy[i][j][k];
39            }
40        }
41        for (j = i; j < n_folds - 1; j++) {
42            split_copy[j] = split_copy[j + 1];
```

```

43     }
44     double** train_set = (double**)malloc(train_size * sizeof(double*));
45     for (k = 0; k < n_folds - 1; k++) {
46         for (l = 0; l < fold_size; l++) {
47             train_set[k*fold_size + l] = (double*)malloc(col * sizeof(double));
48             train_set[k*fold_size + l] = split_copy[k][l];
49         }
50     }
51     double* predicted = (double*)malloc(test_size * sizeof(double));
52     predicted = get_test_prediction(train_set, train_size, test_set, test_size,
col, num_neighbors, l_rate, n_epoch);
53     double* actual = (double*)malloc(test_size * sizeof(double));
54     for (l = 0; l < test_size; l++) {
55         actual[l] = test_set[l][col - 1];
56     }
57     double acc = accuracy_metric(actual, predicted, test_size);
58     score[i] = acc;
59     printf("Scores[%d] = %f%%\n", i, score[i]);
60     free(split_copy);
61 }
62 double total = 0;
63 for (l = 0; l < n_folds; l++) {
64     total += score[l];
65 }
66 printf("mean_accuracy = %f%%\n", total / n_folds);
67 }

```

## 三、算法详解

在本章节中，我们将详细展开讲解算法，从最基础的简单线性回归算法讲起，一直到堆栈泛化算法，囊括线性算法，非线性算法，集成算法我们会从原理讲解每一个算法，并用代码实现加深读者的理解。

### 3.1 Simple Linear Regression

线性回归是一种已有200多年历史的预测方法。简单线性回归是一种可以由训练集来估计属性的机器学习算法，它相对简单，便于初学者理解。在本节中，您将看到如何用C语言一步步实现这个算法。

#### 3.1.1 算法介绍

线性回归可以在输入变量 (X) 与输出变量 (Y) 之间建立一种线性关系。具体的说，输出变量 (Y) 可以由输入变量 (X) 的线性组合计算出来。当输入变量为单一变量时，这种算法就叫做简单线性回归。

在简单线性回归中，我们可以使用训练数据的统计量来估计模型对新数据预测所需的系数。

一个简单线性回归模型可以写成：

$$y = b_0 + b_1 \times x \quad (1.1)$$

其中， $B_0$ 和 $B_1$ 为我们需要从训练集中估计的系数。得到系数后，我们可以利用此方程估计新的输入变量 (X) 对应的输出变量 (Y)。估计系数前，我们需要计算训练集的一些统计量，如平均值、方差和协方差等。

当计算出所需的统计量后，我们可以通过如下公式计算 $B_0$ ， $B_1$ ：

$$B_1 = \frac{\sum_{i=1}^n ((x_i - \text{mean}(x)) \times (y_i - \text{mean}(y)))}{\sum_{i=1}^n (x_i - \text{mean}(x))^2} \quad (1.2)$$

$$B_0 = \text{mean}(y) - B_1 \times \text{mean}(x) \quad (1.3)$$

其中， $i$ 表示训练集中的第 $i$ 个输入变量 $x$ 或输出变量 $y$ 。

#### 3.1.2 算法讲解

本小节中，我们将通过代码来讲解算法。实现简单线性回归算法的步骤可以分为如下7个部分：

- 读取csv
- 计算均值、方差及协方差
- 估计回归系数
- 由回归系数计算测试集的预测值
- 按划分的 $k$ 折交叉验证计算预测所得准确率
- 按划分的训练集和测试集计算预测所得RMSE
- `main`函数设定以上函数的参数并调用

其中，读取csv、划分数据为 $k$ 折、按划分的 $k$ 折交叉验证计算预测所得准确率参考以前章节。

#### 计算均值方差等统计量

##### 计算均值

输入变量 (X) 与输出变量 (Y) 的均值可以由以下公式得到：

$$\text{mean}(x) = \frac{\sum_{i=1} x_i}{\text{count}(x)} \quad (1.4)$$

其中，`count(x)` 表示 $x$ 的个数。

以下 `mean()` 函数可以计算一组数据的均值，它需要一维数组、数组长度作为参数。

```

1 double mean(double* values, int length)
2 {
3     //对一维数组求均值
4     int i;
5     double sum = 0.0;
6     for (i = 0; i < length; i++) {
7         sum += values[i];
8     }
9     double mean = (double)(sum / length);
10    return mean;
11 }

```

## 计算方差

方差是每个值与均值之差的平方和。一组数字的方差可计算为:

$$variance = \sum_{i=1}^n (x_i - mean(x))^2 \quad (1.5)$$

以下 `variance()` 函数可以计算一组数据的方差，它需要一维数组变量、数组的均值、以及输出数组的长度作为参数。

```

1 double variance(double* values, double mean, int length) {
2     //这里求的是平方和，没有除以n
3     double sum = 0.0;
4     int i;
5     for (i = 0; i < length; i++) {
6         sum += (values[i] - mean)*(values[i] - mean);
7     }
8     return sum;
9 }

```

我们利用以下数据集：

```

1 double x[5] = {1, 2, 4, 3, 5};
2 printf("%F", mean(x, 5));
3 printf("%F", variance(x, mean(x, 5), 5))

```

得到结果如下:

```

1 3.000
2 10.000

```

## 计算协方差

协方差在概率论和统计学中用于衡量两个变量的总体误差。而方差是协方差的一种特殊情况，即当两个变量是相同的情况。

协方差表示的是两个变量的总体的误差，这与只表示一个变量误差的方差不同。如果两个变量的变化趋势一致，那么两个变量之间的协方差就是正值。如果两个变量的变化趋势相反，那么两个变量之间的协方差就是负值。

我们可以通过以下公式来计算两个变量的协方差：

$$covariance = \sum_{i=1}^n ((x_i - mean(x)) \times (y_i - mean(y))) \quad (1.6)$$

以下 `covariance()` 函数可以计算两组数据的协方差，它需要输入数组变量 (X)、输入数组的均值、输出数组变量 (Y)、输出数组的均值、数组长度作为参数。

```

1 double covariance(double* x, double mean_x, double* y, double mean_y, int length) {
2     double cov = 0.0;
3     int i = 0;
4     for (i = 0; i < length; i++) {
5         cov += (x[i] - mean_x)*(y[i] - mean_y);
6     }
7     return cov;
8 }

```

我们利用以下数据：

```

1 double x[5] = {1, 2, 4, 3, 5};
2 double y[5] = {1, 3, 3, 2, 5};
3 printf("%f", covariance(x, mean(x, 5), y, mean(y, 5), 5));

```

得到如下结果：

```

1 8.000

```

## 估计回归系数

在简单线性回归中，我们需要估计两个系数的值。第一个是B1，可以利用公式(1.2)估计。

我们可以简化这个公式：

$$B_1 = \frac{\text{covariance}(x, y)}{\text{variance}(x)} \quad (1.7)$$

我们已经有了计算协方差和方差的函数。接下来，我们需要估计B0的值，也称为截距。可以利用公式(1.3)。

以下 `coefficients()` 函数将计算B0、B1并将其存在名为coef的数组。它需要训练集（二维数组），存储B0、B1的数组以及训练集数组长度作为参数。

```

1 //由均值方差估计回归系数
2 void coefficients(double** data, double* coef, int length) {
3     double* x = (double*)malloc(length * sizeof(double));
4     double* y = (double*)malloc(length * sizeof(double));
5     int i;
6     for (i = 0; i < length; i++) {
7         x[i] = data[i][0];
8         y[i] = data[i][1];
9     }
10    double x_mean = mean(x, length);
11    double y_mean = mean(y, length);
12    coef[1] = covariance(x, x_mean, y, y_mean, length) / variance(x, x_mean, length);
13    coef[0] = y_mean - coef[1] * x_mean;
14    for (i = 0; i < 2; i++) {
15        printf("coef[%d]=%f\n", i, coef[i]);
16    }
17 }

```

我们利用如下数据：



```

1 double data[3][2] = {
2     {1,1},
3     {2,2},
4     {3,3}
5 };
6 double coef[2] = {1,1};
7 double* dataptr[3];
8 dataptr[0] = data[0];
9 dataptr[1] = data[1];
10 dataptr[2] = data[2];
11 coefficients(dataptr, coef, 3);

```

coef作为输入的数组，经过函数操作后输出得到如下结果:

```

1 Ccoef[0] = 0.000000
2 coef[1] = 1.000000

```

## 计算测试集的预测值

简单线性回归模型是一条由训练数据估计的系数定义的直线。系数估计出来后，我们就可以用它们来进行预测。用简单的线性回归模型进行预测的方程为公式(1.1)。

以下 `get_test_prediction()` 函数实现了对数据集的预测，它需要数据行数、列数、训练集、测试集、K折交叉验证数组大小作为输入参数。

```

1 double* get_test_prediction(int col,int row,double** train, double** test, int
  n_folds) {
2     double* coef = (double*)malloc(col * sizeof(double));
3     int i;
4     for (i = 0; i < col; i++) {
5         coef[i] = 0.0;
6     }
7     int fold_size = (int)row / n_folds;
8     int train_size = fold_size * (n_folds - 1);
9     coefficients(train, coef, train_size);
10    double* predictions = (double*)malloc(fold_size * sizeof(double));
11    for (i = 0; i < fold_size; i++) {
12        predictions[i] = coef[0] + coef[1] * test[i][0];
13    }
14    return predictions;
15 }

```

## 计算预测准确率

将数据集按划分的k折交叉验证计算预测所得准确率，以下 `evaluate_algorithm()` 函数需要训练集、测试集的二维数组，学习率，epoch数，交叉验证折数，交叉验证fold的长度作为参数输入。

```

1 double evaluate_algorithm(double **dataset, int n_folds, int fold_size, double
  l_rate, int n_epoch)
2 {
3     double*** split = cross_validation_split(double **dataset, int row, int n_folds,
  int fold_size);
4     int i, j, k, l;
5     int test_size = fold_size;
6     int train_size = fold_size * (n_folds - 1); //train_size个一维数组
7     double* score = (double*)malloc(n_folds * sizeof(float));
8     for (i = 0; i < n_folds; i++)
9     {
10        //因为要遍历删除，所以拷贝一份split
11        double*** split_copy = (double***)malloc(n_folds * sizeof(double**));

```

```

12     for (j = 0; j < n_folds; j++) {
13         split_copy[j] = (double**)malloc(fold_size * sizeof(double*));
14         for (k = 0; k < fold_size; k++) {
15             split_copy[j][k] = (double*)malloc(column * sizeof(double));
16         }
17     }
18     for (j = 0; j < n_folds; j++)
19     {
20         for (k = 0; k < fold_size; k++)
21         {
22             for (l = 0; l < column; l++)
23             {
24                 split_copy[j][k][l] = split[j][k][l];
25             }
26         }
27     }
28     double** test_set = (double**)malloc(test_size * sizeof(double*));
29     for (j = 0; j < test_size; j++) { //对test_size中的每一行
30         test_set[j] = (double*)malloc(column * sizeof(double));
31         for (k = 0; k < column; k++) {
32             test_set[j][k] = split_copy[i][j][k];
33         }
34     }
35     for (j = i; j < n_folds - 1; j++) {
36         split_copy[j] = split_copy[j + 1];
37     }
38     double** train_set = (double**)malloc(train_size * sizeof(double*));
39     for (k = 0; k < n_folds - 1; k++) {
40         for (l = 0; l < fold_size; l++) {
41             train_set[k*fold_size + l] = (double*)malloc(column *
sizeof(double));
42             train_set[k*fold_size + l] = split_copy[k][l];
43         }
44     }
45     double* predicted = (double*)malloc(test_size * sizeof(double)); //predicted有
test_size个
46     predicted = get_test_prediction(train_set, test_set, l_rate, n_epoch,
fold_size);
47     double* actual = (double*)malloc(test_size * sizeof(double));
48     for (l = 0; l < test_size; l++) {
49         actual[l] = test_set[l][column - 1];
50     }
51     double accuracy = accuracy_metric(actual, predicted, test_size);
52     score[i] = accuracy;
53     printf("score[%d]=%f\n", i, score[i]);
54     free(split_copy);
55 }
56 double total = 0.0;
57 for (l = 0; l < n_folds; l++) {
58     total += score[l];
59 }
60 printf("mean_accuracy=%lf\n", total / n_folds);
61 return score;
62 }

```

### 3.1.3 算法代码

我们现在知道了如何实现简单线性回归算法，那么我们把应用应用到瑞典保险数据集 [insurance.csv](https://aistudio.baidu.com/aistudio/datasetdetail/105756/0)

我们给出链接：<https://aistudio.baidu.com/aistudio/datasetdetail/105756/0>

## C语言细节讲解

本节假设您已下载数据集 `insurance.csv`，并且它在当前工作目录中可用。下面我们给出一个完整实例，使用C语言详细讲解每一处细节。我们给出每一个.c文件的所有代码：

### 1) read\_csv.c

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <string.h>
4
5  int get_row(char *filename) //获取行数
6  {
7      char line[1024];
8      int i = 0;
9      FILE *stream = fopen(filename, "r");
10     while (fgets(line, 1024, stream))
11     {
12         i++;
13     }
14     fclose(stream);
15     return i;
16 }
17
18 int get_col(char *filename) //获取列数
19 {
20     char line[1024];
21     int i = 0;
22     FILE *stream = fopen(filename, "r");
23     fgets(line, 1024, stream);
24     char *token = strtok(line, ",");
25     while (token)
26     {
27         token = strtok(NULL, ",");
28         i++;
29     }
30     fclose(stream);
31     return i;
32 }
33
34 void get_two_dimension(char *line, double **data, char *filename)
35 {
36     FILE *stream = fopen(filename, "r");
37     int i = 0;
38     while (fgets(line, 1024, stream)) //逐行读取
39     {
40         int j = 0;
41         char *tok;
42         char *tmp = strdup(line);
43         for (tok = strtok(line, ","); tok && *tok; j++, tok = strtok(NULL, ",\n"))
44         {
45             data[i][j] = atof(tok); //转换成浮点数
46             //字符串拆分操作
47         }
48         i++;
49         free(tmp);
50     }
51     fclose(stream); //文件打开后要关闭操作
52 }
```

## 2) k\_fold.c

```
1 #include <stdlib.h>
2 #include <stdio.h>
3
4 double ***cross_validation_split(double **dataset, int row, int n_folds, int
fold_size, int col)
5 {
6     srand(10); //种子
7     double ***split;
8     int i, j = 0, k = 0;
9     int index;
10    double **fold;
11    split = (double ***)malloc(n_folds * sizeof(double **));
12    for (i = 0; i < n_folds; i++)
13    {
14        fold = (double **)malloc(fold_size * sizeof(double *));
15        while (j < fold_size)
16        {
17            fold[j] = (double *)malloc(col * sizeof(double));
18            index = rand() % row;
19            fold[j] = dataset[index];
20            for (k = index; k < row - 1; k++) //for循环删除这个数组中被rand取到的元素
21            {
22                dataset[k] = dataset[k + 1];
23            }
24            row--; //每次随机取出一个后总行数-1, 保证不会重复取某一行
25            j++;
26        }
27        j = 0; //清零j
28        split[i] = fold;
29    }
30    return split;
31 }
```

## 3) rmse.c

```
1 #include <stdlib.h>
2 #include <stdio.h>
3 #include <math.h>
4
5 double rmse_metric(double *actual, double *predicted, int fold_size)
6 {
7     double sum_err = 0.0;
8     int i;
9     for (i = 0; i < fold_size; i++)
10    {
11        double err = predicted[i] - actual[i];
12        sum_err += err * err;
13    }
14    double mean_err = sum_err / fold_size;
15    return sqrt(mean_err);
16 }
```

## 4) test\_prediction.c

```
1 #include <stdlib.h>
2 #include <stdio.h>
3
4 extern void coefficients(double **data, double *coef, int length);
```

```

5 double *get_test_prediction(int col, int row, double **train, double **test, int
n_folds)
6 {
7     double *coef = (double *)malloc(col * sizeof(double));
8     int i;
9     for (i = 0; i < col; i++)
10     {
11         coef[i] = 0.0;
12     }
13     int fold_size = (int)row / n_folds;
14     int train_size = fold_size * (n_folds - 1);
15     coefficients(train, coef, train_size);
16     double *predictions = (double *)malloc(fold_size * sizeof(double));
17     for (i = 0; i < fold_size; i++)
18     {
19         predictions[i] = coef[0] + coef[1] * test[i][0];
20     }
21     return predictions;
22 }

```

## 5) evaluate.c

```

1  #include <stdlib.h>
2  #include <stdio.h>
3
4  extern double *get_test_prediction(int col, int row, double **train, double **test,
int n_folds);
5  extern double rmse_metric(double *actual, double *predicted, int fold_size);
6  extern double ***cross_validation_split(double **dataset, int row, int col, int
n_folds, int fold_size);
7
8  double *evaluate_algorithm(double **dataset, int row, int col, int n_folds)
9  {
10     int fold_size = (int)row / n_folds;
11     double ***split = cross_validation_split(dataset, row, n_folds, fold_size, col);
12     int i, j, k, l;
13     int test_size = fold_size;
14     int train_size = fold_size * (n_folds - 1);
15     double *score = (double *)malloc(n_folds * sizeof(double));
16     for (i = 0; i < n_folds; i++)
17     {
18         //因为要遍历删除，所以拷贝一份split
19         double ***split_copy = (double ***)malloc(n_folds * sizeof(double **));
20         for (j = 0; j < n_folds; j++)
21         {
22             split_copy[j] = (double **)malloc(fold_size * sizeof(double *));
23             for (k = 0; k < fold_size; k++)
24             {
25                 split_copy[j][k] = (double *)malloc(col * sizeof(double));
26             }
27         }
28         for (j = 0; j < n_folds; j++)
29         {
30             for (k = 0; k < fold_size; k++)
31             {
32                 for (l = 0; l < col; l++)
33                 {
34                     split_copy[j][k][l] = split[j][k][l];
35                 }
36             }
37         }
38     }

```

```

39     double **test_set = (double **)malloc(test_size * sizeof(double *));
40     for (j = 0; j < test_size; j++)
41     {
42         test_set[j] = (double *)malloc(col * sizeof(double));
43         for (k = 0; k < col; k++)
44         {
45             test_set[j][k] = split_copy[i][j][k];
46         }
47     }
48     for (j = i; j < n_folds - 1; j++)
49     {
50         split_copy[j] = split_copy[j + 1];
51     } //删除取出来的fold
52
53     double **train_set = (double **)malloc(train_size * sizeof(double *));
54     for (k = 0; k < n_folds - 1; k++)
55     {
56         for (l = 0; l < fold_size; l++)
57         {
58             train_set[k * fold_size + l] = (double *)malloc(col *
sizeof(double));
59             train_set[k * fold_size + l] = split_copy[k][l];
60             //printf("split_copy[%d][%d]=%f\n", k,l,split_copy[k][l]);
61         }
62     }
63
64     double *predicted = (double *)malloc(test_size * sizeof(double));
65     predicted = get_test_prediction(col, row, train_set, test_set, n_folds);
66     double *actual = (double *)malloc(test_size * sizeof(double));
67     for (l = 0; l < test_size; l++)
68     {
69         actual[l] = test_set[l][col - 1];
70     }
71     double rmse = rmse_metric(actual, predicted, test_size);
72     score[i] = rmse;
73     printf("score[%d] = %lf\n", i, score[i]);
74     free(split_copy);
75 }
76 double total = 0;
77 for (l = 0; l < n_folds; l++)
78 {
79     total += score[l];
80 }
81 printf("mean_rmse = %lf\n", total / n_folds);
82 return score;
83 }

```

## 6) main.c

```

1  #include <stdlib.h>
2  #include <stdio.h>
3
4  double **dataset;
5  int row, col;
6
7  extern int get_row(char *filename);
8  extern int get_col(char *filename);
9  extern void get_two_dimension(char *line, double **dataset, char *filename);
10 extern double *evaluate_algorithm(double **dataset, int row, int col, int n_folds);
11
12 double mean(double *values, int length);
13 double covariance(double *x, double mean_x, double *y, double mean_y, int length);

```

```

14 double variance(double *values, double mean, int length);
15 void coefficients(double **data, double *coef, int length);
16
17 //计算均值方差等统计量 (多个函数)
18 double mean(double *values, int length)
19 { //对一维数组求均值
20     int i;
21     double sum = 0.0;
22     for (i = 0; i < length; i++)
23     {
24         sum += values[i];
25     }
26     double mean = (double)(sum / length);
27     return mean;
28 }
29
30 double covariance(double *x, double mean_x, double *y, double mean_y, int length)
31 {
32     double cov = 0.0;
33     int i = 0;
34     for (i = 0; i < length; i++)
35     {
36         cov += (x[i] - mean_x) * (y[i] - mean_y);
37     }
38     return cov;
39 }
40
41 double variance(double *values, double mean, int length)
42 { //这里求的是平方和, 没有除以n
43     double sum = 0.0;
44     int i;
45     for (i = 0; i < length; i++)
46     {
47         sum += (values[i] - mean) * (values[i] - mean);
48     }
49     return sum;
50 }
51
52 //由均值方差估计回归系数
53 void coefficients(double **data, double *coef, int length)
54 {
55     double *x = (double *)malloc(length * sizeof(double));
56     double *y = (double *)malloc(length * sizeof(double));
57     int i;
58     for (i = 0; i < length; i++)
59     {
60         x[i] = data[i][0];
61         y[i] = data[i][1];
62     }
63     double x_mean = mean(x, length);
64     double y_mean = mean(y, length);
65     coef[1] = covariance(x, x_mean, y, y_mean, length) / variance(x, x_mean, length);
66     coef[0] = y_mean - coef[1] * x_mean;
67 }
68
69 int main()
70 {
71     char filename[] = "Auto insurance.csv";
72     char line[1024];
73     row = get_row(filename);
74     col = get_col(filename);
75     dataset = (double **)malloc(row * sizeof(double *));
76     for (int i = 0; i < row; ++i)

```

```

77     {
78         dataset[i] = (double *)malloc(col * sizeof(double));
79     } //动态申请二维数组
80     get_two_dimension(line, dataset, filename);
81     int n_folds = 10;
82     int fold_size = (int)row / n_folds;
83     evaluate_algorithm(dataset, row, col, n_folds);
84     return 0;
85 }

```

## 7) compile.sh

```
1 gcc main.c read_csv.c k_fold.c evaluate.c rmse.c test_prediction.c -o run -lm && ./run
```

编译&运行:

```
1 bash compile.sh
```

最终输出结果如下:

```

1 score[0] = 33.263512
2 score[1] = 30.319399
3 score[2] = 22.835829
4 score[3] = 38.080193
5 score[4] = 23.662033
6 score[5] = 25.166845
7 score[6] = 47.085342
8 score[7] = 46.614182
9 score[8] = 40.007351
10 score[9] = 28.511198
11 mean_rmse = 33.554588

```

## Python语言实战

本节同样假设您已经下载数据集，我们使用著名机器学习开源库sklearn高效实现简单线性回归算法，以便您在实战中使用该算法：

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import KFold
4 from sklearn.linear_model import LinearRegression
5
6
7 def rmse_metric(actual, predicted):
8     sum_err = 0.0
9     for i in range(len(actual)):
10         err = predicted[i] - actual[i]
11         sum_err += err ** 2
12     mean_err = sum_err / (len(actual)-1)
13     return np.sqrt(mean_err)
14
15
16 if __name__ == '__main__':
17     dataset = np.array(pd.read_csv("insurance.csv", sep=',', header=None))
18     k_Cross = KFold(n_splits=10, random_state=1, shuffle=True)
19     index = 0
20     score = np.array([])
21     for train_index, test_index in k_Cross.split(dataset):
22         train_data, train_label = dataset[train_index, :-1], dataset[train_index,
-1]

```



```

23 test_data, test_label = dataset[test_index, :-1], dataset[test_index, -1]
24 model = LinearRegression()
25 model.fit(train_data, train_label)
26 pred = model.predict(test_data)
27 rmse = rmse_metric(test_label, pred)
28 score = np.append(score, rmse)
29 print('score[{}] = {}'.format(index, rmse))
30 index+=1
31 print('mean_rmse = {}'.format(np.mean(score)))

```

输出结果如下:

```

1 score[0] = 31.081476539821356
2 score[1] = 31.903964258437547
3 score[2] = 37.76453473731135
4 score[3] = 52.46285733249147
5 score[4] = 46.256226601172855
6 score[5] = 25.094234805956997
7 score[6] = 27.19738282646511
8 score[7] = 59.134038915742195
9 score[8] = 34.08824724550272
10 score[9] = 39.795062610664274
11 mean_rmse = 38.477802587356585

```

## 3.2 Multivariate Linear Regression

许多机器学习算法的核心是优化。优化算法通常被用来从给定的数据集中寻找一个好的模型参数。机器学习中最常用的优化算法是随机梯度下降法。在本教程中，您将了解如何使用C语言从头开始实现随机梯度下降，并以此来优化线性回归算法。

### 3.2.1 算法介绍

#### 多元线性回归

线性回归是一种预测真实值的算法。这些需要预测真实值的问题叫做回归问题。线性回归是一种使用直线来建立输入值和输出值之间关系模型的算法。在二维以上的空间中，这条直线被称为一个平面或超平面。

预测就是通过输入值的组合来预测输出值。每个输入属性(x)都使用一个系数(b)对其进行加权，学习算法的目标遍是寻找这组能产生良好预测(y)的系数。

$$y = b_0 + b_1 \times x_1 + b_2 \times x_1 + \dots \quad (1.1)$$

这组系数可以用随机梯度下降法求得。

#### 随机梯度下降

梯度下降是沿着函数的斜率或梯度最小化函数的过程。在机器学习中，我们可以使用一种算法来评估和更新每次迭代的效率，以最小化我们的模型在训练数据上的误差，这种算法被称为随机梯度下降法。

这种优化算法的工作原理是每次向模型提供一个训练样本，模型对这个训练样本进行预测，计算误差并更新模型以减少下一次预测的误差。此迭代过程将重复进行固定次数。

当我们得到某个模型的一个系数后，以下方程可以用来由已知系数计算新系数，从而使模型在训练数据上的误差最小。每次迭代，该方程都会使模型中的系数(b)更新。

$$b = b - learning\ rate \times error \times x \quad (1.2)$$

其中，b是被优化的系数或权值，学习率是你需要设定的参数(例如0.01)，error是模型由于权值问题导致在训练集上的预测误差，x是输入变量。

## 3.2.2 算法讲解

### 训练集梯度下降估计回归系数

我们可以使用随机梯度下降法来估计训练数据的系数。随机梯度下降法需要两个参数：

- **Learning Rate**:用于限制每次更新时每个系数的修正量。
- **Epochs**:更新系数时遍历训练数据的次数。

以上参数及训练数据将作为函数的输入参数。我们需要在函数中执行以下3个循环：

- 循环每个Epoch。
- 对于每个Epoch，循环训练数据中的每一行。
- 对于每个Epoch中的每一行，循环并更新每个系数。

我们可以看到，对于每个Epochs，我们更新了训练数据中每一行的系数。这种更新是根据模型的误差产生的。其中，误差是用候选系数的预测值与预期输出值的差来计算的。

$$error = prediction - expected \quad (1.3)$$

每个输入属性都有一个权重系数，并且以相同的方式更新，例如：

$$b_1(t+1) = b_1(t) - learning\ rate \times error(t) \times x_1(t) \quad (1.4)$$

迭代最初的系数，也称为截距或偏置，同样以类似的方式更新。最终结果与系数的初始输入值无关。

$$b_0(t+1) = b_0(t) - learning\ rate \times error(t) \quad (1.5)$$

下面是一个名为coefficients\_sgd()的函数，它使用随机梯度下降法计算训练数据集的系数。

- 功能——估计回归系数
- 以训练集数组、数组列数、系数存储数组、学习率、epoch、训练集长度作为输入参数。
- 最终输出系数存储数组。

```
1 double* coefficients_sgd(double** dataset,int col,double coef[], double l_rate, int
  n_epoch, int train_size) {
2     int i;
3     for (i = 0; i < n_epoch; i++) {
4         int j = 0; //遍历每一行
5         for (j = 0; j < train_size; j++) {
6             double yhat = predict(col,dataset[j], coef);
7             double err = yhat - dataset[j][col - 1];
8             coef[0] -= l_rate * err;
9             int k;
10            for (k = 0; k < col - 1; k++) {
11                coef[k + 1] -= l_rate * err*dataset[j][k];
12            }
13        }
14    }
15    for (i = 0; i < col; i++) {
16        printf("coef[i]=%f\n", coef[i]);
17    }
18    return coef;
19 }
```

我们利用如下数据集：

```

1  int main()
2  {
3      double data[5][5];
4      double* ptr[5];
5      double weight[5]={1,2,3,4,5};
6      for(int i=0;i<5;i++)
7      {
8          for(int j=0;j<5;j++)
9              data[i][j] = i+j;
10         ptr[i] = data[i];
11     }
12     coefficients_sgd(ptr, 5,weight, 0.1,100, 4);
13 }

```

得到结果如下:

```

1  weights[0] = 181.000
2  weights[1] = 322.000
3  weights[2] = 503.000
4  weights[3] = 684.000
5  weights[4] = 865.000

```

## 由回归系数计算预测值

在随机梯度下降中评估候选系数值时，在模型完成并且我们想要开始对测试数据或新数据进行预测时，我们都需要一个预测函数。

例如，有一个输入值(x)和两个系数值(b0和b1)。这个问题建模的预测方程为：

$$y = b_0 + b_1 \times x \quad (1.6)$$

以下是一个名为predict()的预测函数，给定系数后，它可以预测一组输入值(x)对应的输出值(y)。

- 功能——预测输出值(y)
  - 以输入值的属性个数、输入值、系数数组为输入参数
  - 最终输出预测值
- ```

1  double predict(int col,double array[], double coefficients[]) { //预测某一行的值
2      double yhat = coefficients[0];
3      int i;
4      for (i = 0; i < col - 1; i++)
5          yhat += coefficients[i + 1] * array[i];
6      return yhat;
7  }

```

我们利用如下数据：

```

1  int main()
2  {
3      double data[5]={0,1,2,3,4};
4      double weights[5]={1,2,3,4,5};
5      predict(5,data,weights);
6  }

```

得到如下结果：

```

1  1.0000

```

### 3.2.3 算法代码

我们现在知道了如何实现多元线性回归算法，那么我们把它应用到葡萄酒质量数据集 [winequize-white.csv](#)

我们给出链接：<https://aistudio.baidu.com/aistudio/datasetdetail/105756/0>

#### C语言细节讲解

本节假设您已下载数据集 `winequize-white.csv`，并且它在当前工作目录中可用。下面我们给出一个完整实例，使用C语言详细讲解每一处细节。我们给出每一个.c文件的所有代码：

##### 1) read\_csv.c

该文件代码与前面代码一致，不再重复给出。

##### 2) k\_fold.c

该文件代码与前面代码一致，不再重复给出。

##### 3) rmse.c

该文件代码与前面代码一致，不再重复给出。

##### 4) normalize.c

```
1 void normalize_dataset(double **dataset,int row, int col)
2 {
3     // 先 对列循环
4     double maximum, minimum;
5     for (int i = 0; i < col; i++)
6     {
7         // 第一行为标题，值为0，不能参与计算最大最小值
8         maximum = dataset[0][i];
9         minimum = dataset[0][i];
10        //再 对行循环
11        for (int j = 0; j < row; j++)
12        {
13            maximum = (dataset[j][i]>maximum)?dataset[j][i]:maximum;
14            minimum = (dataset[j][i]<minimum)?dataset[j][i]:minimum;
15        }
16        // 归一化处理
17        for (int j = 0; j < row; j++)
18        {
19            dataset[j][i] = (dataset[j][i] - minimum) / (maximum - minimum);
20        }
21    }
22 }
```

##### 5) test\_prediction.c

```
1 #include <stdlib.h>
2 #include <stdio.h>
3
4 extern double *coefficients_sgd(double **dataset, int col, double coef[], double
l_rate, int n_epoch, int train_size);
5 extern double predict(int col, double array[], double coefficients[]);
6
7 double *get_test_prediction(double **dataset, int row, int col, double **train,
double **test, double l_rate, int n_epoch, int n_folds)
8 {
9     double *coef = (double *)malloc(col * sizeof(double));
10    int i;
11    for (i = 0; i < col; i++)
12    {
```

```

13     coef[i] = 0.0;
14 }
15 int fold_size = (int)row / n_folds;
16 int train_size = fold_size * (n_folds - 1);
17 coefficients_sgd(train, col, coef, l_rate, n_epoch, train_size);
18 double *predictions = (double *)malloc(fold_size * sizeof(double));
19 for (i = 0; i < fold_size; i++)
20 {
21     predictions[i] = predict(col, test[i], coef);
22 }
23 return predictions;
24 }

```

## 6) evaluate.c

```

1  #include <stdlib.h>
2  #include <stdio.h>
3
4  extern double *get_test_prediction(double **dataset, int row, int col, double
  **train, double **test, double l_rate, int n_epoch, int n_folds);
5  extern double rmse_metric(double *actual, double *predicted, int fold_size);
6  extern double ***cross_validation_split(double **dataset, int row, int col, int
  n_folds, int fold_size);
7
8  double *evaluate_algorithm(double **dataset, int row, int col, int n_folds, int
  n_epoch, double l_rate)
9  {
10     int fold_size = (int)row / n_folds;
11     double ***split = cross_validation_split(dataset, row, n_folds, fold_size, col);
12     int i, j, k, l;
13     int test_size = fold_size;
14     int train_size = fold_size * (n_folds - 1);
15     double *score = (double *)malloc(n_folds * sizeof(double));
16
17     for (i = 0; i < n_folds; i++)
18     {
19         //因为要遍历删除，所以拷贝一份split
20         double ***split_copy = (double ***)malloc(n_folds * sizeof(double **));
21         for (j = 0; j < n_folds; j++)
22         {
23             split_copy[j] = (double **)malloc(fold_size * sizeof(double **));
24             for (k = 0; k < fold_size; k++)
25             {
26                 split_copy[j][k] = (double *)malloc(col * sizeof(double));
27             }
28         }
29         for (j = 0; j < n_folds; j++)
30         {
31             for (k = 0; k < fold_size; k++)
32             {
33                 for (l = 0; l < col; l++)
34                 {
35                     split_copy[j][k][l] = split[j][k][l];
36                 }
37             }
38         }
39         double **test_set = (double **)malloc(test_size * sizeof(double **));
40         for (j = 0; j < test_size; j++)
41         {
42             test_set[j] = (double *)malloc(col * sizeof(double));
43             for (k = 0; k < col; k++)
44             {

```

```

45         test_set[j][k] = split_copy[i][j][k];
46     }
47 }
48 for (j = i; j < n_folds - 1; j++)
49 {
50     split_copy[j] = split_copy[j + 1];
51 } //删除取出来的fold
52
53 double **train_set = (double **)malloc(train_size * sizeof(double *));
54 for (k = 0; k < n_folds - 1; k++)
55 {
56     for (l = 0; l < fold_size; l++)
57     {
58         train_set[k * fold_size + l] = (double *)malloc(col *
sizeof(double));
59         train_set[k * fold_size + l] = split_copy[k][l];
60     }
61 }
62 double *predicted = (double *)malloc(test_size * sizeof(double));
63 predicted = get_test_prediction(dataset, row, col, train_set, test_set,
l_rate, n_epoch, n_folds);
64 double *actual = (double *)malloc(test_size * sizeof(double));
65 for (l = 0; l < test_size; l++)
66 {
67     actual[l] = test_set[l][col - 1];
68 }
69 double rmse = rmse_metric(actual, predicted, test_size);
70 score[i] = rmse;
71 printf("score[%d] = %f\n", i, score[i]);
72 free(split_copy);
73 }
74 double total = 0;
75 for (l = 0; l < n_folds; l++)
76 {
77     total += score[l];
78 }
79 printf("mean_rmse = %f\n", total / n_folds);
80 return score;
81 }

```

## 7) main.c

```

1  #include <stdlib.h>
2  #include <stdio.h>
3
4  extern int get_row(char *filename);
5  extern int get_col(char *filename);
6  extern void get_two_dimension(char *line, double **dataset, char *filename);
7  extern double *evaluate_algorithm(double **dataset, int row, int col, int n_folds,
int n_epoch, double l_rate);
8  extern void normalize_dataset(double **dataset, int row, int col);
9
10 double *coefficients_sgd(double **dataset, int col, double coef[], double l_rate, int
n_epoch, int train_size)
11 {
12     int i;
13     for (i = 0; i < n_epoch; i++)
14     {
15         int j = 0;
16         for (j = 0; j < train_size; j++)
17         {
18             double yhat = predict(col, dataset[j], coef);

```

```

19         double err = yhat - dataset[j][col - 1];
20         coef[0] -= l_rate * err;
21         int k;
22         for (k = 0; k < col - 1; k++)
23         {
24             coef[k + 1] -= l_rate * err * dataset[j][k];
25         }
26     }
27 }
28 return coef;
29 }
30
31 double predict(int col, double array[], double coefficients[])
32 {
33     double yhat = coefficients[0];
34     int i;
35     for (i = 0; i < col - 1; i++)
36         yhat += coefficients[i + 1] * array[i];
37     return yhat;
38 }
39
40 int main()
41 {
42     char filename[] = "winequality-white.csv";
43     char line[1024];
44     double **dataset;
45     int row, col;
46     row = get_row(filename);
47     col = get_col(filename);
48     dataset = (double **)malloc(row * sizeof(double *));
49     for (int i = 0; i < row; ++i)
50     {
51         dataset[i] = (double *)malloc(col * sizeof(double));
52     }
53     get_two_dimension(line, dataset, filename);
54     normalize_dataset(dataset, row, col);
55
56     int n_folds = 10;
57     double l_rate = 0.001f;
58     int n_epoch = 50;
59     evaluate_algorithm(dataset, row, col, n_folds, n_epoch, l_rate);
60     return 0;
61 }

```

## 8) compile.sh

```

1 gcc main.c read_csv.c normalize.c k_fold.c evaluate.c rmse.c test_prediction.c -o run
  -lm && ./run

```

编译&运行:

```

1 bash compile.sh

```

最终输出结果如下:

```

1 score[0] = 0.221540
2 score[1] = 0.209277
3 score[2] = 0.221540
4 score[3] = 0.219608
5 score[4] = 0.219479
6 score[5] = 0.216744
7 score[6] = 0.205718
8 score[7] = 0.202798
9 score[8] = 0.214637
10 score[9] = 0.207231
11 mean_rmse = 0.213857

```

## Python语言实战

本节同样假设您已经下载数据集，我们使用著名机器学习开源库sklearn高效实现多元线性回归算法，以便您在实战中使用该算法：

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import KFold
4 from sklearn.linear_model import SGDRegressor
5 from sklearn.preprocessing import MinMaxScaler
6
7
8 def rmse_metric(actual, predicted):
9     sum_err = 0.0
10    for i in range(len(actual)):
11        err = predicted[i] - actual[i]
12        sum_err += err ** 2
13    mean_err = sum_err / (len(actual)-1)
14    return np.sqrt(mean_err)
15
16
17 if __name__ == '__main__':
18     dataset = np.array(pd.read_csv("winequality-white.csv", sep=',', header=None))
19     k_Cross = KFold(n_splits=10, random_state=0, shuffle=True)
20     index = 0
21     score = np.array([])
22     scaler = MinMaxScaler()
23     data,label = dataset[:, :-1], dataset[:, -1]
24     data = scaler.fit_transform(data)
25     for train_index, test_index in k_Cross.split(dataset):
26         train_data, train_label = data[train_index, :], label[train_index]
27         test_data, test_label = data[test_index, :], label[test_index]
28         model = SGDRegressor()
29         model.fit(train_data, train_label)
30         pred = model.predict(test_data)
31         rmse = rmse_metric(test_label, pred)
32         score = np.append(score, rmse)
33         print('score[{}] = {}'.format(index, rmse))
34         index+=1
35     print('mean_rmse = {}'.format(np.mean(score)))

```

输出结果如下，读者可以尝试分析一下为何结果会存在差异？



```

1 score[0] = 0.8419014234018158
2 score[1] = 0.8408919161041173
3 score[2] = 0.7311825499558641
4 score[3] = 0.8147707681816518
5 score[4] = 0.7276314042865725
6 score[5] = 0.7403970929333936
7 score[6] = 0.7865008610855795
8 score[7] = 0.822294388008359
9 score[8] = 0.7899616477361368
10 score[9] = 0.7447500726966548
11 mean_rmse = 0.7840282124390145

```

## 3.3 Logistic Regression

逻辑回归是一种针对两类问题的归一化线性分类算法。它易于实现，易于理解，并且在各种各样的问题上都能得到很好的结果。在本教程中，您将了解如何借助C语言，使用随机梯度下降从零开始实现逻辑回归

### 3.3.1 算法介绍

#### 逻辑回归

逻辑回归是根据该方法的核心函数——Logistic函数而命名的。逻辑回归可以使用一个方程来表示，与线性回归非常相似。输入值(x)加以权重值或系数值，然后通过它们的线性组合来预测输出值(y)。逻辑与线性回归的一个关键区别是，建模的输出值是二进制值(0或1)，而不是数值。

$$\hat{y} = \frac{e^{b_0 + b_1 \times x_1}}{1 + e^{b_0 + b_2 \times x_1}} \quad (1.1)$$

可以简化为：

$$\hat{y} = \frac{1.0}{1.0 + e^{-(b_0 + b_2 \times x_1)}} \quad (1.2)$$

其中e为自然对数(欧拉数)的底数，yhat为预测输出，b0为偏置或截距项，b1为单个输入值(x1)的系数。预测值ythat是一个0到1之间的实值，需要四舍五入到一个整数值并映射到一个预测的类值。

输入数据中的每一列都有一个相关的b系数(一个常实值)，它必须从训练数据中获得。您存储在内存或文件中的模型实际是方程中的系数(beta值或b)。Logistic回归算法的系数由训练数据估计而来。

#### 随机梯度下降

Logistic回归使用梯度下降更新系数。梯度下降在8.1.2节中进行了介绍和描述。每次梯度下降迭代时，机器学习语言中的系数(b)根据以下公式更新：

$$b = b + \text{learning rate} \times (y - \hat{y}) \times \hat{y} \times (1 - \hat{y}) \times x \quad (1.3)$$

其中，b是将被优化的系数或权重，learning rate一个学习速率，它需要您的配置(例如0.01)，(y - yhat)是模型在被分配有权重的训练集上的误差，yhat是由预测系数得到的预测值，x是输入值。

### 3.3.2 算法讲解

#### 训练集梯度下降估计回归系数

我们可以使用随机梯度下降估计训练数据的系数值。随机梯度下降需要两个参数：

- **Learning Rate:**用于限制每次更新时每个系数的修正量。
- **Epochs:**更新系数时，训练数据运行的次数。

以上参数及训练数据将作为函数的参数。我们将在函数中执行3个循环：

- 每个Epoch循环一次
- 对于每一个Epoch，循环遍历训练集中的每一行
- 对于每一个Epoch中的每一行，循环遍历每个权重并更新它

如你所见，每个Epoch，我们根据模型产生的误差更新了训练数据中每一行的系数。误差计算为期望输出值与用候选系数得到的预测值之间的差。每个输入属性有一个权重系数，并且以一致的方式更新，例如：

$$b1(t+1) = b1(t) + learning\ rate \times (y(t) - \hat{y}(t)) \times \hat{y}(t) \times (1 - \hat{y}(t)) \times x1(t) \quad (1.4)$$

列表开始处的特殊系数，也称为截距，以类似的方式更新，只是没有输入，因为它与特定的输入值没有关联：

$$b0(t+1) = b0(t) + learning\ rate \times (y(t) - \hat{y}(t)) \times \hat{y}(t) \times (1 - \hat{y}(t)) \quad (1.4)$$

现在我们可以把这些放在一起。下面是一个名为coefficients\_sgd()的函数，它使用随机梯度下降法计算训练数据集的系数值。

```
1 void coefficients_sgd(double ** dataset, int col, double *coef, double l_rate, int
  n_epoch, int train_size) {
2     int i;
3     for (i = 0; i < n_epoch; i++) {
4         int j = 0; //遍历每一行
5         for (j = 0; j < train_size; j++) {
6             double yhat = predict(col, dataset[j], coef);
7             double err = dataset[j][col - 1] - yhat;
8             coef[0] += l_rate * err * yhat * (1 - yhat);
9             int k;
10            for (k = 0; k < col - 1; k++) {
11                coef[k + 1] += l_rate * err * yhat * (1 - yhat) * dataset[j][k];
12            }
13        }
14    }
15    for (i = 0; i < col; i++) {
16        printf("coef[%d]=%f\n", i, coef[i]);
17    }
18 }
```

我们利用如下数据集：

```
1 int main()
2 {
3     double data[5][5];
4     double* ptr[5];
5     double weight[5]={1,2,3,4,5};
6     for(int i=0;i<5;i++)
7     {
8         for(int j=0;j<5;j++)
9             data[i][j] = i+j;
10        ptr[i] = data[i];
11    }
12    coefficients_sgd(ptr, 5, weight, 0.1, 100, 4);
13 }
```

得到结果如下：

```
1 weights[0] = 181.000
2 weights[1] = 322.000
3 weights[2] = 503.000
4 weights[3] = 684.000
5 weights[4] = 865.000
```

## 由回归系数计算预测值

在随机梯度下降中评估候选系数值时，以及在模型完成并且我们希望开始对测试数据或新数据进行预测后，都需要一个可以预测的函数。下面是一个名为predict()的函数，给定一组系数后，它可以输出一行预测值。第一个系数是截距，也称为偏差或b0，因为它是独立的，不对应特定的输入值。

```
1 double predict(int col, double array[], double coefficients[]) { //预测某一行的值
2     double yhat = coefficients[0];
3     int i;
4     for (i = 0; i < col - 1; i++)
5         yhat += coefficients[i + 1] * array[i];
6     printf("%f", 1 / (1 + exp(-yhat)));
7     return 1 / (1 + exp(-yhat));
8 }
```

我们使用如下数据：

```
1 int main()
2 {
3     double data[5]={0,1,2,3,4};
4     double weights[5]={1,2,3,4,5};
5     predict(5,data,weights);
6 }
```

得到如下结果：

```
1 1.0000
```

### 3.3.3 算法代码

我们现在知道了如何实现一个**逻辑回归模型**，那么我们把它应用到[糖尿病数据集 Pima.csv](#)

我们给出链接：<https://aistudio.baidu.com/aistudio/datasetdetail/105756/0>

### C语言细节讲解

本节假设您已下载数据集 `Pima.csv`，并且它在当前工作目录中可用。下面我们给出一个完整实例，使用C语言详细讲解每一处细节。我们给出每一个.c文件的所有代码：

#### 1) read\_csv.c

该文件代码与前面代码一致，不再重复给出。

#### 2) k\_fold.c

该文件代码与前面代码一致，不再重复给出。

#### 3) score.c

```
1 #include <stdlib.h>
2 #include <stdio.h>
3 #include <math.h>
4
5 double accuracy_metric(double *actual, double *predicted, int fold_size)
6 {
7     int correct = 0;
8     int i;
9     for (i = 0; i < fold_size; i++)
10     {
11         if (actual[i] == predicted[i])
12             correct += 1;
13     }
14 }
```

```

13     }
14     return (correct / (double)fold_size) * 100.0;
15 }

```

#### 4) normalize.c

该文件代码与前面代码一致，不再重复给出。

#### 5) test\_prediction.c

```

1  #include <stdlib.h>
2  #include <stdio.h>
3
4  extern void coefficients_sgd(double **dataset, int col, double *coef, double l_rate,
5  int n_epoch, int train_size);
6  extern double predict(int col, double array[], double coefficients[]);
7
8  double *get_test_prediction(double **dataset, int row, int col, double **train,
9  double **test, double l_rate, int n_epoch, int n_folds)
10 {
11     double *coef = (double *)malloc(col * sizeof(double));
12     int i;
13     for (i = 0; i < col; i++)
14     {
15         coef[i] = 0.0;
16     }
17     int fold_size = (int)row / n_folds;
18     int train_size = fold_size * (n_folds - 1);
19     coefficients_sgd(train, col, coef, l_rate, n_epoch, train_size);
20     double *predictions = (double *)malloc(fold_size * sizeof(double));
21     for (i = 0; i < fold_size; i++)
22     {
23         predictions[i] = predict(col, test[i], coef);
24         predictions[i] = (float)(int)(predictions[i] + 0.5);
25     }
26     return predictions;
27 }

```

#### 6) evaluate.c

```

1  #include <stdlib.h>
2  #include <stdio.h>
3
4  extern double *get_test_prediction(double **dataset, int row, int col, double
5  **train, double **test, double l_rate, int n_epoch, int n_folds);
6  extern double accuracy_metric(double *actual, double *predicted, int fold_size);
7  extern double ***cross_validation_split(double **dataset, int row, int col, int
8  n_folds, int fold_size);
9
10 double *evaluate_algorithm(double **dataset, int row, int col, int n_folds, int
11 n_epoch, double l_rate)
12 {
13     int fold_size = (int)row / n_folds;
14     double ***split = cross_validation_split(dataset, row, n_folds, fold_size, col);
15     int i, j, k, l;
16     int test_size = fold_size;
17     int train_size = fold_size * (n_folds - 1);
18     double *score = (double *)malloc(n_folds * sizeof(double));
19
20     for (i = 0; i < n_folds; i++)
21     {
22         // 因为要遍历删除，所以拷贝一份split

```

```

20     double ***split_copy = (double ***)malloc(n_folds * sizeof(double **));
21     for (j = 0; j < n_folds; j++)
22     {
23         split_copy[j] = (double **)malloc(fold_size * sizeof(double *));
24         for (k = 0; k < fold_size; k++)
25         {
26             split_copy[j][k] = (double *)malloc(col * sizeof(double));
27         }
28     }
29     for (j = 0; j < n_folds; j++)
30     {
31         for (k = 0; k < fold_size; k++)
32         {
33             for (l = 0; l < col; l++)
34             {
35                 split_copy[j][k][l] = split[j][k][l];
36             }
37         }
38     }
39     double **test_set = (double **)malloc(test_size * sizeof(double *));
40     for (j = 0; j < test_size; j++)
41     {
42         test_set[j] = (double *)malloc(col * sizeof(double));
43         for (k = 0; k < col; k++)
44         {
45             test_set[j][k] = split_copy[i][j][k];
46         }
47     }
48     for (j = i; j < n_folds - 1; j++)
49     {
50         split_copy[j] = split_copy[j + 1];
51     }
52     // 删除取出来的fold
53     double **train_set = (double **)malloc(train_size * sizeof(double *));
54     for (k = 0; k < n_folds - 1; k++)
55     {
56         for (l = 0; l < fold_size; l++)
57         {
58             train_set[k * fold_size + l] = (double *)malloc(col *
sizeof(double));
59             train_set[k * fold_size + l] = split_copy[k][l];
60         }
61     }
62     double *predicted = (double *)malloc(test_size * sizeof(double));
63     predicted = get_test_prediction(dataset, row, col, train_set, test_set,
l_rate, n_epoch, n_folds);
64     double *actual = (double *)malloc(test_size * sizeof(double));
65     for (l = 0; l < test_size; l++)
66     {
67         actual[l] = test_set[l][col - 1];
68     }
69     double accuracy = accuracy_metric(actual, predicted, test_size);
70     score[i] = accuracy;
71     printf("score[%d]=%f%%\n", i, score[i]);
72     free(split_copy);
73 }
74 double total = 0;
75 for (l = 0; l < n_folds; l++)
76 {
77     total += score[l];
78 }
79 printf("mean_accuracy=%f%%\n", total / n_folds);
80 return score;

```

```
81 | }
```

## 7) main.c

```
1  #include <stdlib.h>
2  #include <stdio.h>
3  #include <math.h>
4
5  double **dataset;
6  int row, col;
7
8  extern int get_row(char *filename);
9  extern int get_col(char *filename);
10 extern void get_two_dimension(char *line, double **dataset, char *filename);
11 extern double *evaluate_algorithm(double **dataset, int row, int col, int n_folds,
12     int n_epoch, double l_rate);
13 extern void normalize_dataset(double **dataset, int row, int col);
14
15 void coefficients_sgd(double **dataset, int col, double *coef, double l_rate, int
16     n_epoch, int train_size)
17 {
18     int i;
19     for (i = 0; i < n_epoch; i++)
20     {
21         int j = 0;
22         for (j = 0; j < train_size; j++)
23         {
24             double yhat = predict(col, dataset[j], coef);
25             double err = dataset[j][col - 1] - yhat;
26             coef[0] += l_rate * err * yhat * (1 - yhat);
27             int k;
28             for (k = 0; k < col - 1; k++)
29             {
30                 coef[k + 1] += l_rate * err * yhat * (1 - yhat) * dataset[j][k];
31             }
32         }
33     }
34 }
35
36 double predict(int col, double array[], double coefficients[])
37 {
38     double yhat = coefficients[0];
39     int i;
40     for (i = 0; i < col - 1; i++)
41         yhat += coefficients[i + 1] * array[i];
42     return 1 / (1 + exp(-yhat));
43 }
44
45 int main()
46 {
47     char filename[] = "Pima.csv";
48     char line[1024];
49     row = get_row(filename);
50     col = get_col(filename);
51     dataset = (double **)malloc(row * sizeof(double *));
52     for (int i = 0; i < row; ++i)
53     {
54         dataset[i] = (double *)malloc(col * sizeof(double));
55     }
56     get_two_dimension(line, dataset, filename);
57     normalize_dataset(dataset, row, col);
58     int n_folds = 5;
```

```

57     double l_rate = 0.1f;
58     int n_epoch = 100;
59     evaluate_algorithm(dataset, row, col, n_folds, n_epoch, l_rate);
60     return 0;
61 }

```

## 8) compile.sh

```

1 gcc main.c read_csv.c normalize.c k_fold.c evaluate.c score.c test_prediction.c -o run
  -lm && ./run

```

编译&运行:

```

1 bash compile.sh

```

最终输出结果如下:

```

1 score[0] = 78.431373%
2 score[1] = 79.738562%
3 score[2] = 72.549020%
4 score[3] = 75.163399%
5 score[4] = 77.124183%
6 mean_accuracy = 76.601307%

```

## Python语言实战

本节同样假设您已经下载数据集，我们使用著名机器学习开源库sklearn高效实现**逻辑回归算法**，以便您在实战中使用该算法：

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import KFold
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.metrics import accuracy_score
6 from sklearn.preprocessing import MinMaxScaler
7
8
9 if __name__ == '__main__':
10     dataset = np.array(pd.read_csv("Pima.csv", sep=',', header=None))
11     k_Cross = KFold(n_splits=5, random_state=0, shuffle=True)
12     index = 0
13     score = np.array([])
14     Scaler = MinMaxScaler()
15     data, label = dataset[:, :-1], dataset[:, -1]
16     data = Scaler.fit_transform(data)
17     for train_index, test_index in k_Cross.split(dataset):
18         train_data, train_label = data[train_index, :], label[train_index]
19         test_data, test_label = data[test_index, :], label[test_index]
20         model = LogisticRegression()
21         model.fit(train_data, train_label)
22         pred = model.predict(test_data)
23         acc = accuracy_score(test_label, pred)
24         score = np.append(score, acc)
25         print('score[{}] = {}'.format(index, acc))
26         index += 1
27     print('mean_accuracy = {}'.format(np.mean(score)))

```

输出结果如下:

```

1 score[0] = 0.8181818181818182%
2 score[1] = 0.7532467532467533%
3 score[2] = 0.7467532467532467%
4 score[3] = 0.7843137254901961%
5 score[4] = 0.7581699346405228%
6 mean_accuracy = 0.7721330956625074%

```

## 3.4 Perceptron

感知器算法是一种最简单的人工神经网络。它是一个单神经元模型，可用于两类分类问题，并为以后开发更大的网络提供了基础。在本教程中，您将了解如何使用C语言从头开始实现感知器算法。

### 3.4.1 算法介绍

#### 感知算法

感知器的灵感来自于单个神经细胞的信息处理过程，这种神经细胞被称为神经元。神经元通过树突接收输入信号，然后树突将电信号传递到细胞体。以类似的方式，感知器从训练数据的例子中接收输入信号，我们将其加权并组合成一个线性方程，这个方程被称为激活函数。

$$activation = bias + \sum_{i=1}^n weight_i \times x_i \quad (1.1)$$

然后利用传递函数(如阶跃传递函数)将激活函数转化为输出值或预测值。

$$prediction = 1.0 IF activation \geq 0.0 ELSE 0.0 \quad (1.2)$$

这样，感知器就是一个解决两类问题（0和1）的分类算法，一个线性方程(如线或超平面)可以用来分离这两个类。它与线性回归和逻辑回归密切相关，后者以类似的方式进行预测(例如输入加权和)。感知器算法的权值必须使用随机梯度下降从训练数据中估计出来。

#### 随机梯度下降

感知器算法使用梯度下降来更新权重。梯度下降在多元线性回归一节中进行了介绍和描述。每次梯度下降迭代，权值w根据公式更新如下：

$$w = w + learning\ rate \times (expected - predicted) \times x \quad (1.3)$$

其中，w是被优化的权重，learning rate是一个需要你配置的学习率（例如0.01），

(如0.01)，(expected - predicted)为模型对训练数据归属权值的预测误差，x为输入值。(expected - predicted)为模型对带有权值的训练数据的预测误差，x为输入值。

### 3.4.2 算法讲解

#### 训练集梯度下降估计回归系数

我们可以使用随机梯度下降来估计训练数据的权值。随机梯度下降需要两个参数：

- **Learning Rate**: 用于限制每个重量的数量，每次更新时修正。
- **Epochs**: 更新权重时，遍历训练集的次数。

以上参数及训练数据将作为函数的参数。我们将在函数中执行3个循环：

- 每个Epoch循环一次
- 对于每一个Epoch，循环遍历训练集中的每一行
- 对于每一个Epoch中的每一行，循环遍历每个权重并更新它

1 如上，对于每一个Epoch，我们都更新了训练数据中每一行的每个权值。这个权值是根据模型产生的误差进行更新的。误差即期望输出值与由候选权值得到的预测值之间的差。

每个输入属性都有一个权重，并且以相同的方式更新这些权重。例如：

$$w_0(t+1) = w_0(t) + Learning\ rate \times (expected - predicted) \times 1 \quad w_1(t+1) = w_1(t) + Learning\ rate \times (expected - predicted) \times x_1 \quad (1.4)$$



$$w(t+1) = w(t) + learning\ rate \times (expected(t)) - predicted(t) \times x(t) \quad (1.4)$$

除了没有输入外，偏差也以类似的方式进行更新，因为它与特定的输入值没有关联：

$$bias(t+1) = bias(t) + learning\ rate \times (expected(t)) - predicted(t) \quad (1.5)$$

现在我们可以把这些放在一起。下面是一个名为train\_weights()的函数，它使用的是随机梯度下降法来计算训练数据集的权值。

```
1 void train_weights(double **data, int col, double *weights, double l_rate, int
  n_epoch, int train_size) {
2     int i;
3     for (i = 0; i < n_epoch; i++) {
4         int j = 0; //遍历每一行
5         for (j = 0; j < train_size; j++) {
6             double yhat = predict(col, data[j], weights);
7             double err = data[j][col - 1] - yhat;
8             weights[0] += l_rate * err;
9             int k;
10            for (k = 0; k < col - 1; k++) {
11                weights[k + 1] += l_rate * err * data[j][k];
12            }
13        }
14    }
15    for (i = 0; i < col; i++) {
16        printf("weights[%d]=%f\n", i, weights[i]);
17    }
18 }
```

我们利用如下数据：

```
1 int main()
2 {
3     double data[5][5];
4     double* ptr[5];
5     double weights[5]={1,2,3,4,5};
6     for(int i=0;i<5;i++)
7     {
8         for(int j=0;j<5;j++)
9             data[i][j] = i+j;
10        ptr[i] = data[i];
11    }
12    train_weights(ptr, 5, weights, 0.1, 100, 4);
13 }
```

得到如下结果：

```
1 weights[0] = 181.000
2 weights[1] = 322.000
3 weights[2] = 503.000
4 weights[3] = 684.000
5 weights[4] = 865.000
```

## 由回归系数计算预测值

开发一个可以进行预测的函数。这在随机梯度下降中评估候选权值时都需要用到，在模型完成后，我们希望对测试数据或新数据进行预测。下面是一个名为predict()的函数，它预测给定一组权重的行的输出值。第一个权重总是偏差，因为它是独立的，不对应任何特定的输入值。

```

1 double predict(int col,double *array, double *weights) {
2     //预测某一行的值
3     double activation = weights[0];
4     int i;
5     for (i = 0; i < col - 1; i++)
6         activation += weights[i + 1] * array[i];
7     double output = 0.0;
8     if (activation >= 0.0)
9         output = 1.0;
10    else
11        output = 0.0;
12    return output;
13 }

```

我们利用如下数据：

```

1 int main()
2 {
3     double data[5]={0,1,2,3,4};
4     double weights[5]={1,2,3,4,5};
5     predict(5,data,weights);
6 }

```

得到如下结果：

```

1 | 1.0000

```

### 3.4.3 算法代码

我们现在知道了如何实现**感知器算法**，那么我们把它应用到[声纳数据集 sonar.csv](#)

我们给出链接：<https://aistudio.baidu.com/aistudio/datasetdetail/105756/0>

### C语言细节讲解

本节假设您已下载数据集 `sonar.csv`，并且它在当前工作目录中可用。下面我们给出一个完整实例，使用C语言详细讲解每一处细节。我们给出每一个.c文件的所有代码：

#### 1) read\_csv.c

该文件代码与前面代码一致，不再重复给出。

#### 2) k\_fold.c

该文件代码与前面代码一致，不再重复给出。

#### 3) score.c

该文件代码与前面代码一致，不再重复给出。

#### 4) normalize.c

该文件代码与前面代码一致，不再重复给出。

#### 5) test\_prediction.c

```

1 #include <stdlib.h>
2 #include <stdio.h>
3
4 extern void train_weights(double **data, int col, double *weights, double l_rate, int
n_epoch, int train_size);
5 extern double predict(int col, double *array, double *weights);

```

```

6
7 double *get_test_prediction(double **train, double **test, int row, int col, double
  l_rate, int n_epoch, int n_folds)
8 {
9     double *weights = (double *)malloc(col * sizeof(double));
10    int i;
11    for (i = 0; i < col; i++)
12    {
13        weights[i] = 0.0;
14    }
15    int fold_size = (int)row / n_folds;
16    int train_size = fold_size * (n_folds - 1);
17    train_weights(train, col, weights, l_rate, n_epoch, train_size);
18    double *predictions = (double *)malloc(fold_size * sizeof(double));
19    for (i = 0; i < fold_size; i++)
20    {
21        predictions[i] = predict(col, test[i], weights);
22    }
23    return predictions;
24 }

```

## 6) evaluate.c

```

1  #include <stdlib.h>
2  #include <stdio.h>
3
4  extern double *get_test_prediction(double **train, double **test, int row, int col,
  double l_rate, int n_epoch, int n_folds);
5  extern double accuracy_metric(double *actual, double *predicted, int fold_size);
6  extern double ***cross_validation_split(double **dataset, int row, int col, int
  n_folds, int fold_size);
7
8  double *evaluate_algorithm(double **dataset, int row, int col, int n_folds, int
  n_epoch, double l_rate)
9  {
10     int fold_size = (int)row / n_folds;
11     double ***split = cross_validation_split(dataset, row, n_folds, fold_size, col);
12     int i, j, k, l;
13     int test_size = fold_size;
14     int train_size = fold_size * (n_folds - 1);
15     double *score = (double *)malloc(n_folds * sizeof(double));
16
17     for (i = 0; i < n_folds; i++)
18     {
19         double ***split_copy = (double ***)malloc(n_folds * sizeof(double **));
20         for (j = 0; j < n_folds; j++)
21         {
22             split_copy[j] = (double **)malloc(fold_size * sizeof(double *));
23             for (k = 0; k < fold_size; k++)
24             {
25                 split_copy[j][k] = (double *)malloc(col * sizeof(double));
26             }
27         }
28         for (j = 0; j < n_folds; j++)
29         {
30             for (k = 0; k < fold_size; k++)
31             {
32                 for (l = 0; l < col; l++)
33                 {
34                     split_copy[j][k][l] = split[j][k][l];
35                 }
36             }

```

```

37     }
38     double **test_set = (double **)malloc(test_size * sizeof(double *));
39     for (j = 0; j < test_size; j++)
40     {
41         test_set[j] = (double *)malloc(col * sizeof(double));
42         for (k = 0; k < col; k++)
43         {
44             test_set[j][k] = split_copy[i][j][k];
45         }
46     }
47     for (j = i; j < n_folds - 1; j++)
48     {
49         split_copy[j] = split_copy[j + 1];
50     }
51     double **train_set = (double **)malloc(train_size * sizeof(double *));
52     for (k = 0; k < n_folds - 1; k++)
53     {
54         for (l = 0; l < fold_size; l++)
55         {
56             train_set[k * fold_size + l] = (double *)malloc(col *
sizeof(double));
57             train_set[k * fold_size + l] = split_copy[k][l];
58         }
59     }
60     double *predicted = (double *)malloc(test_size * sizeof(double));
61     predicted = get_test_prediction(train_set, test_set, row, col, l_rate,
n_epoch, n_folds);
62     double *actual = (double *)malloc(test_size * sizeof(double));
63     for (l = 0; l < test_size; l++)
64     {
65         actual[l] = test_set[l][col - 1];
66     }
67     double accuracy = accuracy_metric(actual, predicted, test_size);
68     score[i] = accuracy;
69     printf("score[%d] = %f%%\n", i, score[i]);
70     free(split_copy);
71 }
72 double total = 0;
73 for (l = 0; l < n_folds; l++)
74 {
75     total += score[l];
76 }
77 printf("mean_accuracy = %f%%\n", total / n_folds);
78 return score;
79 }

```

## 7) main.c

```

1  #include <stdlib.h>
2  #include <stdio.h>
3
4  extern int get_row(char *filename);
5  extern int get_col(char *filename);
6  extern void get_two_dimension(char *line, double **dataset, char *filename);
7  extern double *evaluate_algorithm(double **dataset, int row, int col, int n_folds,
int n_epoch, double l_rate);
8  extern void normalize_dataset(double **dataset, int row, int col);
9
10 double predict(int col, double *array, double *weights)
11 {
12     double activation = weights[0];
13     int i;

```

```

14     for (i = 0; i < col - 1; i++)
15         activation += weights[i + 1] * array[i];
16     double output = 0.0;
17     if (activation >= 0.0)
18         output = 1.0;
19     else
20         output = 0.0;
21     return output;
22 }
23
24 void train_weights(double **data, int col, double *weights, double l_rate, int
n_epoch, int train_size)
25 {
26     int i;
27     for (i = 0; i < n_epoch; i++)
28     {
29         int j = 0;
30         for (j = 0; j < train_size; j++)
31         {
32             double yhat = predict(col, data[j], weights);
33             double err = data[j][col - 1] - yhat;
34             weights[0] += l_rate * err;
35             int k;
36             for (k = 0; k < col - 1; k++)
37             {
38                 weights[k + 1] += l_rate * err * data[j][k];
39             }
40         }
41     }
42 }
43
44 int main()
45 {
46     double **dataset;
47     int row, col;
48     char filename[] = "sonar.csv";
49     char line[1024];
50     row = get_row(filename);
51     col = get_col(filename);
52     dataset = (double **)malloc(row * sizeof(double *));
53     for (int i = 0; i < row; ++i)
54     {
55         dataset[i] = (double *)malloc(col * sizeof(double));
56     }
57     get_two_dimension(line, dataset, filename);
58     normalize_dataset(dataset, row, col);
59     int n_folds = 3;
60     double l_rate = 0.01f;
61     int n_epoch = 500;
62     evaluate_algorithm(dataset, row, col, n_folds, n_epoch, l_rate);
63     return 0;
64 }

```

## 8) compile.sh

```

1 gcc main.c read_csv.c normalize.c k_fold.c evaluate.c score.c test_prediction.c -o run
-lm && ./run

```

### 编译&运行:

```

1 bash compile.sh

```

最终输出结果如下：

```
1 score[0] = 82.608696%
2 score[1] = 79.710145%
3 score[2] = 73.913043%
4 mean_accuracy = 78.743961%
```

## Python语言实战

本节同样假设您已经下载数据集，我们使用著名机器学习开源库sklearn高效实现感知机算法，以便您在实战中使用该算法：

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import KFold
4 from sklearn.linear_model import Perceptron
5 from sklearn.metrics import accuracy_score
6 from sklearn.preprocessing import MinMaxScaler
7
8
9 if __name__ == '__main__':
10     dataset = np.array(pd.read_csv("sonar.csv", sep=',', header=None))
11     k_Cross = KFold(n_splits=3, random_state=8, shuffle=True)
12     index = 0
13     score = np.array([])
14     scaler = MinMaxScaler()
15     data, label = dataset[:, :-1], dataset[:, -1]
16     data = scaler.fit_transform(data)
17     for train_index, test_index in k_Cross.split(dataset):
18         train_data, train_label = data[train_index, :], label[train_index]
19         test_data, test_label = data[test_index, :], label[test_index]
20         model = Perceptron(eta0=0.01, max_iter=500)
21         model.fit(train_data, train_label)
22         pred = model.predict(test_data)
23         acc = accuracy_score(test_label, pred)
24         score = np.append(score, acc)
25         print('score[{}] = {}'.format(index, acc))
26         index += 1
27     print('mean_accuracy = {}'.format(np.mean(score)))
```

输出结果如下：

```
1 score[0] = 0.7571428571428571%
2 score[1] = 0.8405797101449275%
3 score[2] = 0.6956521739130435%
4 mean_accuracy = 0.7644582470669427%
```

## 3.5 Classification and Regression Trees

决策树是一种广受欢迎的、强大的预测方法。它之所以受到欢迎，是因为其最终的模型对于从业人员来说易于理解，给出的决策树可以确切解释为何做出特定的预测。决策树是最简单的机器学习算法，它易于实现，可解释性强，完全符合人类的直观思维，有着广泛的应用。

同时，决策树也是更为高级的集成算法（如bagging，random forests和gradient boosting等）的基础。在本节中，您将了解Gini指数的概念、如何创建数据集的拆分、如何构建一棵树、如何利用构建的树作出分类决策以及如何在Banknote数据集上应用这些知识。

## 3.5.1 算法介绍

Classification and Regression Trees（简称CART），指的是可用于分类或回归预测建模问题的决策树算法。在本节中，我们将重点介绍如何使用CART解决分类问题，并以Banknote数据集为例进行演示。

CART模型的表示形式是一棵二叉树。每个节点表示单个输入变量（X）和该变量的分割点（假定变量是数字化的）。树的叶节点（也称作终端节点）包含用于预测的输出变量（y）。

创建二元决策树实际上是划分输入空间的过程。一般采用贪婪方法对变量进行递归的二进制拆分，使用某个成本函数（通常是Gini指数）测试不同的分割点，选择成本最高的拆分（即拆分完之后，剩余成本降到最低，亦代表这种拆分所含的“信息量”最大）。

## 3.5.2 算法讲解

### 按属性分割数据

- 功能：切分函数，根据切分点将数据分为左右两组
- 输出：从切分点处切分后的数据结果

```
1 struct dataset *test_split(int index, double value, int row, int col, double **data)
2 {
3     // 将切分结果作为结构体返回
4     struct dataset *split = (struct dataset *)malloc(sizeof(struct dataset));
5     int count1=0,count2=0;
6     double ***groups = (double ***)malloc(2 * sizeof(double **));
7     for (int i = 0; i < 2; i++)
8     {
9         groups[i]=(double **)malloc(row * sizeof(double *));
10        for (int j = 0; j < row; j++)
11        {
12            groups[i][j] = (double *)malloc(col * sizeof(double ));
13        }
14    }
15    for (int i = 0; i < row; i++)
16    {
17        if (data[i][index]<value)
18        {
19            groups[0][count1]=data[i];
20            count1 ++;
21        }else{
22            groups[1][count2] = data[i];
23            count2++;
24        }
25    }
26    split->splitdata = groups;
27    split->row1 = count1;
28    split->row2 = count2;
29    return split;
30 }
```

### Gini指数

基尼指数是用于评估数据集中的拆分所常用的成本函数。数据集中的拆分涉及一个输入属性和该属性的一个值。它可以用于将训练模式分为两组。最理想的拆分是使基尼指数变为0，而最坏的情况是在二分类问题中分为每一类的概率都是50%（即基尼指数变为0.5）。

基尼系数的具体计算公式如下：

$$G = 1 - \sum_{i=1}^k p_i^2 \quad (5.1)$$

其中 $k$ 是数据集中样本分类的数量,  $p_i$ 表示第 $i$ 类样本占总样本的比例。如果某一属性取多个值, 则按照每一个值所占的比重进行加权平均。

例如, 对于下面这些样本:

| day | deadline? | party? | lazy? | activity |
|-----|-----------|--------|-------|----------|
| 1   | urgent    | yes    | yes   | party    |
| 2   | urgent    | no     | yes   | study    |
| 3   | near      | yes    | yes   | party    |
| 4   | none      | yes    | no    | party    |
| 5   | none      | no     | yes   | pub      |
| 6   | none      | yes    | no    | party    |
| 7   | near      | no     | no    | study    |
| 8   | near      | no     | yes   | TV       |
| 9   | near      | yes    | yes   | party    |
| 10  | urgent    | no     | no    | study    |

以“deadline?”这个属性为例。首先计算deadline这个属性取每一个值的比例:

$$\begin{aligned}
 P(\text{deadline} = \text{urgent}) &= \frac{3}{10} \\
 P(\text{deadline} = \text{near}) &= \frac{4}{10} \\
 P(\text{deadline} = \text{none}) &= \frac{3}{10}
 \end{aligned} \tag{5.2}$$

然后分别计算deadline这个属性取每一个值下的Gini指数:

$$\begin{aligned}
 P(\text{deadline} = \text{urgent} \&\text{activity} = \text{party}) &= \frac{1}{3} \\
 P(\text{deadline} = \text{urgent} \&\text{activity} = \text{study}) &= \frac{2}{3}
 \end{aligned} \tag{5.3}$$

$$\begin{aligned}
 G(\text{urgent}) &= 1 - ((\frac{1}{3})^2 + (\frac{2}{3})^2) = \frac{4}{9} \\
 P(\text{deadline} = \text{near} \&\text{activity} = \text{party}) &= \frac{2}{4} \\
 P(\text{deadline} = \text{near} \&\text{activity} = \text{study}) &= \frac{1}{4} \\
 P(\text{deadline} = \text{near} \&\text{activity} = \text{TV}) &= \frac{1}{4}
 \end{aligned} \tag{5.4}$$

$$\begin{aligned}
 G(\text{near}) &= 1 - ((\frac{2}{4})^2 + (\frac{1}{4})^2 + (\frac{1}{4})^2) = \frac{5}{8} \\
 P(\text{deadline} = \text{none} \&\text{activity} = \text{party}) &= \frac{2}{3} \\
 P(\text{deadline} = \text{none} \&\text{activity} = \text{pub}) &= \frac{1}{3}
 \end{aligned} \tag{5.5}$$

$$G(\text{none}) = 1 - ((\frac{2}{3})^2 + (\frac{1}{3})^2) = \frac{4}{9}$$

最后按照取每一个值所占的比重对以上三个Gini指数做加权平均:

$$G_1 = G(\text{deadline}) = \frac{3}{10} \times \frac{4}{9} + \frac{4}{10} \times \frac{5}{8} + \frac{3}{10} \times \frac{4}{9} = \frac{31}{60} \tag{5.6}$$



同理可以算出按属性“party?”和“lazy?”切分时的Gini指数：

$$G_2 = G(\text{party}) = \frac{5}{10} \times [1 - (\frac{5}{5})^2] + \frac{5}{10} \times [1 - ((\frac{3}{5})^2 + (\frac{1}{5})^2 + (\frac{1}{5})^2)] = \frac{7}{25} \quad (5.7)$$

$$G_3 = G(\text{lazy}) = \frac{6}{10} \times [1 - ((\frac{3}{6})^2 + (\frac{1}{6})^2 + (\frac{1}{6})^2 + (\frac{1}{6})^2)] + \frac{4}{10} \times [1 - ((\frac{2}{4})^2 + (\frac{2}{4})^2)] = \frac{3}{5} \quad (5.8)$$

由于  $G_2 < G_1 < G_3$

```
1 double gini_index(int index,double value,int row, int col, double **dataset, double
  *class, int classnum)
2 {
3     float *numcount1 = (float *)malloc(classnum * sizeof(float));
4     float *numcount2 = (float *)malloc(classnum * sizeof(float));
5     for (int i = 0; i < classnum; i++)
6     {
7         numcount1[i]=numcount2[i]=0;
8     }
9     float count1 = 0, count2 = 0;
10    double gini1,gini2,gini;
11    gini1=gini2=gini=0;
12    // 计算每一类的个数
13    for (int i = 0; i < row; i++)
14    {
15        if (dataset[i][index] < value)
16        {
17            count1 ++;
18            for (int j = 0; j < classnum; j++)
19                if (dataset[i][col-1]==class[j])
20                    numcount1[j] += 1;
21        }
22        else
23        {
24            count2++;
25            for (int j = 0; j < classnum; j++)
26                if (dataset[i][col - 1] == class[j])
27                    numcount2[j]++;
28        }
29    }
30    // 判断分母是否为0，防止运算错误
31    if (count1==0)
32    {
33        gini1=1;
34        for (int i = 0; i < classnum; i++)
35            gini2 += (numcount2[i] / count2) * (numcount2[i] / count2);
36    }else if (count2==0)
37    {
38        gini2=1;
39        for (int i = 0; i < classnum; i++)
40            gini1 += (numcount1[i] / count1) * (numcount1[i] / count1);
41    }else
42    {
43        for (int i = 0; i < classnum; i++)
44        {
45            gini1 += (numcount1[i] / count1) * (numcount1[i] / count1);
46            gini2 += (numcount2[i] / count2) * (numcount2[i] / count2);
47        }
48    }
49    // 计算Gini指数
50    gini1 = 1 - gini1;
51    gini2 = 1 - gini2;
52    gini = (count1 / row) * gini1 + (count2 / row) * gini2;
```

```

53     free(numcount1);free(numcount2);
54     numcount1=numcount2=NULL;
55     return gini;
56 }

```

## 寻找最佳分割点

我们需要根据计算出的Gini指数来决定最佳的分割点。具体做法是计算所有切分点Gini指数，选出Gini指数最小的切分点作为最后的分割点。

- 功能：选取数据的最优切分点
- 输出：数据中最优切分点下的树结构

```

1  struct treeBranch *get_split(int row, int col, double **dataset, double *class, int
   classnum)
2  {
3      struct treeBranch *tree=(struct treeBranch *)malloc(sizeof(struct treeBranch));
4      int b_index=999;
5      double b_score = 999, b_value = 999,score;
6      // 计算所有切分点Gini系数，选出Gini系数最小的切分点
7      for (int i = 0; i < col-1; i++)
8      {
9          for (int j = 0; j < row; j++)
10         {
11             double value=dataset[j][i];
12             score=gini_index(i,value,row,col,dataset,class,classnum);
13             if (score<b_score)
14             {
15                 b_score=score;
16                 b_value=value;
17                 b_index=i;
18             }
19         }
20     }
21     tree->index=b_index;tree->value=b_value;tree->flag=0;
22     return tree;
23 }

```

## 计算叶子节点结果

我们不能让树一直生长下去，为此我们一般有两种方法来决定何时停止树的生长。

1. 最大树深。这是从树的根节点开始的最大节点数。一旦达到树的最大深度，就必须停止添加新节点。更深的树更复杂，可能更适合训练数据。
2. 最小节点记录。这是给定节点负责的最少训练模式。一旦达到或低于此最小值，我们必须停止拆分和添加新节点。训练模式很少的节点可能过于具体，可能会过度拟合训练数据。

当我们在某个点停止树的生长时，该节点称为终端节点或叶子节点，用于做出最终预测。这是通过获取分配给该节点的行组并选择该组中最常见的类值来完成的。下面这个函数将为一组行选择一个类值，它返回行列表中最常见的输出值。

- 功能：计算叶子节点结果
- 输出：输出最多的一类

```

1  double to_terminal(int row, int col, double **data, double *class, int classnum)
2  {
3      int *num=(int *)malloc(classnum*sizeof(classnum));
4      double maxnum=0;
5      int flag=0;
6      // 计算所有样本中结果最多的一类
7      for (int i = 0; i < classnum; i++)
8          num[i]=0;

```

```

9     for (int i = 0; i < row; i++)
10         for (int j = 0; j < classnum; j++)
11             if (data[i][col-1]==class[j])
12                 num[j]++;
13     for (int j = 0; j < classnum; j++)
14     {
15         if (num[j] > flag)
16         {
17             flag = num[j];
18             maxnum = class[j];
19         }
20     }
21     free(num);
22     num=NULL;
23     return maxnum;
24 }

```

## 分裂左右迭代

通过上述尝试，我们已经知道如何以及何时创建叶子节点。现在我们可以建立我们的树了。构建决策树需要在为每个节点创建的组上反复调用上面的get\_split()函数。添加到现有节点的新节点称为子节点。一个节点可以有零个子节点（叶子节点），一个子节点（在某一侧直接进行预测）或两个子节点。

创建节点后，我们可以通过再次调用相同的函数在拆分后的每一组数据上递归创建子节点。下面是实现此递归过程的函数。

- 功能：创建子树或生成叶子节点
- 输出：生成子树或叶子节点后的树

```

1 void split(struct treeBranch *tree, int row, int col, double **data, double *class,
2           int classnum, int depth, int min_size, int max_depth)
3 {
4     // 判断是否已经达到最大层数
5     if (depth>=max_depth)
6     {
7         tree->flag=1;
8         tree->output = to_terminal(row, col, data, class, classnum);
9         return;
10    }
11    struct dataset *childdata = test_split(tree->index, tree->value, row, col, data);
12    // 判断样本是否已被分为一边
13    if (childdata->row1==0 || childdata->row2==0)
14    {
15        tree->flag = 1;
16        tree->output = to_terminal(row, col, data, class, classnum);
17        return;
18    }
19    // 左子树，判断样本是否达到最小样本数，如不是则继续迭代
20    if (childdata->row1<=min_size)
21    {
22        struct treeBranch *leftchild = (struct treeBranch *)malloc(sizeof(struct
23        treeBranch));
24        leftchild->flag=1;
25        leftchild->output = to_terminal(childdata->row1, col, childdata-
26        >splitdata[0], class, classnum);
27        tree->leftBranch=leftchild;
28    }
29    else
30    {
31        struct treeBranch *leftchild = get_split(childdata->row1, col, childdata-
32        >splitdata[0], class, classnum);
33        tree->leftBranch=leftchild;
34    }
35    // 右子树
36    if (childdata->row2<=min_size)
37    {
38        struct treeBranch *rightchild = (struct treeBranch *)malloc(sizeof(struct
39        treeBranch));
40        rightchild->flag=1;
41        rightchild->output = to_terminal(childdata->row2, col, childdata-
42        >splitdata[1], class, classnum);
43        tree->rightBranch=rightchild;
44    }
45    else
46    {
47        struct treeBranch *rightchild = get_split(childdata->row2, col, childdata-
48        >splitdata[1], class, classnum);
49        tree->rightBranch=rightchild;
50    }
51    return tree;
52 }

```

```

30     split(leftchild, childdata->row1, col, childdata->splitdata[0], class,
classnum, depth+1, min_size, max_depth);
31 }
32 // 右子树, 判断样本是否达到最小样本数, 如不是则继续迭代
33 if (childdata->row2 <= min_size)
34 {
35     struct treeBranch *rightchild = (struct treeBranch *)malloc(sizeof(struct
treeBranch));
36     rightchild->flag = 1;
37     rightchild->output = to_terminal(childdata->row2, col, childdata-
>splitdata[1], class, classnum);
38     tree->rightBranch = rightchild;
39 }
40 else
41 {
42     struct treeBranch *rightchild = get_split(childdata->row2, col, childdata-
>splitdata[1], class, classnum);
43     tree->rightBranch = rightchild;
44     split(rightchild, childdata->row2, col, childdata->splitdata[1], class,
classnum, depth + 1, min_size, max_depth);
45 }
46 free(childdata->splitdata); childdata->splitdata=NULL;
47 free(childdata); childdata=NULL;
48 return;
49 }

```

## 建立决策树

下面, 我们就可以利用上面编写的函数构建根节点并调用split()函数, 然后进行递归调用以构建整个树。

- 功能: 生成决策树
- 输出: 生成后的决策树

```

1 struct treeBranch *build_tree(int row, int col, double **data, int min_size, int
max_depth)
2 {
3     int count1 = 0, flag1 = 0;
4     // 判断结果一共有多少类别, 此处classes[20]仅仅是取一个较大的数20, 默认类别不可能超过20类
5     double classes[20];
6     for (int i = 0; i < row; i++)
7     {
8         if (count1 == 0)
9         {
10             classes[0] = data[i][col - 1];
11             count1++;
12         }
13         else
14         {
15             flag1 = 0;
16             for (int j = 0; j < count1; j++)
17                 if (classes[j] == data[i][col - 1])
18                     flag1 = 1;
19             if (flag1 == 0)
20             {
21                 classes[count1] = data[i][col - 1];
22                 count1++;
23             }
24         }
25     }
26     // 生成切分点
27     struct treeBranch *result = get_split(row, col, data, classes, count1);
28     // 进入迭代, 不断生成子树

```

```

29     split(result, row, col, data, classes, count1, 1, min_size, max_depth);
30     return result;
31 }

```

### 3.5.3 算法代码

我们现在知道了如何实现**决策树算法**，那么我们把它应用到**钞票数据集 banknote.csv**

我们给出链接：<https://aistudio.baidu.com/aistudio/datasetdetail/105756/0>

### C语言细节讲解

本节假设您已下载数据集 `banknote.csv`，并且它在当前工作目录中可用。下面我们给出一个完整实例，使用C语言详细讲解每一处细节。我们给出每一个.c文件的所有代码：

#### 1) read\_csv.c

该文件代码对头文件作出调整

```

1  #include "DT.h"
2
3  //获取行数
4  int get_row(char *filename)
5  {
6      char line[1024];
7      int i = 0;
8      FILE *stream = fopen(filename, "r");
9      while (fgets(line, 1024, stream))
10     {
11         i++;
12     }
13     fclose(stream);
14     return i;
15 }
16
17 //获取列数
18 int get_col(char *filename)
19 {
20     char line[1024];
21     int i = 0;
22     FILE *stream = fopen(filename, "r");
23     fgets(line, 1024, stream);
24     char *token = strtok(line, ",");
25     while (token)
26     {
27         token = strtok(NULL, ",");
28         i++;
29     }
30     fclose(stream);
31     return i;
32 }
33
34 // 获取完整数据集
35 void get_two_dimension(char *line, double **data, char *filename)
36 {
37     FILE *stream = fopen(filename, "r");
38     int i = 0;
39     while (fgets(line, 1024, stream)) //逐行读取
40     {
41         int j = 0;
42         char *tok;
43         char *tmp = strdup(line);
44         for (tok = strtok(line, ","); tok && *tok; j++, tok = strtok(NULL, ",\n"))

```

```

45     {
46         data[i][j] = atof(tok); //转换成浮点数
47     } //字符串拆分操作
48     i++;
49     free(tmp);
50 }
51 fclose(stream); //文件打开后要进行关闭操作
52 }

```

## 2) k\_fold.c

该文件代码对头文件作出调整

```

1  #include "DT.h"
2
3  double ***cross_validation_split(double **dataset, int row, int n_folds, int
fold_size)
4  {
5      srand(10); //种子
6      double ***split;
7      int i, j = 0, k = 0;
8      int index;
9      double **fold;
10     split = (double ***)malloc(n_folds * sizeof(double **));
11     for (i = 0; i < n_folds; i++)
12     {
13         fold = (double **)malloc(fold_size * sizeof(double *));
14         while (j < fold_size)
15         {
16             fold[j] = (double *)malloc(col * sizeof(double));
17             index = rand() % row;
18             fold[j] = dataset[index];
19             for (k = index; k < row - 1; k++) //for循环删除这个数组中被rand取到的元素
20             {
21                 dataset[k] = dataset[k + 1];
22             }
23             row--; //每次随机取出一个后总行数-1, 保证不会重复取某一行
24             j++;
25         }
26         j = 0; //清零j
27         split[i] = fold;
28     }
29     return split;
30 }

```

## 3) DT.c

核心函数部分，用以构建整棵决策树，并给出决策树的预测结果。

```

1  #include "DT.h"
2
3  // 切分函数，根据切分点将数据分为左右两组
4  struct dataset *test_split(int index, double value, int row, int col, double **data)
5  {
6      // 将切分结果作为结构体返回
7      struct dataset *split = (struct dataset *)malloc(sizeof(struct dataset));
8      int count1 = 0, count2 = 0;
9      double ***groups = (double ***)malloc(2 * sizeof(double **));
10     for (int i = 0; i < 2; i++)
11     {
12         groups[i] = (double **)malloc(row * sizeof(double *));
13         for (int j = 0; j < row; j++)

```

```

14     {
15         groups[i][j] = (double *)malloc(col * sizeof(double));
16     }
17 }
18 for (int i = 0; i < row; i++)
19 {
20     if (data[i][index] < value)
21     {
22         groups[0][count1] = data[i];
23         count1++;
24     }
25     else
26     {
27         groups[1][count2] = data[i];
28         count2++;
29     }
30 }
31 split->splitdata = groups;
32 split->row1 = count1;
33 split->row2 = count2;
34 return split;
35 }
36
37 // 计算Gini系数
38 double gini_index(int index, double value, int row, int col, double **dataset,
39 double *class, int classnum)
40 {
41     float *numcount1 = (float *)malloc(classnum * sizeof(float));
42     float *numcount2 = (float *)malloc(classnum * sizeof(float));
43     for (int i = 0; i < classnum; i++)
44         numcount1[i] = numcount2[i] = 0;
45
46     float count1 = 0, count2 = 0;
47     double gini1, gini2, gini;
48     gini1 = gini2 = gini = 0;
49     // 计算每一类的个数
50     for (int i = 0; i < row; i++)
51     {
52         if (dataset[i][index] < value)
53         {
54             count1++;
55             for (int j = 0; j < classnum; j++)
56                 if (dataset[i][col - 1] == class[j])
57                     numcount1[j] += 1;
58         }
59         else
60         {
61             count2++;
62             for (int j = 0; j < classnum; j++)
63                 if (dataset[i][col - 1] == class[j])
64                     numcount2[j]++;
65         }
66     }
67     // 判断分母是否为0，防止运算错误
68     if (count1 == 0)
69     {
70         gini1 = 1;
71         for (int i = 0; i < classnum; i++)
72             gini2 += (numcount2[i] / count2) * (numcount2[i] / count2);
73     }
74     else if (count2 == 0)
75     {
76         gini2 = 1;

```

```

76         for (int i = 0; i < classnum; i++)
77             gini1 += (numcount1[i] / count1) * (numcount1[i] / count1);
78     }
79     else
80     {
81         for (int i = 0; i < classnum; i++)
82         {
83             gini1 += (numcount1[i] / count1) * (numcount1[i] / count1);
84             gini2 += (numcount2[i] / count2) * (numcount2[i] / count2);
85         }
86     }
87     // 计算Gini系数
88     gini1 = 1 - gini1;
89     gini2 = 1 - gini2;
90     gini = (count1 / row) * gini1 + (count2 / row) * gini2;
91     free(numcount1);
92     free(numcount2);
93     numcount1 = numcount2 = NULL;
94     return gini;
95 }
96
97 // 选取数据的最优切分点
98 struct treeBranch *get_split(int row, int col, double **dataset, double *class, int
classnum)
99 {
100     struct treeBranch *tree = (struct treeBranch *)malloc(sizeof(struct
treeBranch));
101     int b_index = 999;
102     double b_score = 999, b_value = 999, score;
103     // 计算所有切分点Gini系数, 选出Gini系数最小的切分点
104     for (int i = 0; i < col - 1; i++)
105     {
106         for (int j = 0; j < row; j++)
107         {
108             double value = dataset[j][i];
109             score = gini_index(i, value, row, col, dataset, class, classnum);
110             if (score < b_score)
111             {
112                 b_score = score;
113                 b_value = value;
114                 b_index = i;
115             }
116         }
117     }
118     tree->index = b_index;
119     tree->value = b_value;
120     tree->flag = 0;
121     return tree;
122 }
123
124 // 计算叶节点结果
125 double to_terminal(int row, int col, double **data, double *class, int classnum)
126 {
127     int *num = (int *)malloc(classnum * sizeof(classnum));
128     double maxnum = 0;
129     int flag = 0;
130     // 计算所有样本中结果最多的一类
131     for (int i = 0; i < classnum; i++)
132         num[i] = 0;
133     for (int i = 0; i < row; i++)
134         for (int j = 0; j < classnum; j++)
135             if (data[i][col - 1] == class[j])
136                 num[j]++;

```



```

137     for (int j = 0; j < classnum; j++)
138     {
139         if (num[j] > flag)
140         {
141             flag = num[j];
142             maxnum = class[j];
143         }
144     }
145     free(num);
146     num = NULL;
147     return maxnum;
148 }
149
150 // 创建子树或生成叶节点
151 void split(struct treeBranch *tree, int row, int col, double **data, double *class,
152 int classnum, int depth, int min_size, int max_depth)
153 {
154     // 判断是否已经达到最大层数
155     if (depth >= max_depth)
156     {
157         tree->flag = 1;
158         tree->output = to_terminal(row, col, data, class, classnum);
159         return;
160     }
161     struct dataset *childdata = test_split(tree->index, tree->value, row, col,
162 data);
163     // 判断样本是否已被分为一边
164     if (childdata->row1 == 0 || childdata->row2 == 0)
165     {
166         tree->flag = 1;
167         tree->output = to_terminal(row, col, data, class, classnum);
168         return;
169     }
170     // 左子树, 判断样本是否达到最小样本数, 如不是则继续迭代
171     if (childdata->row1 <= min_size)
172     {
173         struct treeBranch *leftchild = (struct treeBranch *)malloc(sizeof(struct
174 treeBranch));
175         leftchild->flag = 1;
176         leftchild->output = to_terminal(childdata->row1, col, childdata-
177 >splitdata[0], class, classnum);
178         tree->leftBranch = leftchild;
179     }
180     else
181     {
182         struct treeBranch *leftchild = get_split(childdata->row1, col, childdata-
183 >splitdata[0], class, classnum);
184         tree->leftBranch = leftchild;
185         split(leftchild, childdata->row1, col, childdata->splitdata[0], class,
186 classnum, depth + 1, min_size, max_depth);
187     }
188     // 右子树, 判断样本是否达到最小样本数, 如不是则继续迭代
189     if (childdata->row2 <= min_size)
190     {
191         struct treeBranch *rightchild = (struct treeBranch *)malloc(sizeof(struct
192 treeBranch));
193         rightchild->flag = 1;
194         rightchild->output = to_terminal(childdata->row2, col, childdata-
195 >splitdata[1], class, classnum);
196         tree->rightBranch = rightchild;
197     }
198     else
199     {
200         struct treeBranch *rightchild = get_split(childdata->row2, col, childdata-
201 >splitdata[1], class, classnum);
202         tree->rightBranch = rightchild;
203         split(rightchild, childdata->row2, col, childdata->splitdata[1], class,
204 classnum, depth + 1, min_size, max_depth);
205     }
206 }

```

```

192     struct treeBranch *rightchild = get_split(childdata->row2, col, childdata->
>splitdata[1], class, classnum);
193     tree->rightBranch = rightchild;
194     split(rightchild, childdata->row2, col, childdata->splitdata[1], class,
classnum, depth + 1, min_size, max_depth);
195 }
196 free(childdata->splitdata);
197 childdata->splitdata = NULL;
198 free(childdata);
199 childdata = NULL;
200 return;
201 }
202
203 // 生成决策树
204 struct treeBranch *build_tree(int row, int col, double **data, int min_size, int
max_depth)
205 {
206     int count1 = 0, flag1 = 0;
207     // 判断结果一共有多少类别, 此处classes[20]仅仅是取一个较大的数20, 默认类别不可能超过20类
208     double classes[20];
209     for (int i = 0; i < row; i++)
210     {
211         if (count1 == 0)
212         {
213             classes[0] = data[i][col - 1];
214             count1++;
215         }
216         else
217         {
218             flag1 = 0;
219             for (int j = 0; j < count1; j++)
220                 if (classes[j] == data[i][col - 1])
221                     flag1 = 1;
222             if (flag1 == 0)
223             {
224                 classes[count1] = data[i][col - 1];
225                 count1++;
226             }
227         }
228     }
229     // 生成切分点
230     struct treeBranch *result = get_split(row, col, data, classes, count1);
231     // 进入迭代, 不断生成子树
232     split(result, row, col, data, classes, count1, 1, min_size, max_depth);
233     return result;
234 }
235
236 // 决策树预测
237 double predict(double *test, struct treeBranch *tree)
238 {
239     double output;
240     // 判断是否达到叶节点, flag=1时为叶节点, flag=0时则继续判断
241     if (tree->flag == 1)
242     {
243         output = tree->output;
244         return output;
245     }
246     else
247     {
248         if (test[tree->index] < tree->value)
249         {
250             output = predict(test, tree->leftBranch);
251             return output;

```

```

252     }
253     else
254     {
255         output = predict(test, tree->rightBranch);
256         return output;
257     }
258 }
259 }

```

#### 4) DT.h

构建决策树所需要包含的头文件。

```

1  #ifndef DT
2  #define DT
3
4  #include <stdio.h>
5  #include <string.h>
6  #include <stdlib.h>
7
8  double **dataset;
9  int row, col;
10
11 struct treeBranch
12 {
13     int flag;
14     int index;
15     double value;
16     double output;
17     struct treeBranch *leftBranch;
18     struct treeBranch *rightBranch;
19 };
20
21 struct dataset
22 {
23     int row1;
24     int row2;
25     double ***splitdata;
26 };
27
28 int get_row(char *filename);
29 int get_col(char *filename);
30 void get_two_dimension(char *line, double **dataset, char *filename);
31 double ***cross_validation_split(double **dataset, int row, int n_folds, int
fold_size);
32 struct dataset *test_split(int index, double value, int row, int col, double **data);
33 double gini_index(int index, double value, int row, int col, double **dataset, double
*class, int classnum);
34 struct treeBranch *get_split(int row, int col, double **dataset, double *class, int
classnum);
35 double to_terminal(int row, int col, double **data, double *class, int classnum);
36 void split(struct treeBranch *tree, int row, int col, double **data, double *class,
int classnum, int depth, int min_size, int max_depth);
37 struct treeBranch *build_tree(int row, int col, double **data, int min_size, int
max_depth);
38 double predict(double *test, struct treeBranch *tree);
39 float *evaluate_algorithm(double **dataset, int column, int n_folds, int fold_size,
int min_size, int max_depth);
40 float accuracy_metric(double *actual, double *predicted, int fold_size);
41 double *get_test_prediction(double **train, double **test, int column, int min_size,
int max_depth, int fold_size, int train_size);
42

```

## 5) score.c

该文件代码对头文件作出调整

```

1  #include "DT.h"
2
3  float accuracy_metric(double *actual, double *predicted, int fold_size)
4  {
5      int correct = 0;
6      for (int i = 0; i < fold_size; i++)
7      {
8          if (actual[i] == predicted[i])
9              correct += 1;
10     }
11     return (correct / (float)fold_size) * 100.0;
12 }
```

## 6) test\_prediction.c

我们将上述的预测在测试集上也进行一遍，由此判断模型对于没有见过的数据会做出怎样的预测，方便进一步对模型的好坏作出评估。

```

1  #include "DT.h"
2
3  double *get_test_prediction(double **train, double **test, int column, int min_size,
4  int max_depth, int fold_size, int train_size)
5  {
6      double *predictions = (double *)malloc(fold_size * sizeof(double)); //预测集的行数
7      //就是数组prediction的长度
8      struct treeBranch *tree = build_tree(train_size, column, train, min_size,
9      max_depth);
10     for (int i = 0; i < fold_size; i++)
11     {
12         predictions[i] = predict(test[i], tree);
13     }
14     return predictions; //返回对test的预测数组
15 }
```

## 7) evaluate.c

```

1  #include "DT.h"
2
3  float *evaluate_algorithm(double **dataset, int column, int n_folds, int fold_size,
4  int min_size, int max_depth)
5  {
6      double ***split = cross_validation_split(dataset, row, n_folds, fold_size);
7      int i, j, k, l;
8      int test_size = fold_size;
9      int train_size = fold_size * (n_folds - 1); //train_size个一维数组
10     float *score = (float *)malloc(n_folds * sizeof(float));
11     for (i = 0; i < n_folds; i++)
12     { //因为要遍历删除，所以拷贝一份split
13         double ***split_copy = (double ***)malloc(n_folds * sizeof(double **));
14         for (j = 0; j < n_folds; j++)
15         {
16             split_copy[j] = (double **)malloc(fold_size * sizeof(double **));
17             for (k = 0; k < fold_size; k++)
18             {
19                 split_copy[j][k] = (double *)malloc(column * sizeof(double));
20             }
21         }
22     }
23 }
```

```

19     }
20 }
21 for (j = 0; j < n_folds; j++)
22 {
23     for (k = 0; k < fold_size; k++)
24     {
25         for (l = 0; l < column; l++)
26         {
27             split_copy[j][k][l] = split[j][k][l];
28         }
29     }
30 }
31 double **test_set = (double **)malloc(test_size * sizeof(double *));
32 for (j = 0; j < test_size; j++)
33 { //对test_size中的每一行
34     test_set[j] = (double *)malloc(column * sizeof(double));
35     for (k = 0; k < column; k++)
36     {
37         test_set[j][k] = split_copy[i][j][k];
38     }
39 }
40 for (j = i; j < n_folds - 1; j++)
41 {
42     split_copy[j] = split_copy[j + 1];
43 }
44 double **train_set = (double **)malloc(train_size * sizeof(double *));
45 for (k = 0; k < n_folds - 1; k++)
46 {
47     for (l = 0; l < fold_size; l++)
48     {
49         train_set[k * fold_size + l] = (double *)malloc(column *
sizeof(double));
50         train_set[k * fold_size + l] = split_copy[k][l];
51     }
52 }
53 double *predicted = (double *)malloc(test_size * sizeof(double)); //predicted
有test_size个
54 predicted = get_test_prediction(train_set, test_set, column, min_size,
max_depth, fold_size, train_size);
55 double *actual = (double *)malloc(test_size * sizeof(double));
56 for (l = 0; l < test_size; l++)
57 {
58     actual[l] = test_set[l][column - 1];
59 }
60 float accuracy = accuracy_metric(actual, predicted, test_size);
61 score[i] = accuracy;
62 printf("score[%d] = %f%%\n", i, score[i]);
63 free(split_copy);
64 }
65 float total = 0.0;
66 for (l = 0; l < n_folds; l++)
67 {
68     total += score[l];
69 }
70 printf("mean_accuracy = %f%%\n", total / n_folds);
71 return score;
72 }

```

## 8) main.c

```
1  #include "DT.h"
2
3  int main()
4  {
5      char filename[] = "banknote.csv";
6      char line[1024];
7      row = get_row(filename);
8      col = get_col(filename);
9      dataset = (double **)malloc(row * sizeof(int *));
10     for (int i = 0; i < row; ++i)
11     {
12         dataset[i] = (double *)malloc(col * sizeof(double));
13     } //动态申请二维数组
14     get_two_dimension(line, dataset, filename);
15
16     // CART参数, 分别为叶节点最小样本数和树最大层数
17     int min_size = 5, max_depth = 10;
18     int n_folds = 5;
19     int fold_size = (int)(row / n_folds);
20
21     // CART决策树, 返回交叉验证正确率
22     float* score = evaluate_algorithm(dataset, col, n_folds, fold_size, min_size,
23                                     max_depth);
24 }
```

## 9) compile.sh

```
1  gcc main.c read_csv.c normalize.c k_fold.c evaluate.c score.c test_prediction.c -o run
   -lm && ./run
```

编译&运行:

```
1  bash compile.sh
```

最终输出结果如下:

```
1  score[0] = 97.080292%
2  score[1] = 97.810219%
3  score[2] = 96.715332%
4  score[3] = 98.905113%
5  score[4] = 98.175179%
6  mean_accuracy = 97.737228%
```

## Python语言实战

本节同样假设您已经下载数据集, 我们使用著名机器学习开源库sklearn高效实现**决策树算法**, 以便您在实战中使用该算法:

```
1  import pandas as pd
2  import numpy as np
3  from sklearn.model_selection import KFold
4  from sklearn.tree import DecisionTreeClassifier
5  from sklearn.metrics import accuracy_score
6
7
8  if __name__ == '__main__':
9      dataset = np.array(pd.read_csv("banknote.csv", sep=',', header=None))
10     k_cross = KFold(n_splits=5, random_state=0, shuffle=True)
```

```

11     index = 0
12     score = np.array([])
13     data, label = dataset[:, :-1], dataset[:, -1]
14     for train_index, test_index in k_Cross.split(dataset):
15         train_data, train_label = data[train_index, :], label[train_index]
16         test_data, test_label = data[test_index, :], label[test_index]
17         model = DecisionTreeClassifier()
18         model.fit(train_data, train_label)
19         pred = model.predict(test_data)
20         acc = accuracy_score(test_label, pred)
21         score = np.append(score, acc)
22         print('score[{}] = {}'.format(index, acc))
23         index += 1
24     print('mean_accuracy = {}'.format(np.mean(score)))

```

输出结果如下：

```

1 score[0] = 0.9927272727272727%
2 score[1] = 0.9854545454545455%
3 score[2] = 0.9671532846715328%
4 score[3] = 0.9963503649635036%
5 score[4] = 0.9635036496350365%
6 mean_accuracy = 0.9810378234903782%

```

## 3.6 Naive Bayes

朴素贝叶斯算法，是应用最为广泛的分类算法之一。该算法利用贝叶斯定理与特征条件独立假设做预测，直接且易于理解。该算法在实际运用中，往往能得到意想不到的好结果。

### 3.6.1 算法介绍

朴素贝叶斯算法的其本质就是计算 $P(class|data)$ ，即数据 $data$ 属于某一类别 $class$ 的概率。

朴素贝叶斯算法的核心就是贝叶斯公式，贝叶斯公式为我们提供了一个适用于计算一些数据属于某一类别的概率的计算方法。贝叶斯公式如下：

$$P(class|data) = \frac{P(data|class) \times P(class)}{P(data)}$$

其中， $P(class|data)$ 表示 $data$ 属于某个 $class$ 的概率。同时，上式假设各个特征条件是独立的。

我们认为，使得 $P(class|data)$ 最大的 $class$ 就是该 $data$ 所属的 $class$ ，从而我们可以预测出该数据所属类别。

下面，我们将结合代码讲解，朴素贝叶斯算法是如何计算 $P(class|data)$ ，进而对 $data$ 做预测的。

### 3.6.2 算法讲解

#### 数据分类

我们需把所有的数据按照各自的类别进行分类，组成一个新的数组。下面先给出将数据的分类函数：

```

1 double*** separate_by_class(double **dataset, int class_num, int *class_num_list, int
  row, int col) {
2     double ***separated;
3     separated = (double***)malloc(class_num * sizeof(double**));
4     int i, j;
5     for (i = 0; i < class_num; i++) {
6         separated[i] = (double**)malloc(class_num_list[i] * sizeof(double *));
7         for (j = 0; j < class_num_list[i]; j++) {
8             separated[i][j] = (double*)malloc(col * sizeof(double));
9         }

```

```

10     }
11     int* index = (int *)malloc(class_num * sizeof(int));
12     for (i = 0; i < class_num; i++) {
13         index[i] = 0;
14     }
15     for (i = 0; i < row; i++) {
16         for (j = 0; j < class_num; j++) {
17             if (dataset[i][col - 1] == j) {
18                 separated[j][index[j]] = dataset[i];
19                 index[j]++;
20             }
21         }
22     }
23     return separated;
24 }

```

以下面的10条数据为例，利用上述函数对数据进行分类

|    |          |          |          |
|----|----------|----------|----------|
| 1  | x1       | x2       | Lable    |
| 2  | 2.000000 | 2.000000 | 0.000000 |
| 3  | 2.005000 | 1.995000 | 0.000000 |
| 4  | 2.010000 | 1.990000 | 0.000000 |
| 5  | 2.015000 | 1.985000 | 0.000000 |
| 6  | 2.020000 | 1.980000 | 0.000000 |
| 7  |          |          |          |
| 8  | 5.000000 | 5.000000 | 1.000000 |
| 9  | 5.005000 | 4.995000 | 1.000000 |
| 10 | 5.010000 | 4.990000 | 1.000000 |
| 11 | 5.015000 | 4.985000 | 1.000000 |
| 12 | 5.020000 | 4.980000 | 1.000000 |

代码如下：

```

1  double*** separate_by_class(double **dataset,int class_num, int *class_num_list, int
row, int col) {
2      double ***separated;
3      separated = (double***)malloc(class_num * sizeof(double**));
4      int i, j;
5      for (i = 0; i < class_num; i++) {
6          separated[i] = (double**)malloc(class_num_list[i] * sizeof(double *));
7          for (j = 0; j < class_num_list[i]; j++) {
8              separated[i][j] = (double*)malloc(col * sizeof(double));
9          }
10     }
11     int* index = (int *)malloc(class_num * sizeof(int));
12     for (i = 0; i < class_num; i++) {
13         index[i] = 0;
14     }
15     for (i = 0; i < row; i++) {
16         for (j = 0; j < class_num; j++) {
17             if (dataset[i][col - 1] == j) {
18                 separated[j][index[j]] = dataset[i];
19                 index[j]++;
20             }
21         }
22     }
23     return separated;
24 }
25
26 void main() {
27     double **dataset;

```



```

28     dataset = (double **)malloc(row * sizeof(double *));
29     for (int i = 0; i < row; ++i) {
30         dataset[i] = (double *)malloc(col * sizeof(double));
31     }
32     for (int i = 0; i < 5; i++) {
33         dataset[i][0] = 2 + i * 0.005;
34         dataset[i][1] = 2 - i * 0.005;
35         dataset[i][2] = 0;
36     }
37     for (int i = 0; i < 5; i++) {
38         dataset[i+5][0] = 5 + i * 0.005;
39         dataset[i+5][1] = 5 - i * 0.005;
40         dataset[i + 5][2] = 1;
41     }
42     int class_num_list[2] = {5,5};
43     double ***separated = separate_by_class(dataset, 2, class_num_list, 10, 3);
44     // 输出结果
45     for (int i = 0; i < 2; i++) {
46         //先按照类别输出
47         for (int j = 0; j < 5; j++) {
48             for (int k = 0; k < 3; k++) {
49                 printf("%f\t", separated[i][j][k]);
50             }
51             printf("\n");
52         }
53         printf("\n");
54     }
55 }

```

输出结果如下:

|    |          |          |          |
|----|----------|----------|----------|
| 1  | 2.000000 | 2.000000 | 0.000000 |
| 2  | 2.005000 | 1.995000 | 0.000000 |
| 3  | 2.010000 | 1.990000 | 0.000000 |
| 4  | 2.015000 | 1.985000 | 0.000000 |
| 5  | 2.020000 | 1.980000 | 0.000000 |
| 6  |          |          |          |
| 7  | 5.000000 | 5.000000 | 1.000000 |
| 8  | 5.005000 | 4.995000 | 1.000000 |
| 9  | 5.010000 | 4.990000 | 1.000000 |
| 10 | 5.015000 | 4.985000 | 1.000000 |
| 11 | 5.020000 | 4.980000 | 1.000000 |

## 计算统计量

我们需要计算数据的均值与标准差, 其公式与前文中提到的一致, 如下所示:

$$mean = \frac{\sum_{i=1}^n x_i}{count(x)}$$

$$standard\ deviation = \sqrt{\frac{\sum_{i=1}^n (x_i - mean(x))^2}{count(x) - 1}}$$

其代码如下:

```

1 double get_mean(double**dataset, int row, int col) {
2     int i;
3     double mean = 0;
4     for (i = 0; i < row; i++) {
5         mean += dataset[i][col];
6     }
7     return mean / row;

```

```

8 }
9
10 double get_std(double**dataset, int row, int col) {
11     int i;
12     double mean = 0;
13     double std = 0;
14     for (i = 0; i < row; i++) {
15         mean += dataset[i][col];
16     }
17     mean /= row;
18     for (i = 0; i < row; i++) {
19         std += pow((dataset[i][col]-mean),2);
20     }
21     return sqrt(std / (row - 1));
22 }

```

仍然以下面的10条数据为例，利用上述函数按照类别计算数据的统计量

|    |          |          |          |
|----|----------|----------|----------|
| 1  | x1       | x2       | Lable    |
| 2  | 2.000000 | 2.000000 | 0.000000 |
| 3  | 2.005000 | 1.995000 | 0.000000 |
| 4  | 2.010000 | 1.990000 | 0.000000 |
| 5  | 2.015000 | 1.985000 | 0.000000 |
| 6  | 2.020000 | 1.980000 | 0.000000 |
| 7  |          |          |          |
| 8  | 5.000000 | 5.000000 | 1.000000 |
| 9  | 5.005000 | 4.995000 | 1.000000 |
| 10 | 5.010000 | 4.990000 | 1.000000 |
| 11 | 5.015000 | 4.985000 | 1.000000 |
| 12 | 5.020000 | 4.980000 | 1.000000 |

代码如下：

```

1 double get_mean(double**dataset, int row, int col) {
2     int i;
3     double mean = 0;
4     for (i = 0; i < row; i++) {
5         mean += dataset[i][col];
6     }
7     return mean / row;
8 }
9
10 double get_std(double**dataset, int row, int col) {
11     int i;
12     double mean = 0;
13     double std = 0;
14     for (i = 0; i < row; i++) {
15         mean += dataset[i][col];
16     }
17     mean /= row;
18     for (i = 0; i < row; i++) {
19         std += pow((dataset[i][col]-mean),2);
20     }
21     return sqrt(std / (row - 1));
22 }
23
24 double*** separate_by_class(double **dataset,int class_num, int *class_num_list, int
row, int col) {
25     double ***separated;
26     separated = (double***)malloc(class_num * sizeof(double**));
27     int i, j;

```

```

28     for (i = 0; i < class_num; i++) {
29         separated[i] = (double**)malloc(class_num_list[i] * sizeof(double *));
30         for (j = 0; j < class_num_list[i]; j++) {
31             separated[i][j] = (double*)malloc(col * sizeof(double));
32         }
33     }
34     int* index = (int *)malloc(class_num * sizeof(int));
35     for (i = 0; i < class_num; i++) {
36         index[i] = 0;
37     }
38     for (i = 0; i < row; i++) {
39         for (j = 0; j < class_num; j++) {
40             if (dataset[i][col - 1] == j) {
41                 separated[j][index[j]] = dataset[i];
42                 index[j]++;
43             }
44         }
45     }
46     return separated;
47 }
48
49 double** summarize_dataset(double **dataset, int row, int col) {
50     int i;
51     double **summary = (double**)malloc((col - 1) * sizeof(double *));
52     for (i = 0; i < (col - 1); i++) {
53         summary[i] = (double*)malloc(2 * sizeof(double));
54         summary[i][0] = get_mean(dataset, row, i);
55         summary[i][1] = get_std(dataset, row, i);
56     }
57     return summary;
58 }
59
60 double*** summarize_by_class(double **train, int class_num, int *class_num_list, int
row, int col) {
61     int i;
62     double ***summarize;
63     summarize = (double***)malloc(class_num * sizeof(double**));
64     double ***separate = separate_by_class(train, class_num, class_num_list, row,
col);
65     for (i = 0; i < class_num; i++) {
66         summarize[i] = summarize_dataset(separate[i], class_num_list[i], col);
67     }
68     return summarize;
69 }
70
71 void main() {
72     int row = 10;
73     int col = 3;
74     int class_num = 2;
75     int class_num_list[2] = { 5, 5 };
76     double** dataset;
77     dataset = (double**)malloc(row * sizeof(double*));
78     for (int i = 0; i < row; ++i) {
79         dataset[i] = (double*)malloc(col * sizeof(double));
80     }
81     for (int i = 0; i < 5; i++) {
82         dataset[i][0] = 2 + i * 0.005;
83         dataset[i][1] = 2 - i * 0.005;
84         dataset[i][2] = 0;
85     }
86     for (int i = 0; i < 5; i++) {
87         dataset[i + 5][0] = 5 + i * 0.005;
88         dataset[i + 5][1] = 5 - i * 0.005;

```

```

89     dataset[i + 5][2] = 1;
90 }
91 double*** summarize = summarize_by_class(dataset, class_num, class_num_list,
row, col);
92 for (int i = 0; i < 2; i++) {
93     //先按照类别输出
94     for (int j = 0; j < 2; j++) {
95         for (int k = 0; k < 2; k++) {
96             printf("%f\t", summarize[i][j][k]);
97         }
98         printf("\n");
99     }
100     printf("\n");
101 }
102 }

```

按照类别依次得到每列数据的均值与方差：

|   |          |          |
|---|----------|----------|
| 1 | 2.010000 | 0.007906 |
| 2 | 1.990000 | 0.007906 |
| 3 |          |          |
| 4 | 5.010000 | 0.007906 |
| 5 | 4.990000 | 0.007906 |

## 高斯概率分布函数

高斯概率密度函数表达式为：

$$probability(x) = \frac{e^{-\frac{(x - mean(x))^2}{2 \times standard\_deviation^2}}}{\sqrt{2 \times PI \times standard\_deviation}}$$

计算代码如下：

```

1 double calculate_probability(double x, double mean, double std)
2 {
3     double pi = acos(-1.0);
4     double p = 1 / (pow(2 * pi, 0.5) * std) * exp(-(pow((x - mean), 2) / (2 * pow(std,
2)))));
5     return p;
6 }

```

## 类别概率

下面就是朴素贝叶斯的关键——计算数据属于某一类别的概率。

代码如下：

```

1 double* calculate_class_probabilities(double ***summaries, double *test_row, int
  class_num, int *class_num_list, int row, int col) {
2     int i, j;
3     double *probabilities = (double *)malloc(class_num * sizeof(double));
4     for (i = 0; i < class_num; i++) {
5         probabilities[i] = (double)class_num_list[i] / row;
6     }
7     for (i = 0; i < class_num; i++) {
8         for (j = 0; j < col-1; j++) {
9             probabilities[i] *= calculate_probability(test_row[j], summaries[i][j]
  [0], summaries[i][j][1]);
10        }
11    }
12    return probabilities;
13 }

```

### 3.6.3 算法代码

我们现在知道了如何实现朴素贝叶斯算法，那么我们把它应用到[鸢尾花数据集 iris.csv](#)

我们给出链接：<https://aistudio.baidu.com/aistudio/datasetdetail/105756/0>

### C语言细节讲解

本节假设您已下载数据集 `iris.csv`，并且它在当前工作目录中可用。下面我们给出一个完整实例，使用C语言详细讲解每一处细节。我们给出每一个.c文件的所有代码：

#### 1) read\_csv.c

该文件代码与前面代码一致，不再重复给出。

#### 2) normalize.c

该文件代码与前面代码一致，不再重复给出。

#### 3) k\_fold.c

该文件代码与前面代码一致，不再重复给出。

#### 4) score.c

该文件代码与前面代码一致，不再重复给出。

#### 5) test\_prediction.c

```

1 #include <stdlib.h>
2 #include <stdio.h>
3
4 extern double predict(double ***summaries, double *test_row, int class_num, int
  *class_num_list, int row, int col);
5 extern int get_class_num(double **dataset, int row, int col);
6 extern int *get_class_num_list(double **dataset, int class_num, int row, int col);
7 extern double ***summarize_by_class(double **train, int class_num, int
  *class_num_list, int row, int col);
8
9 double *get_test_prediction(double **train, int train_size, double **test, int
  test_size, int col)
10 {
11     int class_num = get_class_num(train, train_size, col);
12     int *class_num_list = get_class_num_list(train, class_num, train_size, col);
13     double *predictions = (double *)malloc(test_size * sizeof(double)); //因为
  test_size和fold_size一样大
14     double ***summaries = summarize_by_class(train, class_num, class_num_list,
  train_size, col);

```

```

15     for (int i = 0; i < test_size; i++)
16     {
17         predictions[i] = predict(summaries, test[i], class_num, class_num_list,
18         train_size, col);
19     }
20     return predictions; //返回对test的预测数组
21 }

```

## 6) evaluate.c

```

1  #include <stdlib.h>
2  #include <stdio.h>
3
4  extern double *get_test_prediction(double **train, int train_size, double **test, int
5  test_size, int col);
6  extern double accuracy_metric(double *actual, double *predicted, int fold_size);
7  extern double ***cross_validation_split(double **dataset, int row, int col, int
8  n_folds, int fold_size);
9
10 void evaluate_algorithm(double **dataset, int row, int col, int n_folds)
11 {
12     int fold_size = (int)row / n_folds;
13     double ***split = cross_validation_split(dataset, row, n_folds, fold_size, col);
14     int i, j, k, l;
15     int test_size = fold_size;
16     int train_size = fold_size * (n_folds - 1);
17     double *score = (double *)malloc(n_folds * sizeof(double));
18
19     for (i = 0; i < n_folds; i++)
20     {
21         double ***split_copy = (double ***)malloc(n_folds * sizeof(double **));
22         for (j = 0; j < n_folds; j++)
23         {
24             split_copy[j] = (double **)malloc(fold_size * sizeof(double *));
25             for (k = 0; k < fold_size; k++)
26             {
27                 split_copy[j][k] = (double *)malloc(col * sizeof(double));
28             }
29         }
30         for (j = 0; j < n_folds; j++)
31         {
32             for (k = 0; k < fold_size; k++)
33             {
34                 for (l = 0; l < col; l++)
35                 {
36                     split_copy[j][k][l] = split[j][k][l];
37                 }
38             }
39         }
40         double **test_set = (double **)malloc(test_size * sizeof(double *));
41         for (j = 0; j < test_size; j++)
42         {
43             test_set[j] = (double *)malloc(col * sizeof(double));
44             for (k = 0; k < col; k++)
45             {
46                 test_set[j][k] = split_copy[i][j][k];
47             }
48         }
49         for (j = i; j < n_folds - 1; j++)
50         {
51             split_copy[j] = split_copy[j + 1];
52         }
53     }
54 }

```

```

51     double **train_set = (double **)malloc(train_size * sizeof(double *));
52     for (k = 0; k < n_folds - 1; k++)
53     {
54         for (l = 0; l < fold_size; l++)
55         {
56             train_set[k * fold_size + l] = (double *)malloc(col *
sizeof(double));
57             train_set[k * fold_size + l] = split_copy[k][l];
58         }
59     }
60     double *predicted = (double *)malloc(test_size * sizeof(double));
61     predicted = get_test_prediction(train_set, train_size, test_set, test_size,
col);
62     double *actual = (double *)malloc(test_size * sizeof(double));
63     for (l = 0; l < test_size; l++)
64     {
65         actual[l] = test_set[l][col - 1];
66     }
67
68     double acc = accuracy_metric(actual, predicted, test_size);
69     score[i] = acc;
70     printf("Scores[%d] = %f%%\n", i, score[i]);
71     free(split_copy);
72 }
73 double total = 0;
74 for (l = 0; l < n_folds; l++)
75 {
76     total += score[l];
77 }
78 printf("mean_accuracy = %f%%\n", total / n_folds);
79 }

```

## 7) main.c

```

1  #include <stdlib.h>
2  #include <string.h>
3  #include <stdio.h>
4  #include <math.h>
5
6  extern int get_row(char *filename);
7  extern int get_col(char *filename);
8  extern void get_two_dimension(char *line, double **dataset, char *filename);
9  extern void evaluate_algorithm(double **dataset, int row, int col, int n_folds);
10
11 void quicksort(double *arr, int L, int R)
12 {
13     int i = L;
14     int j = R;
15     //支点
16     int kk = (L + R) / 2;
17     double pivot = arr[kk];
18     //左右两端进行扫描，只要两端还没有交替，就一直扫描
19     while (i <= j)
20     { //寻找直到比支点大的数
21         while (pivot > arr[i])
22         {
23             i++;
24         } //寻找直到比支点小的数
25         while (pivot < arr[j])
26         {
27             j--;
28         } //此时已经分别找到了比支点小的数(右边)、比支点大的数(左边)，它们进行交换

```

```

29     if (i <= j)
30     {
31         double temp = arr[i];
32         arr[i] = arr[j];
33         arr[j] = temp;
34         i++;
35         j--;
36     }
37 } //上面一个while保证了第一趟排序支点的左边比支点小，支点的右边比支点大了。
38 //“左边”再做排序，直到左边剩下一个数(递归出口)
39 if (L < j)
40 {
41     quicksort(arr, L, j);
42 }
43 //“右边”再做排序，直到右边剩下一个数(递归出口)
44 if (i < R)
45 {
46     quicksort(arr, i, R);
47 }
48 }
49 double get_mean(double **dataset, int row, int col)
50 {
51     int i;
52     double mean = 0;
53     for (i = 0; i < row; i++)
54     {
55         mean += dataset[i][col];
56     }
57     return mean / row;
58 }
59 double get_std(double **dataset, int row, int col)
60 {
61     int i;
62     double mean = 0;
63     double std = 0;
64     for (i = 0; i < row; i++)
65     {
66         mean += dataset[i][col];
67     }
68     mean /= row;
69     for (i = 0; i < row; i++)
70     {
71         std += pow((dataset[i][col] - mean), 2);
72     }
73     return sqrt(std / (row - 1));
74 }
75
76 int get_class_num(double **dataset, int row, int col)
77 {
78     int i;
79     int num = 1;
80     double *class_data = (double *)malloc(row * sizeof(double));
81     for (i = 0; i < row; i++)
82     {
83         class_data[i] = dataset[i][col - 1];
84     }
85     quicksort(class_data, 0, row - 1);
86     for (i = 0; i < row - 1; i++)
87     {
88         if (class_data[i] != class_data[i + 1])
89         {
90             num += 1;
91         }

```



```

92     }
93     return num;
94 }
95 int *get_class_num_list(double **dataset, int class_num, int row, int col)
96 {
97     int i, j;
98     int *class_num_list = (int *)malloc(class_num * sizeof(int));
99     for (j = 0; j < class_num; j++)
100     {
101         class_num_list[j] = 0;
102     }
103     for (j = 0; j < class_num; j++)
104     {
105         for (i = 0; i < row; i++)
106         {
107             if (dataset[i][col - 1] == j)
108             {
109                 class_num_list[j] += 1;
110             }
111         }
112     }
113     return class_num_list;
114 }
115
116 double ***separate_by_class(double **dataset, int class_num, int *class_num_list,
117 int row, int col)
118 {
119     double ***separated;
120     separated = (double ***)malloc(class_num * sizeof(double **));
121     int i, j;
122     for (i = 0; i < class_num; i++)
123     {
124         separated[i] = (double **)malloc(class_num_list[i] * sizeof(double *));
125         for (j = 0; j < class_num_list[i]; j++)
126         {
127             separated[i][j] = (double *)malloc(col * sizeof(double));
128         }
129     }
130     int *index = (int *)malloc(class_num * sizeof(int));
131     for (i = 0; i < class_num; i++)
132     {
133         index[i] = 0;
134     }
135     for (i = 0; i < row; i++)
136     {
137         for (j = 0; j < class_num; j++)
138         {
139             if (dataset[i][col - 1] == j)
140             {
141                 separated[j][index[j]] = dataset[i];
142                 index[j]++;
143             }
144         }
145     }
146     return separated;
147 }
148 double **summarize_dataset(double **dataset, int row, int col)
149 {
150     int i;
151     double **summary = (double **)malloc((col - 1) * sizeof(double *));
152     for (i = 0; i < (col - 1); i++)
153     {
154         summary[i] = (double *)malloc(2 * sizeof(double));
155     }

```

```

154     summary[i][0] = get_mean(dataset, row, i);
155     summary[i][1] = get_std(dataset, row, i);
156 }
157 return summary;
158 }
159 double ***summarize_by_class(double **train, int class_num, int *class_num_list, int
row, int col)
160 {
161     int i;
162     double ***summarize;
163     summarize = (double ***)malloc(class_num * sizeof(double **));
164     double ***separate = separate_by_class(train, class_num, class_num_list, row,
col);
165     for (i = 0; i < class_num; i++)
166     {
167         summarize[i] = summarize_dataset(separate[i], class_num_list[i], col);
168     }
169     return summarize;
170 }
171
172 double calculate_probability(double x, double mean, double std)
173 {
174     double pi = acos(-1.0);
175     double p = 1 / (pow(2 * pi, 0.5) * std) *
176         exp(-(pow((x - mean), 2) / (2 * pow(std, 2))));
177     return p;
178 }
179 double *calculate_class_probabilities(double ***summaries, double *test_row, int
class_num, int *class_num_list, int row, int col)
180 {
181     int i, j;
182     double *probabilities = (double *)malloc(class_num * sizeof(double));
183     for (i = 0; i < class_num; i++)
184     {
185         probabilities[i] = (double)class_num_list[i] / row;
186     }
187     for (i = 0; i < class_num; i++)
188     {
189         for (j = 0; j < col - 1; j++)
190         {
191             probabilities[i] *= calculate_probability(test_row[j], summaries[i][j]
[0], summaries[i][j][1]);
192         }
193     }
194     return probabilities;
195 }
196 double predict(double ***summaries, double *test_row, int class_num, int
*class_num_list, int row, int col)
197 {
198     int i;
199     double *probabilities = calculate_class_probabilities(summaries, test_row,
class_num, class_num_list, row, col);
200     double label = 0;
201     double best_prob = probabilities[0];
202     for (i = 1; i < class_num; i++)
203     {
204         if (probabilities[i] > best_prob)
205         {
206             label = i;
207             best_prob = probabilities[i];
208         }
209     }
210     return label;

```

```

211 }
212
213 void main()
214 {
215     char filename[] = "iris.csv";
216     char line[1024];
217     int row = get_row(filename);
218     int col = get_col(filename);
219     //printf("row = %d, col = %d\n", row, col);
220     double **dataset;
221     dataset = (double **)malloc(row * sizeof(double *));
222     for (int i = 0; i < row; ++i)
223     {
224         dataset[i] = (double *)malloc(col * sizeof(double));
225     }
226     get_two_dimension(line, dataset, filename);
227     int n_folds = 5;
228     evaluate_algorithm(dataset, row, col, n_folds);
229 }

```

## 8) compile.sh

```

1 gcc main.c read_csv.c normalize.c k_fold.c evaluate.c score.c test_prediction.c -o run
  -lm && ./run

```

### 编译&运行:

```

1 bash compile.sh

```

最终输出结果如下:

```

1 Scores[0] = 96.666667%
2 Scores[1] = 93.333333%
3 Scores[2] = 96.666667%
4 Scores[3] = 100.000000%
5 Scores[4] = 93.333333%
6 mean_accuracy = 96.000000%

```

## Python语言实战

本节同样假设您已经下载数据集，我们使用著名机器学习开源库sklearn高效实现朴素贝叶斯算法，以便您在实战中使用该算法：

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import KFold
4 from sklearn.naive_bayes import GaussianNB
5 from sklearn.metrics import accuracy_score
6 from sklearn.preprocessing import MinMaxScaler
7
8
9 if __name__ == '__main__':
10     dataset = np.array(pd.read_csv("iris.csv", sep=',', header=None))
11     k_Cross = KFold(n_splits=5, random_state=0, shuffle=True)
12     index = 0
13     score = np.array([])
14     scaler = MinMaxScaler()
15     data, label = dataset[:, :-1], dataset[:, -1]
16     data = scaler.fit_transform(data)
17     for train_index, test_index in k_Cross.split(dataset):

```

```

18 train_data, train_label = data[train_index, :], label[train_index]
19 test_data, test_label = data[test_index, :], label[test_index]
20 model = GaussianNB()
21 model.fit(train_data, train_label)
22 pred = model.predict(test_data)
23 acc = accuracy_score(test_label, pred)
24 score = np.append(score, acc)
25 print('score[{}] = {}'.format(index, acc))
26 index+=1
27 print('mean_accuracy = {}'.format(np.mean(score)))

```

输出结果如下：

```

1 score[0] = 0.9666666666666667%
2 score[1] = 0.9%
3 score[2] = 0.9666666666666667%
4 score[3] = 1.0%
5 score[4] = 0.9333333333333333%
6 mean_accuracy = 0.9533333333333334%

```

## 3.7 $k$ -Nearest Neighbors

$k$ -近邻 ( $k$ -Nearest Neighbour) , 简称KNN。KNN算法最初由Cover和Hart于1968年提出, 是一个理论上比较成熟的方法, 也是最简单的机器学习算法之一。

### 3.7.1 算法介绍

KNN算法是一种有监督学习。KNN算法的核心就是从训练集中选取  $k$  个与新样本**相似度最高**的样本 (  $k$  个近邻) , 通过这  $k$  个近邻的类别来确定待新样本的类别。其中,  $k$  的大小是可以自由选取的。

如何衡量样本之间的相似度呢? 下面, 引入欧式距离公式:

我们知道, 两点  $(x_0, y_0), (x_1, y_1)$  之间的欧几里得距离公式如下:

$$d = \sqrt{(x_0 - x_1)^2 + (y_0 - y_1)^2}$$

然而, 多数用于机器学习的样本可能是高维的。因此, 我们将其推广至  $n$  维, 设两点分别是  $(x_1^{(1)}, x_2^{(1)}, \dots, x_n^{(1)})$  与  $(x_1^{(2)}, x_2^{(2)}, \dots, x_n^{(2)})$ , 则两点之间的欧几里得距离公式为:

$$d = \sqrt{\sum_{i=1}^n (x_i^{(1)} - x_i^{(2)})^2}$$

若训练集共  $m$  个样本, 我们得到每一个训练样本与新样本的距离序列  $\{d_i\} (i = 1, 2, \dots, m)$ 。

**我们认为两个样本的距离最小时, 它们最相似。**因此我们从序列  $\{d_i\}$  中选取最小的  $k$  个样本, 设它们的类别分别是  $y_1, y_2, \dots, y_k$ , 我们求出这  $k$  个数的众数  $l$ ,  $l$  即为KNN算法对新样本的分类。

**当然, KNN算法不仅可用于分类任务, 也可以用于回归任务。当KNN算法用于回归任务时, 我们可以求出新样本的  $k$  个近邻类别的均值  $m$ ,  $m$  即为KNN算法对新样本的预测。**

### 3.7.2 算法讲解

#### 欧几里得距离

实现KNN算法的第一步是计算同一份数据集中的任意两行数据的距离。

下面给出一个C语言自定义函数 `euclidean_distance()`, 完成欧几里得距离的计算。

```

1 #include<stdio.h>
2 #include<math.h>
3

```

```

4 double euclidean_distance(double *row1, double *row2, int col) {
5     double distance = 0;
6     for (int i = 0; i < col - 1; i++) {
7         distance += pow((row1[i] - row2[i]), 2);
8     }
9     return sqrt(distance);
10 }
11
12 void main(){
13     double test_data[10][3] = {
14         {2.56373457, 2.63727045, 0},
15         {1.62548536, 2.26342507, 0},
16         {3.69634668, 4.34629352, 0},
17         {1.45607019, 1.84562031, 0},
18         {3.06407232, 3.00530597, 0},
19         {7.54753121, 2.98926223, 1},
20         {5.12422124, 2.08862677, 1},
21         {6.86549671, 1.77106367, 1},
22         {8.67541865, -0.24206865, 1},
23         {7.67375646, 3.76356301, 1}
24     };
25     double** dataset;
26     dataset = (double**)malloc(10 * sizeof(double*));
27     for (int i = 0; i < 10; ++i) {
28         dataset[i] = test_data[i];
29     }
30     double result;
31     for (int i = 0; i < 10; i++) {
32         result = euclidean_distance(test_data[0], test_data[i], 3);
33         printf("%f\n", result);
34     }
35 }

```

运行代码，得到每一行数据与第一行数据的距离：

```

1 0.000000
2 1.009986
3 2.050261
4 1.361481
5 0.621118
6 4.996211
7 2.618607
8 4.388106
9 6.755981
10 5.232672

```

## 获取近邻

通过计算某一条数据 $x_i$ 与其他所有数据之间的距离，得到一个距离集合。将该距离集合中的元素由小到大排序，找到与 $x_i$ 最近的 $k$ 个数据，即 $k$ 个近邻，返回他们的标签。

下面以3个近邻为例，沿用上面的数据，求出与 $x_0$ 距离最近的3条数据的标签。

```

1 void QuickSort(double **arr, int L, int R) {
2     int i = L;
3     int j = R;
4     int kk = (L + R) / 2; //支点
5     double pivot = arr[kk][0];
6     //左右两端进行扫描，直到两端交替
7     while (i <= j) {
8         //寻找比支点大的数
9         while (pivot > arr[i][0])

```

```

10     {
11         i++;
12     } //寻找比支点小的数
13     while (pivot < arr[j][0])
14     {
15         j--;
16     } //此时已经分别找到了比支点小的数(右边)、比支点大的数(左边), 交换他们
17     if (i <= j) {
18         double *temp = arr[i];
19         arr[i] = arr[j];
20         arr[j] = temp;
21         i++; j--;
22     }
23 } //上面的while保证了第一次排序支点的左边比支点小, 支点的右边比支点大了。
24 //左边再做排序, 直到左边剩下一个数(递归出口)
25 if (L < j)
26 {
27     QuickSort(arr, L, j);
28 }
29 //“右边”再做排序, 直到右边剩下一个数(递归出口)
30 if (i < R)
31 {
32     QuickSort(arr, i, R);
33 }
34 }
35
36 double* get_neighbors(double **train_data, int train_row, int col, double *test_row,
37 int num_neighbors) {
38     double *neighbors = (double *)malloc(num_neighbors * sizeof(double));
39     double **distances = (double **)malloc(train_row * sizeof(double *));
40     for (int i = 0; i < train_row; i++) {
41         distances[i] = (double *)malloc(2 * sizeof(double));
42         distances[i][0] = euclidean_distance(train_data[i], test_row, col);
43         distances[i][1] = train_data[i][col - 1];
44     }
45     QuickSort(distances, 0, train_row - 1);
46     for (int i = 0; i < num_neighbors; i++) {
47         neighbors[i] = distances[i][1];
48     }
49     return neighbors;
50 }
51
52 void main(){
53     double result;
54     for (int i = 0; i < 10; i++) {
55         result = euclidean_distance(dataset[0], dataset[i], 3);
56         printf("%f\n", result);
57     }
58     int num_neighbors = 3;
59     double* neighbors = get_neighbors(dataset, 10, 3, dataset[0], num_neighbors);
60     for (int i = 0; i < num_neighbors; i++) {
61         printf("%f\n", neighbors[i]);
62     }
63 }

```

得到近邻的标签为:

```

1  0.000000
2  0.000000
3  0.000000

```

## 预测结果

```
1 double predict(double **train_data, int train_row, int col, double *test_row, int
  num_neighbors) {
2     double* neighbors = get_neighbors(train_data, train_row, col, test_row,
  num_neighbors);
3     double result = 0;
4     for (int i = 0; i < num_neighbors; i++) {
5         result += neighbors[i];
6     }
7     return result / num_neighbors;
8 }
9
10 void main() {
11     double test_data[10][3] = {
12         {2.56373457, 2.63727045, 0},
13         {1.62548536, 2.26342507, 0},
14         {3.69634668, 4.34629352, 0},
15         {1.45607019, 1.84562031, 0},
16         {3.06407232, 3.00530597, 0},
17         {7.54753121, 2.98926223, 1},
18         {5.12422124, 2.08862677, 1},
19         {6.86549671, 1.77106367, 1},
20         {8.67541865, -0.24206865, 1},
21         {7.67375646, 3.76356301, 1}
22     };
23     double** dataset;
24     dataset = (double**)malloc(10 * sizeof(double*));
25     for (int i = 0; i < 10; ++i) {
26         dataset[i] = test_data[i];
27     }
28     int num_neighbors = 3;
29     double result;
30     result = predict(dataset, 10, 3, dataset[0], num_neighbors);
31     printf("%f\n", result);
32 }
```

该数据的标签为0，而结果为：

1 | 0

### 3.7.3 算法代码

我们现在知道了如何实现**KNN算法**，那么我们把它应用到[鲍鱼数据集 abalone.csv](https://aistudio.baidu.com/aistudio/datasetdetail/105756/0)

我们给出链接：<https://aistudio.baidu.com/aistudio/datasetdetail/105756/0>

### C语言细节讲解

本节假设您已下载数据集 `abalone.csv`，并且它在当前工作目录中可用。下面我们给出一个完整实例，使用C语言详细讲解每一处细节。我们给出每一个.c文件的所有代码：

#### 1) read\_csv.c

该文件代码与前面代码一致，不再重复给出。

#### 1) normalize.c

该文件代码与前面代码一致，不再重复给出。

## 2) k\_fold.c

该文件代码与前面代码一致，不再重复给出。

## 3) rmse.c

该文件代码与前面代码一致，不再重复给出。

## 4) test\_prediction.c

```
1  #include<stdlib.h>
2  #include<stdio.h>
3  #include<math.h>
4
5  extern void Quicksort(double **arr, int L, int R);
6  extern double euclidean_distance(double *row1, double *row2, int col);
7  extern double* get_neighbors(double **train_data, int train_row, int col, double
    *test_row, int num_neighbors);
8  extern double predict(double **train_data, int train_row, int col, double *test_row,
    int num_neighbors);
9
10 double* get_test_prediction(double **train, int train_size, double **test, int
    test_size, int col, int num_neighbors)
11 {
12     double* predictions = (double*)malloc(test_size * sizeof(double));
13     for (int i = 0; i < test_size; i++)
14     {
15         predictions[i] = predict(train, train_size,col,test[i],num_neighbors);
16     }
17     return predictions;//返回对test的预测数组
18 }
```

## 5) evaluate.c

```
1  #include <stdlib.h>
2  #include <stdio.h>
3
4  extern double *get_test_prediction(double **train, int train_size, double **test, int
    test_size, int col, int num_neighbors);
5  extern double rmse_metric(double *actual, double *predicted, int fold_size);
6  extern double ***cross_validation_split(double **dataset, int row, int col, int
    n_folds, int fold_size);
7
8  void evaluate_algorithm(double **dataset, int row, int col, int n_folds, int
    num_neighbors)
9  {
10     int fold_size = (int)row / n_folds;
11     double ***split = cross_validation_split(dataset, row, n_folds, fold_size, col);
12     int i, j, k, l;
13     int test_size = fold_size;
14     int train_size = fold_size * (n_folds - 1);
15     double *score = (double *)malloc(n_folds * sizeof(double));
16
17     for (i = 0; i < n_folds; i++)
18     {
19         double ***split_copy = (double ***)malloc(n_folds * sizeof(double **));
20         for (j = 0; j < n_folds; j++)
21         {
22             split_copy[j] = (double **)malloc(fold_size * sizeof(double **));
23             for (k = 0; k < fold_size; k++)
24             {
25                 split_copy[j][k] = (double *)malloc(col * sizeof(double));
26             }
27         }
28     }
```



```

27     }
28     for (j = 0; j < n_folds; j++)
29     {
30         for (k = 0; k < fold_size; k++)
31         {
32             for (l = 0; l < col; l++)
33             {
34                 split_copy[j][k][l] = split[j][k][l];
35             }
36         }
37     }
38     double **test_set = (double **)malloc(test_size * sizeof(double *));
39     for (j = 0; j < test_size; j++)
40     {
41         test_set[j] = (double *)malloc(col * sizeof(double));
42         for (k = 0; k < col; k++)
43         {
44             test_set[j][k] = split_copy[i][j][k];
45         }
46     }
47     for (j = i; j < n_folds - 1; j++)
48     {
49         split_copy[j] = split_copy[j + 1];
50     }
51     double **train_set = (double **)malloc(train_size * sizeof(double *));
52     for (k = 0; k < n_folds - 1; k++)
53     {
54         for (l = 0; l < fold_size; l++)
55         {
56             train_set[k * fold_size + l] = (double *)malloc(col *
sizeof(double));
57             train_set[k * fold_size + l] = split_copy[k][l];
58         }
59     }
60     double *predicted = (double *)malloc(test_size * sizeof(double));
61     predicted = get_test_prediction(train_set, train_size, test_set, test_size,
col, num_neighbors);
62     double *actual = (double *)malloc(test_size * sizeof(double));
63     for (l = 0; l < test_size; l++)
64     {
65         actual[l] = test_set[l][col - 1];
66     }
67     double rmse = rmse_metric(actual, predicted, test_size);
68     score[i] = rmse;
69     printf("rmse[%d]=%f\n", i, score[i]);
70     free(split_copy);
71 }
72 double total = 0;
73 for (l = 0; l < n_folds; l++)
74 {
75     total += score[l];
76 }
77 printf("mean_rmse=%f\n", total / n_folds);
78 }

```

## 6) main.c

```

1  #define _CRT_SECURE_NO_WARNINGS
2  #include <stdlib.h>
3  #include <string.h>
4  #include <stdio.h>
5  #include <math.h>

```

```

6
7 extern int get_row(char *filename);
8 extern int get_col(char *filename);
9 extern void get_two_dimension(char *line, double **dataset, char *filename);
10 extern void normalize_dataset(double **dataset, int row, int col);
11 extern void evaluate_algorithm(double **dataset, int row, int col, int n_folds, int
    num_neighbors);
12
13 void QuickSort(double **arr, int L, int R)
14 {
15     int i = L;
16     int j = R;
17     //支点
18     int kk = (L + R) / 2;
19     double pivot = arr[kk][0];
20     //左右两端进行扫描，只要两端还没有交替，就一直扫描
21     while (i <= j)
22     { //寻找直到比支点大的数
23         while (pivot > arr[i][0])
24         {
25             i++;
26         } //寻找直到比支点小的数
27         while (pivot < arr[j][0])
28         {
29             j--;
30         } //此时已经分别找到了比支点小的数(右边)、比支点大的数(左边)，它们进行交换
31         if (i <= j)
32         {
33             double *temp = arr[i];
34             arr[i] = arr[j];
35             arr[j] = temp;
36             i++;
37             j--;
38         }
39     } //上面一个while保证了第一趟排序支点的左边比支点小，支点的右边比支点大了。
40     //“左边”再做排序，直到左边剩下一个数(递归出口)
41     if (L < j)
42     {
43         QuickSort(arr, L, j);
44     }
45     //“右边”再做排序，直到右边剩下一个数(递归出口)
46     if (i < R)
47     {
48         QuickSort(arr, i, R);
49     }
50 }
51 // Calculate the Euclidean distance between two vectors
52 double euclidean_distance(double *row1, double *row2, int col)
53 {
54     double distance = 0;
55     for (int i = 0; i < col - 1; i++)
56     {
57         distance += pow((row1[i] - row2[i]), 2);
58     }
59     return sqrt(distance);
60 }
61 // Locate the most similar neighbors
62 double *get_neighbors(double **train_data, int train_row, int col, double *test_row,
    int num_neighbors)
63 {
64     double *neighbors = (double *)malloc(num_neighbors * sizeof(double));
65     double **distances = (double **)malloc(train_row * sizeof(double *));
66     for (int i = 0; i < train_row; i++)

```

```

67     {
68         distances[i] = (double *)malloc(2 * sizeof(double));
69         distances[i][0] = euclidean_distance(train_data[i], test_row, col);
70         distances[i][1] = train_data[i][col - 1];
71     }
72     quickSort(distances, 0, train_row - 1);
73     for (int i = 0; i < num_neighbors; i++)
74     {
75         neighbors[i] = distances[i][1];
76     }
77     return neighbors;
78 }
79 double predict(double **train_data, int train_row, int col, double *test_row, int
num_neighbors)
80 {
81     double *neighbors = get_neighbors(train_data, train_row, col, test_row,
num_neighbors);
82     double result = 0;
83     for (int i = 0; i < num_neighbors; i++)
84     {
85         result += neighbors[i];
86     }
87     return result / num_neighbors;
88 }
89
90 void main()
91 {
92     char filename[] = "abalone.csv";
93     char line[1024];
94     int row = get_row(filename);
95     int col = get_col(filename);
96     //printf("row = %d, col = %d\n", row, col);
97     double **dataset;
98     dataset = (double **)malloc(row * sizeof(double *));
99     for (int i = 0; i < row; ++i)
100     {
101         dataset[i] = (double *)malloc(col * sizeof(double));
102     }
103     get_two_dimension(line, dataset, filename);
104     normalize_dataset(dataset, row, col);
105     int k_fold = 5;
106     int num_neighbors = 5;
107     evaluate_algorithm(dataset, row, col, k_fold, num_neighbors);
108 }

```

## 7) compile.sh

```

1 | gcc main.c read_csv.c normalize.c k_fold.c evaluate.c rmse.c test_prediction.c -o run
   | -lm && ./run

```

### 编译&运行:

```

1 | bash compile.sh

```

最终输出结果如下:

```

1 rmse[0] = 0.081334
2 rmse[1] = 0.083535
3 rmse[2] = 0.080164
4 rmse[3] = 0.081941
5 rmse[4] = 0.079612
6 mean_rmse = 0.081317

```

## Python语言实战

本节同样假设您已经下载数据集，我们使用著名机器学习开源库sklearn高效实现**KNN算法**，以便您在实战中使用该算法：

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import KFold
4 from sklearn.neighbors import KNeighborsClassifier
5 from sklearn.preprocessing import MinMaxScaler
6
7
8 def rmse_metric(actual, predicted):
9     sum_err = 0.0
10    for i in range(len(actual)):
11        err = predicted[i] - actual[i]
12        sum_err += err ** 2
13    mean_err = sum_err / (len(actual)-1)
14    return np.sqrt(mean_err)
15
16
17 if __name__ == '__main__':
18     dataset = np.array(pd.read_csv("abalone.csv", sep=',', header=None))
19     k_Cross = KFold(n_splits=5, random_state=0, shuffle=True)
20     index = 0
21     scores = np.array([])
22     scaler = MinMaxScaler()
23     data,label = dataset[:, :-1], dataset[:, -1]
24     data = scaler.fit_transform(data)
25     for train_index, test_index in k_Cross.split(dataset):
26         train_data, train_label = data[train_index, :], label[train_index]
27         test_data, test_label = data[test_index, :], label[test_index]
28         model = KNeighborsClassifier(n_neighbors=5)
29         model.fit(train_data, train_label)
30         pred = model.predict(test_data)
31         score = rmse_metric(test_label, pred)
32         scores = np.append(scores, score)
33         print('score[{}] = {}'.format(index, score))
34         index+=1
35     print('mean_rmse = {}'.format(np.mean(scores)))

```

输出结果如下，读者可以尝试分析一下为何结果会存在差异？

```

1 score[0] = 2.852249873576536
2 score[1] = 2.764618196506543
3 score[2] = 2.6721316902024443
4 score[3] = 2.6719073204392374
5 score[4] = 2.9138304126803454
6 mean_rmse = 2.7749474986810214

```

## 3.8 Learning Vector Quantization

## 3.8.1 算法简介

学习向量量化 (Learning Vector Quantization) 与K-Mean算法类似, 其为试图找到一组原型向量来刻画聚类结构, 但与一般的聚类算法不同的是, LVQ假设数据样本带有类别标记, 学习过程利用样本的这些监督信息来辅助聚类, 从而克服自组织网络采用无监督学习算法带来的缺乏分类信息的弱点。

向量量化的思路是, 将高维输入空间分成若干个不同的区域, 对每个区域确定一个中心向量作为聚类中心, 与其处于同一个区域的输入向量可用作该中心向量来代表, 从而形成了以各中心向量为聚类中心的点集。

### LVQ网络结构与工作原理

其结构分为输入层、竞争层、输出层, 竞争层和输出层之间完全连接。输出层每个神经元只与竞争层中的一组神经元连接, 连接权重固定为1, 训练过程中输入层和竞争层之间的权值逐渐被调整为聚类中心。当一个样本输入LVQ网络时, 竞争层的神经元通过胜者为王学习规则产生获胜神经元, 容许其输出为1, 其它神经元输出为0。与获胜神经元所在组相连的输出神经元输出为1, 而其它输出神经元为0, 从而给出当前输入样本的模式类。将竞争层学习得到的类成为子类, 而将输出层学习得到的类成为目标类。

### LVQ网络学习算法

LVQ网络的学习规则结合了竞争学习规则和有导师学习规则, 所以样本集应当为 $\{(x_i, d_i)\}$ 。其中 $d_i$ 为 $l$ 维, 对应输出层的 $l$ 个神经元, 它只有一个分量为1, 其他分量均为0。通常把竞争层的每个神经元指定给一个输出神经元, 相应的权值为1, 从而得到输出层的权值。比如某LVQ网络竞争层6个神经元, 输出层3个神经元, 代表3类。若将竞争层的1, 3指定为第一个输出神经元, 2, 5指定为第二个输出神经元, 3, 6指定为第三个输出神经元。

训练前预先定义好竞争层到输出层权重, 从而指定了输出神经元类别, 训练中不再改变。网络的学习通过改变输入层到竞争层的权重来进行。根据输入样本类别和获胜神经元所属类别, 可判断当前分类是否正确。若分类正确, 则将获胜神经元的权向量向输入向量方向调整, 分类错误则向相反方向调整。

## 3.8.2 算法讲解

### 算法流程

输入: 样本集 $D = (x_1, y_1), (x_2, y_2) \dots (x_m, y_m)$ ; 原型向量个数为 $q$ , 各原型向量预设的类别标记 $t_1, t_2 \dots t_q$ , 学习率 $\delta \in (0, 1)$

1. 初始化一些原型向量 $p_1, p_2 \dots p_q$
2. repeat
3. 从样本集 $D$ 随机选取样本 $(x_j, y_j)$
4. 计算样本 $x_j$ 与 $p_i (1 < i < q)$
5. 找出与 $x_j$ 距离最近的原型向量 $p_{i^*}, i^* = \operatorname{argmin}_{i \in \{1, 2, \dots, q\}} d_{ji}$
6. if  $y_j = t_{i^*}$ , then
7.  $p' = p_{i^*} + \delta(x_j - p_{i^*})$
8. else
9.  $p' = p_{i^*} - \delta(x_j - p_{i^*})$
10. end if
11. 将原型向量 $p_{i^*}$ 更新为 $p'$
12. until 满足停止条件

输出: 原型向量 $p_1, p_2, \dots, p_q$

## 核心理想

1.对原型向量进行迭代优化，每一轮随机选择一个有标记的训练样本，找出与其距离最近的原型向量，根据两者的类别标记是否一致来对原型向量进行相应的更新。

2.LVQ的关键在于第6-10行如何更新原型向量，对于样本

$x_j$ ，若最近的原型向量 $p_i$ 与 $x_j$ 的类别标记相同，则令 $p_i$ 向 $x_j$ 方向靠近，否则远离其方向，学习率为 $\delta$

## 计算欧式距离

```
1 double euclidean_distance(double*row1, double*row2){
2     int i;
3     double distance = 0.0;
4     for (i=0;i<col-1;i++){
5         distance =distance+ (row1[i] - row2[i])*(row1[i] - row2[i]);
6     }
7     return sqrt(distance);
8     //其返回的是两个标志的欧氏距离的绝对值
9 }
```

```
1 input:
2 row1:  2  4
3 row2:  1  3
4 output: 1
```

## 2.4 确定最佳匹配位置

```
1 int get_best_matching_unit(double**codebooks, double*test_row,int n_codebooks){
2     double dist_min,dist;
3     int i,min=0;
4     dist_min = euclidean_distance(codebooks[0], test_row);
5     for (i=0;i< n_codebooks;i++){
6         dist = euclidean_distance(codebooks[i], test_row);
7         if(dist < dist_min){
8             dist_min=dist;
9             min=i;
10        }
11    }
12    //bmu=codebooks[min];
13    return min;
14 }//其返回欧氏距离最小的标志的index
```

```
1 input:
2 3  6  7  2
3 output:
4 3
```

## 2.5 初始化原型向量

```
1 double** random_codebook(double**train,int n_codebooks,int fold_size){
2     int i,j,r;
3     int n_folds=(int)(row/fold_size);
4     double **codebooks=(double **)malloc(n_codebooks * sizeof(int*));
5     for ( i=0;i < n_codebooks; ++i){
6         codebooks[i] = (double *)malloc(col * sizeof(double));
7     };
8     srand((unsigned)time(NULL));
9     for(i=0;i<n_codebooks;i++){
10        for(j=0;j<col;j++){
```

```

11         //srand((unsigned int)time(0));
12         r=rand()%((n_folds-1)*fold_size);
13         //printf(" r%d",r);
14         codebooks[i][j]=train[r][j];
15     }
16 }
17 return codebooks;
18 }//产生初始化原型向量

```

```

1 output:
2 0.001251
3 0.563585
4 0.193304
5 0.808741
6 0.585009
7 0.479873
8 0.350291
9 0.895962
10 0.822840
11 0.746605
12 0.174108
13 0.858943
14 0.710501
15 0.513535
16 0.303995
17 0.014985
18 0.091403
19 0.364452
20 0.147313
21 0.165899
22 0.988525
23 0.445692
24 0.119083
25 0.004669
26 .....

```

## 预测神经网络

```

1 float* get_test_prediction(double **train, double **test, float l_rate, int n_epoch,
2 int fold_size)
3 {
4     int i;
5     double **codebooks=(double **)malloc(n_codebooks * sizeof(int*));
6     for ( i=0;i < n_codebooks; ++i){
7         codebooks[i] = (double *)malloc(col * sizeof(double));
8     };
9     float *predictions=(float*)malloc(fold_size*sizeof(float));//预测集的行数就是数组
10    prediction的长度
11
12    codebooks=train_codebooks(train,l_rate,n_epoch,n_codebooks,fold_size);
13    for(i=0;i<fold_size;i++)
14    {
15        predictions[i]=predict(codebooks,test[i]);
16    }
17    return predictions;
18 }//返回聚类预测分类结果

```

```

1 output:
2 2
3 2
4 1

```

```

5 | 1
6 | 2
7 | 1
8 | 0
9 | 0
10 | 1
11 | 2
12 | 2
13 | 2
14 | 0
15 | 2
16 | .....

```

### 3.8.3 算法代码

我们现在知道了如何实现**学习向量量化算法**，那么我们把它应用到**电离数据集 ionosphere-full.csv**

我们给出链接：<https://aistudio.baidu.com/aistudio/datasetdetail/105756/0>

### C语言细节讲解

本节假设您已下载数据集 `ionosphere-full.csv`，并且它在当前工作目录中可用。下面我们给出一个完整实例，使用C语言详细讲解每一处细节。我们给出每一个.c文件的所有代码：

#### 1) read\_csv.c

该文件代码与前面代码一致，不再重复给出。

#### 2) k\_fold.c

该文件代码与前面代码一致，不再重复给出。

#### 3) score.c

该文件代码与前面代码一致，不再重复给出。

#### 4) test\_prediction.c

```

1 | #include <stdlib.h>
2 | #include <stdio.h>
3 | #include <math.h>
4 |
5 | extern double predict(int col, double **codebooks, double *test_row, int
  n_codebooks);
6 | extern double **train_codebooks(double **train, int row, int col, double l_rate, int
  n_epoch, int n_codebooks, int fold_size);
7 |
8 | double *get_test_prediction(double **train, double **test, int row, int col, double
  l_rate, int n_epoch, int fold_size, int n_codebooks)
9 | {
10 |     int i;
11 |     double **codebooks = (double **)malloc(n_codebooks * sizeof(int *));
12 |     for (i = 0; i < n_codebooks; ++i)
13 |     {
14 |         codebooks[i] = (double *)malloc(col * sizeof(double));
15 |     };
16 |     double *predictions = (double *)malloc(fold_size * sizeof(double)); //预测集的行数
    就是数组prediction的长度
17 |     codebooks = train_codebooks(train, row, col, l_rate, n_epoch, n_codebooks,
  fold_size);
18 |     for (i = 0; i < fold_size; i++)
19 |     {
20 |         predictions[i] = predict(col, codebooks, test[i], n_codebooks);
21 |     }

```



```

22     return predictions;
23 }

```

## 5) evaluate.c

```

1  #include <stdlib.h>
2  #include <stdio.h>
3
4  extern double accuracy_metric(double *actual, double *predicted, int fold_size);
5  extern double ***cross_validation_split(double **dataset, int row, int n_folds, int
fold_size, int col);
6  extern double *get_test_prediction(double **train, double **test, int row, int col,
double l_rate, int n_epoch, int fold_size, int n_codebooks);
7
8  void evaluate_algorithm(double **dataset, int row, int col, int n_folds, int
fold_size, double l_rate, int n_epoch, int n_codebooks)
9  {
10     double ***split;
11     split = cross_validation_split(dataset, row, n_folds, fold_size, col);
12     int i, j, k, l;
13     int test_size = fold_size;
14     int train_size = fold_size * (n_folds - 1);
15     double *score = (double *)malloc(n_folds * sizeof(double));
16     for (i = 0; i < n_folds; i++)
17     {
18         double ***split_copy = (double ***)malloc(n_folds * sizeof(int **));
19         for (j = 0; j < n_folds; j++)
20         {
21             split_copy[j] = (double **)malloc(fold_size * sizeof(int *));
22             for (k = 0; k < fold_size; k++)
23             {
24                 split_copy[j][k] = (double *)malloc(col * sizeof(double));
25             }
26         }
27         for (j = 0; j < n_folds; j++)
28         {
29             for (k = 0; k < fold_size; k++)
30             {
31                 for (l = 0; l < col; l++)
32                 {
33                     split_copy[j][k][l] = split[j][k][l];
34                 }
35             }
36         }
37         double **test_set = (double **)malloc(test_size * sizeof(int *));
38         for (j = 0; j < test_size; j++)
39         {
40             test_set[j] = (double *)malloc(col * sizeof(double));
41             for (k = 0; k < col; k++)
42             {
43                 test_set[j][k] = split_copy[i][j][k];
44             }
45         }
46         for (j = i; j < n_folds - 1; j++)
47         {
48             split_copy[j] = split_copy[j + 1];
49         }
50         double **train_set = (double **)malloc(train_size * sizeof(int *));
51         for (k = 0; k < n_folds - 1; k++)
52         {
53             for (l = 0; l < fold_size; l++)
54             {

```

```

55         train_set[k * fold_size + 1] = (double *)malloc(col *
sizeof(double));
56         train_set[k * fold_size + 1] = split_copy[k][1];
57     }
58 }
59     double *predicted = (double *)malloc(test_size * sizeof(double));
60     predicted = get_test_prediction(train_set, test_set, row, col, l_rate,
n_epoch, fold_size, n_codebooks);
61     double *actual = (double *)malloc(test_size * sizeof(double));
62     for (l = 0; l < test_size; l++)
63     {
64         actual[l] = (double)test_set[l][col - 1];
65     }
66     double accuracy = accuracy_metric(actual, predicted, test_size);
67     score[i] = accuracy;
68     printf("score[%d]=%.2f%%\n", i, score[i]);
69     free(split_copy);
70 }
71 double total = 0.0;
72 for (l = 0; l < n_folds; l++)
73 {
74     total += score[l];
75 }
76 printf("mean_accuracy=%.2f%%\n", total / n_folds);
77 }

```

## 6) main.c

```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <string.h>
4  #include <malloc.h>
5  #include <time.h>
6  #include <math.h>
7
8  extern int get_row(char *filename);
9  extern int get_col(char *filename);
10 extern void get_two_dimension(char *line, double **dataset, char *filename);
11 extern void evaluate_algorithm(double **dataset, int row, int col, int n_folds, int
fold_size, double l_rate, int n_epoch, int n_codebooks);
12
13 double euclidean_distance(int col, double *row1, double *row2)
14 {
15     int i;
16     double distance = 0.0;
17     for (i = 0; i < col - 1; i++)
18     {
19         distance = distance + (row1[i] - row2[i]) * (row1[i] - row2[i]);
20     }
21     return sqrt(distance);
22 }
23
24 //Locate the best matching unit
25 int get_best_matching_unit(int col, double **codebooks, double *test_row, int
n_codebooks)
26 {
27     double dist_min, dist;
28     int i, min = 0;
29     dist_min = euclidean_distance(col, codebooks[0], test_row);
30     for (i = 0; i < n_codebooks; i++)
31     {
32         dist = euclidean_distance(col, codebooks[i], test_row);

```

```

33         if (dist < dist_min)
34         {
35             dist_min = dist;
36             min = i;
37         }
38     }
39     return min;
40 }
41
42 // Make a prediction with codebook vectors
43 double predict(int col, double **codebooks, double *test_row, int n_codebooks)
44 {
45     int min;
46     min = get_best_matching_unit(col, codebooks, test_row, n_codebooks);
47     return (double)codebooks[min][col - 1];
48 }
49
50 // Create random codebook vectors
51 double **random_codebook(double **train, int row, int col, int n_codebooks, int
fold_size)
52 {
53     int i, j, r;
54     int n_folds = (int)(row / fold_size);
55     double **codebooks = (double **)malloc(n_codebooks * sizeof(int *));
56     for (i = 0; i < n_codebooks; ++i)
57     {
58         codebooks[i] = (double *)malloc(col * sizeof(double));
59     };
60     srand((unsigned)time(NULL));
61     for (i = 0; i < n_codebooks; i++)
62     {
63         for (j = 0; j < col; j++)
64         {
65             r = rand() % ((n_folds - 1) * fold_size);
66             codebooks[i][j] = train[r][j];
67         }
68     }
69     return codebooks;
70 }
71
72 double **train_codebooks(double **train, int row, int col, double l_rate, int
n_epoch, int n_codebooks, int fold_size)
73 {
74     int i, j, k, min = 0;
75     double error, rate = 0.0;
76     int n_folds = (int)(row / fold_size);
77     double **codebooks = (double **)malloc(n_codebooks * sizeof(int *));
78     for (i = 0; i < n_codebooks; ++i)
79     {
80         codebooks[i] = (double *)malloc(col * sizeof(double));
81     };
82     codebooks = random_codebook(train, row, col, n_codebooks, fold_size);
83     for (i = 0; i < n_epoch; i++)
84     {
85         rate = l_rate * (1.0 - (i / (double)n_epoch));
86         for (j = 0; j < fold_size * (n_folds - 1); j++)
87         {
88             min = get_best_matching_unit(col, codebooks, train[j], n_codebooks);
89             for (k = 0; k < col - 1; k++)
90             {
91                 error = train[j][k] - codebooks[min][k];
92                 if (fabs(codebooks[min][col - 1] - train[j][col - 1]) < 1e-13)
93                     {

```

```

94         codebooks[min][k] = codebooks[min][k] + rate * error;
95     }
96     else
97     {
98         codebooks[min][k] = codebooks[min][k] - rate * error;
99     }
100 }
101 }
102 }
103 return codebooks;
104 }
105
106 int main()
107 {
108     char filename[] = "ionosphere-full.csv";
109     char line[1024];
110     int row = get_row(filename);
111     int col = get_col(filename);
112     int i;
113     double **dataset = (double **)malloc(row * sizeof(int *));
114     for (i = 0; i < row; ++i)
115     {
116         dataset[i] = (double *)malloc(col * sizeof(double));
117     }
118     get_two_dimension(line, dataset, filename);
119     int n_folds = 5;
120     double l_rate = 0.3;
121     int n_epoch = 50;
122     int fold_size = (int)(row / n_folds);
123     int n_codebooks = 20;
124     evaluate_algorithm(dataset, row, col, n_folds, fold_size, l_rate, n_epoch,
n_codebooks);
125     return 0;
126 }

```

## 7) compile.sh

```

1 | gcc main.c read_csv.c normalize.c k_fold.c evaluate.c score.c test_prediction.c -o run
   | -lm && ./run

```

### 编译&运行:

```

1 | bash compile.sh

```

最终输出结果如下:

```

1 | score[0] = 91.4286%
2 | score[1] = 90.0000%
3 | score[2] = 87.1429%
4 | score[3] = 81.4286%
5 | score[4] = 87.1429%
6 | mean_accuracy = 87.4286%

```

## Python语言实战

本节同样假设您已经下载数据集，我们使用著名机器学习开源库sklearn与sklearn\_lvq高效实现学习向量量化算法，以便您在实战中使用该算法：

```

1 | import pandas as pd
2 | import numpy as np

```

```

3 from sklearn.model_selection import KFold
4 from sklearn.metrics import accuracy_score
5 from sklearn_lvq import LgmlvqModel
6
7
8 if __name__ == '__main__':
9     dataset = np.array(pd.read_csv("ionosphere-full.csv", sep=',', header=None))
10    k_Cross = KFold(n_splits=5, random_state=8, shuffle=True)
11    index = 0
12    score = np.array([])
13    data, label = dataset[:, :-1], dataset[:, -1]
14    for train_index, test_index in k_Cross.split(dataset):
15        train_data, train_label = data[train_index, :], label[train_index]
16        test_data, test_label = data[test_index, :], label[test_index]
17        model = LgmlvqModel()
18        model.fit(train_data, train_label)
19        pred = model.predict(test_data)
20        acc = accuracy_score(test_label, pred)
21        score = np.append(score, acc)
22        print('score[{}] = {}'.format(index, acc))
23        index += 1
24    print('mean_accuracy = {}'.format(np.mean(score)))

```

输出结果如下：

```

1 score[0] = 0.9014084507042254%
2 score[1] = 0.8714285714285714%
3 score[2] = 0.8428571428571429%
4 score[3] = 0.9%
5 score[4] = 0.9285714285714286%
6 mean_accuracy = 0.8888531187122737%

```

## 3.9 Support Vector Machine

支持向量机 (Support Vector Machine, SVM) 是一类按监督学习方式对数据进行二元分类的广义线性分类器，其决策边界是对学习样本求解的最大边距超平面。SVM算法最早是由 Vladimir N. Vapnik 和 Alexey Ya. Chervonenkis 在1963年提出；目前的版本(soft margin)是由Corinna Cortes 和 Vapnik在1993年提出，并在1995年发表；深度学习 (2012) 出现之前，SVM被认为机器学习中近十几年来最成功，表现最好的算法。

### 3.9.1 算法简介

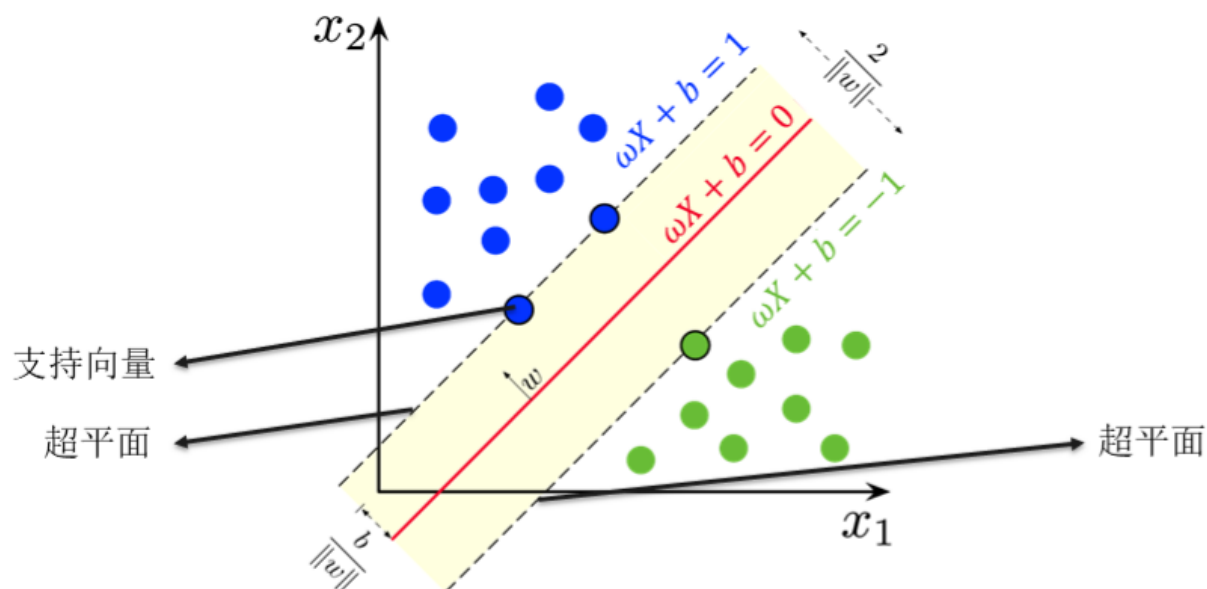
支持向量机 (SVM) 是用于分类与回归分析中分析数据的监督式学习模型与相关的学习算法。给定一组数据集，每个训练数据被标记为属于两个类别中的一个或另一个，SVM训练算法创建一个根据训练集的模型，使该模型成为非概率二元线性分类器。除了进行线性分类之外，SVM还可以使用所谓的核技巧有效地进行非线性分类，将其输入隐式映射到高维特征空间中。

本书中，我们将主要介绍数据集**线性可分**情况下的SVM算法。

### 3.9.2 算法讲解

#### 超平面

SVM算法的目标是找到一个**超平面**  $\omega X + b = 0$ ，以此超平面为边界，可以将数据集分为两类，如下图所示：



因此，我们需要先找到各个分类的样本点离这个超平面最近的点，使得这个点到超平面的距离最大化，最近的点就是在虚线上的点，这个点也被称之为**支持向量**。于是，我们认为 $\omega X - b > 1$ 为蓝色点， $\omega X - b > -1$ 为绿色点，于是完成二分类。

## 线性可分

在数据集式线性可分的情况下，数据集应分布于超平面的两侧。那么，使得支持向量到超平面的距离最大化，可以使用如下等价的两个优化式表示：

$$\begin{aligned} \max_{w,b} \quad & \frac{2}{\|w\|} & \iff \min_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y_i (w^\top X_i + b) \geq 1 & \text{s.t.} \quad & y_i (w^\top X_i + b) \geq 1 \end{aligned}$$

其中， $w = \sum_{i=1}^n \alpha_i t_i x_i$ ， $x_i$ 表示数据集特征， $t_i$ 表示label。

我们可以通过复杂的数学推导（此处省略），求解出 $w$ 和 $b$ 的值，进而确定我们的算法参数。为此，我们介绍简化版的SMO算法，该算法可以求解参数 $\alpha$ 和 $b$ ，进而我们可以确定 $w$ 。

## SMO求解

我们首先给出算法的输入输出：

### Algorithm: Simplified SMO

#### Input:

$C$ : regularization parameter

$tol$ : numerical tolerance

$max\_passes$ : max number of times to iterate over  $\alpha$ 's without changing

$((x_1, t_1), \dots, (x_n, t_n))$ : training data

#### Output:

$\alpha \in \mathbb{R}^n$ : Lagrange multipliers for solution

$b \in \mathbb{R}$ : threshold for solution

进而，我们给出算法的伪代码：

- Initialize  $\alpha_i = 0, \forall i; \quad b = 0$ .
- Initialize  $passes = 0$ .
- **while** ( $passes < max\_passes$ )
  - $num\_changed\_alphas = 0$ ;
  - **for**  $i = 1, \dots, n$  :
    - \* Calculate  $E_i = f(x_i) - t_i$  using (45);
    - \* **if** ( $(t_i E_i < -tol \ \&\& \ \alpha_i < C) \ || \ (t_i E_i > tol \ \&\& \ \alpha_i > 0)$ )
      - Select  $j \neq i$  randomly;
      - Calculate  $E_j = f(x_j) - t_j$  using (45);
      - Save old  $\alpha$ 's:  $\alpha_i^{old} = \alpha_i, \ \alpha_j^{old} = \alpha_j$ ;
      - Compute  $L, H$  by (46) or (47);
      - **if** ( $L == H$ ):
        - **continue** to next  $i$ ;
      - Compute  $\eta$  by (48);
      - **if** ( $\eta \geq 0$ ):
        - **continue** to next  $i$ ;
      - Compute and clip new value for  $\alpha_j$  using (49) and (50);
      - **if** ( $|\alpha_j - \alpha_j^{old}| < 10^{-3}$ ): (originally  $10^{-5}$ )
        - **continue** to next  $i$ ;
      - Determine value for  $\alpha_i$  using (51);
      - Compute  $b_1, b_2$  using (52) and (53);
      - Compute  $b$  by (54);
      - $num\_changed\_alphas++$ ;
    - \* **end if**
  - **end for**
  - **if** ( $num\_changed\_alphas == 0$ ):
    - $passes++$ ;
  - **else**
    - $passes = 0$ ;
- **end while**

其中，上述为代码中使用的公式如下所示：

$$f(x) = \omega x^T + b = \sum_{i=1}^n \alpha_i^{old} t_i x_i x^T + b \quad (45)$$

$$t_i \neq t_j \implies L = \max(0, \alpha_j^{old} - \alpha_i^{old}), \quad H = \min(C, C + \alpha_j^{old} - \alpha_i^{old}) \quad (46)$$

$$t_i = t_j \implies L = \max(0, \alpha_j v + \alpha_i^{old} - C), \quad H = \min(C, \alpha_j^{old} + \alpha_i^{old}) \quad (47)$$

$$E_k = f(x_k) - t_k, \quad \eta = 2x_i x_j^T - x_i x_i^T - x_j x_j^T \quad (48)$$

$$\alpha_j^{raw} = \alpha_j^{old} - \frac{t_j (E_i - E_j)}{\eta} \quad (49)$$

$$\alpha_j^{new} = \begin{cases} H & \text{if } \alpha_j^{raw} > H \\ \alpha_j^{raw} & \text{if } L \leq \alpha_j^{raw} \leq H \\ L & \text{if } \alpha_j^{raw} < L. \end{cases} \quad (50)$$

$$\alpha_i^{new} = \alpha_i^{old} + t_i t_j (\alpha_j^{old} - \alpha_j^{new}) \quad (51)$$

$$b_1 = b - E_i - t_i (\alpha_i^{new} - \alpha_i^{old}) x_i x_i^T - t_j (\alpha_j^{new} - \alpha_j^{old}) x_i x_j^T \quad (52)$$

$$b_2 = b - E_j - t_i (\alpha_i^{new} - \alpha_i^{old}) x_i x_j^T - t_j (\alpha_j^{new} - \alpha_j^{old}) x_j x_j^T \quad (53)$$

$$b^{new} = \begin{cases} b_1 & \text{if } 0 < \alpha_i^{new} < C \\ b_2 & \text{if } 0 < \alpha_j^{new} < C \\ (b_1 + b_2)/2 & \text{otherwise.} \end{cases} \quad (54)$$

我们给出该函数 `Svm_Smo` 的代码:

```

1 void Svm_Smo(double *b, double *alpha, int m_passes, double **train_data, double
2 *label, double tol, double C, double change_limit, int row, int col)
3 {
4     srand((unsigned)time(NULL));
5     int p_num = 0;
6     while (p_num < m_passes)
7     {
8         int num_chaged_alpha = 0;
9         for (int i = 0; i < row; i++)
10         {
11             double error_i = calculate_error(*b, label, train_data, alpha, row, col,
12             i);
13             if (((label[i] * error_i < (-tol)) && (alpha[i] < C)) || ((label[i] *
14             error_i > tol) && (alpha[i] > 0)))
15             {
16                 int j = rand() % row;
17                 while (j == i)
18                 {
19                     j = rand() % row;
20                 }
21                 double error_j = calculate_error(*b, label, train_data, alpha, row,
22                 col, j);
23                 // save old alpha i, j
24                 double alpha_old_i = alpha[i];
25                 double alpha_old_j = alpha[j];
26                 // compute L and H
27                 double L = 0, H = C;
28                 if (label[i] != label[j])
29                 {
30                     double L = 0 > (alpha[j] - alpha[i]) ? 0 : (alpha[j] - alpha[i]);
31                     double H = C < (C + alpha[j] - alpha[i]) ? C : (C + alpha[j] -
32                     alpha[i]);
33                     alpha[i] = L;
34                 }
35                 else
36                 {
37                     double L = 0 > (alpha[j] + alpha[i] - C) ? 0 : (alpha[j] +
38                     alpha[i] - C);
39                     double H = C < (alpha[j] + alpha[i]) ? C : (alpha[j] + alpha[i]);
40                     alpha[i] = H;
41                 }
42                 num_chaged_alpha++;
43             }
44         }
45         p_num++;
46     }
47 }

```



```

34         if (L == H)
35         {
36             continue;
37         }
38         // compute eta, in order to be convenient to judge
39         double eta = 2 * array_dot(train_data[i], train_data[j], col - 1) -
40             array_dot(train_data[i], train_data[i], col - 1) -
41             array_dot(train_data[j], train_data[j], col - 1);
42         if (eta >= 0)
43         {
44             continue;
45         }
46         // compute and clip new value for alpha_raw_j
47         alpha[j] -= (label[j] * (error_i - error_j) / eta);
48         // compute alpha_new_j
49         if (alpha[j] > H)
50         {
51             alpha[j] = H;
52         }
53         else if (alpha[j] < L)
54         {
55             alpha[j] = L;
56         }
57         // Check
58         if (fabs(alpha[j] - alpha_old_j) < change_limit)
59         {
60             continue;
61         }
62         // compute alpha_new_i
63         alpha[i] += label[i] * label[j] * (alpha_old_j - alpha[j]);
64         // compute b1, b2
65         double b1 = *b - error_i - label[i] * (alpha[i] - alpha_old_i) *
array_dot(train_data[i], train_data[i], col - 1) -
66             label[j] * (alpha[j] - alpha_old_j) * array_dot(train_data[i],
train_data[j], col - 1);
67         double b2 = *b - error_j - label[i] * (alpha[i] - alpha_old_i) *
array_dot(train_data[i], train_data[j], col - 1) -
68             label[j] * (alpha[j] - alpha_old_j) * array_dot(train_data[j],
train_data[j], col - 1);
69         if ((0 < alpha[i]) && (alpha[i] < C))
70         {
71             *b = b1;
72         }
73         else if ((0 < alpha[j]) && (alpha[j] < C))
74         {
75             *b = b2;
76         }
77         else
78         {
79             *b = (b1 + b2) / 2;
80             num_chaged_alpha += 1;
81         }
82     }
83     else
84     {
85         continue;
86     }
87 }
88 if (num_chaged_alpha == 0)
89 {
90     p_num += 1;
91 }
92 else

```

```

93     {
94         p_num = 0;
95     }
96 }
97 }

```

如此我们求得参数 $\alpha$ 和 $b$ ，再通过函数 `get_weight`，即可求得 $w$ 的值：

```

1  double *get_weight(double *alpha, double *label, double **train_data, int row, int
   col)
2  {
3      double *weight;
4      weight = (double *)malloc((col - 1) * sizeof(double));
5      for (int j = 0; j < col - 1; j++)
6      {
7          weight[j] = 0;
8          for (int i = 0; i < row; i++)
9          {
10             weight[j] += alpha[i] * label[i] * train_data[i][j];
11         }
12     }
13     return weight;
14 }

```

### 3.9.3 算法代码

我们现在知道了如何实现支持向量机算法，那么我们把它应用到[声纳数据集 sonar.csv](#)

我们给出链接：<https://aistudio.baidu.com/aistudio/datasetdetail/105756/0>

### C语言细节讲解

本节假设您已下载数据集 `sonar.csv`，并且它在当前工作目录中可用。下面我们给出一个完整实例，使用C语言详细讲解每一处细节。我们给出每一个.c文件的所有代码：

#### 1) read\_csv.c

该文件代码与前面代码一致，不再重复给出。

#### 2) normalize.c

该文件代码与前面代码一致，不再重复给出。

#### 3) k\_fold.c

该文件代码与前面代码一致，不再重复给出。

#### 4) score.c

该文件代码与前面代码一致，不再重复给出。

#### 5) test\_prediction.c

```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <time.h>
4  #include <math.h>
5
6  double *get_label(double **dataset, int row, int col)
7  {
8      double *label = (double *)malloc(row * sizeof(double));
9      for (int i = 0; i < row; i++)
10     {
11         label[i] = dataset[i][col - 1];

```

```

12     }
13     return label;
14 }
15
16 double array_dot(double *row1, double *row2, int col)
17 {
18     double res = 0;
19     for (int i = 0; i < col; i++)
20     {
21         res += row1[i] * row2[i];
22     }
23     return res;
24 }
25
26 double calculate_error(double b, double *label, double **data, double *alpha, int
row, int column, int index)
27 {
28     double error = 0;
29     double *dot_res;
30     dot_res = (double *)malloc(row * sizeof(double));
31     for (int i = 0; i < row; i++)
32     {
33         dot_res[i] = array_dot(data[i], data[index], column - 1);
34
35         dot_res[i] = dot_res[i] * alpha[i] * label[i];
36         error += dot_res[i];
37     }
38
39     error += b - label[index];
40
41     return error;
42 }
43
44 void Svm_Smo(double *b, double *alpha, int m_passes, double **train_data, double
*label, double tol, double C, double change_limit, int row, int col)
45 {
46     srand((unsigned)time(NULL));
47     int p_num = 0;
48     while (p_num < m_passes)
49     {
50         int num_chaged_alpha = 0;
51         for (int i = 0; i < row; i++)
52         {
53             double error_i = calculate_error(*b, label, train_data, alpha, row, col,
i);
54             if (((label[i] * error_i < (-tol)) && (alpha[i] < C)) || ((label[i] *
error_i > tol) && (alpha[i] > 0)))
55             {
56                 int j = rand() % row;
57                 while (j == i)
58                 {
59                     j = rand() % row;
60                 }
61                 double error_j = calculate_error(*b, label, train_data, alpha, row,
col, j);
62                 // save old alpha i, j
63                 double alpha_old_i = alpha[i];
64                 double alpha_old_j = alpha[j];
65                 // compute L and H
66                 double L = 0, H = C;
67                 if (label[i] != label[j])
68                 {

```

```

69         double L = 0 > (alpha[j] - alpha[i]) ? 0 : (alpha[j] -
alpha[i]);
70         double H = C < (C + alpha[j] - alpha[i]) ? C : (C + alpha[j] -
alpha[i]);
71     }
72     else
73     {
74         double L = 0 > (alpha[j] + alpha[i] - C) ? 0 : (alpha[j] +
alpha[i] - C);
75         double H = C < (alpha[j] + alpha[i]) ? C : (alpha[j] +
alpha[i]);
76     }
77     if (L == H)
78     {
79         continue;
80     }
81     // compute eta, in order to be convenient to judge
82     double eta = 2 * array_dot(train_data[i], train_data[j], col - 1) -
array_dot(train_data[i], train_data[i], col - 1) -
83     array_dot(train_data[j], train_data[j], col - 1);
84     if (eta >= 0)
85     {
86         continue;
87     }
88     // compute and clip new value for alpha_raw_j
89     alpha[j] -= (label[j] * (error_i - error_j) / eta);
90     // compute alpha_new_j
91     if (alpha[j] > H)
92     {
93         alpha[j] = H;
94     }
95     else if (alpha[j] < L)
96     {
97         alpha[j] = L;
98     }
99     // Check
100     if (fabs(alpha[j] - alpha_old_j) < change_limit)
101     {
102         continue;
103     }
104     // compute alpha_new_i
105     alpha[i] += label[i] * label[j] * (alpha_old_j - alpha[j]);
106     // compute b1, b2
107     double b1 = *b - error_i - label[i] * (alpha[i] - alpha_old_i) *
array_dot(train_data[i], train_data[i], col - 1) -
108     label[j] * (alpha[j] - alpha_old_j) * array_dot(train_data[i],
train_data[j], col - 1);
109     double b2 = *b - error_j - label[i] * (alpha[i] - alpha_old_i) *
array_dot(train_data[i], train_data[j], col - 1) -
110     label[j] * (alpha[j] - alpha_old_j) * array_dot(train_data[j],
train_data[j], col - 1);
111     if ((0 < alpha[i]) && (alpha[i] < C))
112     {
113         *b = b1;
114     }
115     else if ((0 < alpha[j]) && (alpha[j] < C))
116     {
117         *b = b2;
118     }
119     else
120     {
121         *b = (b1 + b2) / 2;
122         num_chaged_alpha += 1;
123     }

```

```

124         }
125     }
126     else
127     {
128         continue;
129     }
130 }
131 if (num_chaged_alpha == 0)
132 {
133     p_num += 1;
134 }
135 else
136 {
137     p_num = 0;
138 }
139 }
140 }
141
142 double *get_weight(double *alpha, double *label, double **train_data, int row, int
col)
143 {
144     double *weight;
145     weight = (double *)malloc((col - 1) * sizeof(double));
146     for (int j = 0; j < col - 1; j++)
147     {
148         weight[j] = 0;
149         for (int i = 0; i < row; i++)
150         {
151             weight[j] += alpha[i] * label[i] * train_data[i][j];
152         }
153     }
154     return weight;
155 }
156
157 double predict(double *w, double *test_data, double b, int col)
158 {
159     return array_dot(test_data, w, col - 1) + b;
160 }
161
162 double *get_test_prediction(double **train, int train_size, double **test_data, int
test_size, int m_passes, double tol, double C, double change_limit, int col)
163 {
164     double b = 0;
165     double *alpha = (double *)malloc(train_size * sizeof(double));
166     for (int i = 0; i < train_size; i++)
167     {
168         alpha[i] = 0;
169     }
170     double *label = get_label(train, train_size, col);
171     svm_smo(&b, alpha, m_passes, train, label, tol, C, change_limit, train_size,
col);
172     double *w = get_weight(alpha, label, train, train_size, col);
173     double *predictions = (double *)malloc(test_size * sizeof(double));
174     for (int i = 0; i < test_size; i++)
175     {
176         predictions[i] = predict(w, test_data[i], b, col);
177         if (predictions[i] >= 0)
178         {
179             predictions[i] = 1;
180         }
181         else
182         {
183             predictions[i] = -1;

```

```

184     }
185 }
186 return predictions;
187 }

```

## 6) evaluate.c

```

1  #include <stdlib.h>
2  #include <stdio.h>
3
4  extern double *get_test_prediction(double **train, int train_size, double
   **test_data, int test_size, int m_passes, double tol, double C, double change_limit,
   int col);
5  extern double accuracy_metric(double *actual, double *predicted, int fold_size);
6  extern double ***cross_validation_split(double **dataset, int row, int col, int
   n_folds, int fold_size);
7
8  void evaluate_algorithm(double **dataset, int row, int col, int n_folds, int
   m_passes, double tol, double C, double change_limit)
9  {
10     int fold_size = (int)row / n_folds;
11     double ***split = cross_validation_split(dataset, row, n_folds, fold_size, col);
12     int i, j, k, l;
13     int test_size = fold_size;
14     int train_size = fold_size * (n_folds - 1);
15     double *score = (double *)malloc(n_folds * sizeof(double));
16
17     for (i = 0; i < n_folds; i++)
18     {
19         double ***split_copy = (double ***)malloc(n_folds * sizeof(double **));
20         for (j = 0; j < n_folds; j++)
21         {
22             split_copy[j] = (double **)malloc(fold_size * sizeof(double *));
23             for (k = 0; k < fold_size; k++)
24             {
25                 split_copy[j][k] = (double *)malloc(col * sizeof(double));
26             }
27         }
28         for (j = 0; j < n_folds; j++)
29         {
30             for (k = 0; k < fold_size; k++)
31             {
32                 for (l = 0; l < col; l++)
33                 {
34                     split_copy[j][k][l] = split[j][k][l];
35                 }
36             }
37         }
38         double **test_set = (double **)malloc(test_size * sizeof(double *));
39         for (j = 0; j < test_size; j++)
40         {
41             test_set[j] = (double *)malloc(col * sizeof(double));
42             for (k = 0; k < col; k++)
43             {
44                 test_set[j][k] = split_copy[i][j][k];
45             }
46         }
47         for (j = i; j < n_folds - 1; j++)
48         {
49             split_copy[j] = split_copy[j + 1];
50         }
51         double **train_set = (double **)malloc(train_size * sizeof(double *));

```

```

52     for (k = 0; k < n_folds - 1; k++)
53     {
54         for (l = 0; l < fold_size; l++)
55         {
56             train_set[k * fold_size + l] = (double *)malloc(col *
sizeof(double));
57             train_set[k * fold_size + l] = split_copy[k][l];
58         }
59     }
60     double *predicted = (double *)malloc(test_size * sizeof(double));
61     predicted = get_test_prediction(train_set, train_size, test_set, test_size,
m_passes, tol, C, change_limit, col);
62     double *actual = (double *)malloc(test_size * sizeof(double));
63     for (l = 0; l < test_size; l++)
64     {
65         actual[l] = test_set[l][col - 1];
66     }
67
68     double acc = accuracy_metric(actual, predicted, test_size);
69     score[i] = acc;
70     printf("Scores[%d] = %f%%\n", i, score[i]);
71     free(split_copy);
72 }
73 double total = 0;
74 for (l = 0; l < n_folds; l++)
75 {
76     total += score[l];
77 }
78 printf("mean_accuracy = %f%%\n", total / n_folds);
79 }

```

## 7) main.c

```

1  #include <stdio.h>
2  #include <stdlib.h>
3
4  extern int get_row(char *filename);
5  extern int get_col(char *filename);
6  extern void get_two_dimension(char *line, double **dataset, char *filename);
7  extern void normalize_dataset(double **dataset, int row, int col);
8  extern void evaluate_algorithm(double **dataset, int row, int col, int n_folds, int
m_passes, double tol, double c, double change_limit);
9
10 void main()
11 {
12     char filename[] = "sonar.csv";
13     char line[1024];
14     int row = get_row(filename);
15     int col = get_col(filename);
16     double **dataset;
17     dataset = (double **)malloc(row * sizeof(double *));
18     for (int i = 0; i < row; ++i)
19     {
20         dataset[i] = (double *)malloc(col * sizeof(double));
21     }
22     get_two_dimension(line, dataset, filename);
23     normalize_dataset(dataset, row, col);
24     int n_folds = 5;
25     int max_passes = 1;
26     double C = 0.05;
27     double tolerance = 0.001;
28     double change_limit = 0.001;

```

```

29     evaluate_algorithm(dataset, row, col, n_folds, max_passes, tolerance, C,
30     change_limit);
    }

```

## 8) compile.sh

```

1  gcc main.c read_csv.c normalize.c k_fold.c evaluate.c score.c test_prediction.c -o run
   -lm && ./run

```

编译&运行:

```

1  bash compile.sh

```

最终输出结果如下:

```

1  Scores[0] = 56.097561%
2  Scores[1] = 53.658537%
3  Scores[2] = 56.097561%
4  Scores[3] = 51.219512%
5  Scores[4] = 53.658537%
6  mean_accuracy = 54.146341%

```

## Python语言实战

本节同样假设您已经下载数据集，我们使用著名机器学习开源库sklearn高效实现**支持向量机算法**，以便您在实战中使用该算法：

```

1  import pandas as pd
2  import numpy as np
3  from sklearn.model_selection import KFold
4  from sklearn.svm import SVC
5  from sklearn.metrics import accuracy_score
6
7
8  if __name__ == '__main__':
9      dataset = np.array(pd.read_csv("sonar.csv", sep=',', header=None))
10     k_Cross = KFold(n_splits=5, random_state=0, shuffle=True)
11     index = 0
12     score = np.array([])
13     data,label = dataset[:, :-1], dataset[:, -1]
14     for train_index, test_index in k_Cross.split(dataset):
15         train_data, train_label = data[train_index, :], label[train_index]
16         test_data, test_label = data[test_index, :], label[test_index]
17         model = SVC()
18         model.fit(train_data, train_label)
19         pred = model.predict(test_data)
20         acc = accuracy_score(test_label, pred)
21         score = np.append(score, acc)
22         print('score[{}] = {}'.format(index, acc))
23         index+=1
24     print('mean_accuracy = {}'.format(np.mean(score)))

```

输出结果如下:



```

1 | score[0] = 0.7857142857142857%
2 | score[1] = 0.6904761904761905%
3 | score[2] = 0.7804878048780488%
4 | score[3] = 0.8536585365853658%
5 | score[4] = 0.8292682926829268%
6 | mean_accuracy = 0.7879210220673635%

```

## 3.10 Backpropagation

### 3.10.1 算法介绍

反向传播算法（Backpropagation）是一种适合于多层神经网络的学习算法，通常用于训练大规模的深度学习网络。反向传播算法主要基于梯度下降法，其过程由前向传播、反向传播、权重更新这三步构成。

下面将结合代码，详细阐述反向传播算法在MLP（多层感知机算法）中的应用过程。

### 3.10.2 算法讲解

#### 初始化网络参数

首先我们需要确定网络的层数与节点数，以本文为例，MLP的层数为两层，隐藏层（第一层）节点数为12，输出层（第二层）节点数为3：

```

1 | #define node1 12 //第一层节点数
2 | #define node2 3  //第二层节点数

```

之后定义sigmoid函数及其导数：

```

1 | //激活函数
2 | double sigmoid(double x)
3 | {
4 |     return 1.0 / (1.0 + exp(-x));
5 | }
6 |
7 | //激活函数的导数，y为激活函数值
8 | double dsigmoid(double y)
9 | {
10 |    return y * (1.0 - y);
11 | }

```

之后，根据层数和节点数初始化权重矩阵：

```

1 | double w_1[node1][col]; //第一层权重
2 | double w_2[node2][node1+1]; //第二层权重
3 |
4 | // 初始化权重
5 | for(j=0;j<node1;j++){
6 |     for(k=0;k<col;k++){
7 |         w_1[j][k] = 0.1;
8 |     }
9 | }
10 | for(j=0;j<node2;j++){
11 |     for(k=0;k<node1+1;k++){
12 |         w_2[j][k] = 0.1;
13 |     }
14 | }

```

将权重打印出来如下：

```
1 w_1[0][0] = 0.100000
2 w_1[0][1] = 0.100000
3 w_1[0][2] = 0.100000
4 w_1[0][3] = 0.100000
5 w_1[0][4] = 0.100000
6 w_1[0][5] = 0.100000
7 w_1[0][6] = 0.100000
8 w_1[0][7] = 0.100000
9 w_1[1][0] = 0.100000
10 w_1[1][1] = 0.100000
11 w_1[1][2] = 0.100000
12 w_1[1][3] = 0.100000
13 w_1[1][4] = 0.100000
14 w_1[1][5] = 0.100000
15 w_1[1][6] = 0.100000
16 w_1[1][7] = 0.100000
17 w_1[2][0] = 0.100000
18 w_1[2][1] = 0.100000
19 w_1[2][2] = 0.100000
20 w_1[2][3] = 0.100000
21 w_1[2][4] = 0.100000
22 w_1[2][5] = 0.100000
23 w_1[2][6] = 0.100000
24 w_1[2][7] = 0.100000
25 w_1[3][0] = 0.100000
26 w_1[3][1] = 0.100000
27 w_1[3][2] = 0.100000
28 w_1[3][3] = 0.100000
29 w_1[3][4] = 0.100000
30 w_1[3][5] = 0.100000
31 w_1[3][6] = 0.100000
32 w_1[3][7] = 0.100000
33 w_1[4][0] = 0.100000
34 w_1[4][1] = 0.100000
35 w_1[4][2] = 0.100000
36 w_1[4][3] = 0.100000
37 w_1[4][4] = 0.100000
38 w_1[4][5] = 0.100000
39 w_1[4][6] = 0.100000
40 w_1[4][7] = 0.100000
41 w_1[5][0] = 0.100000
42 w_1[5][1] = 0.100000
43 w_1[5][2] = 0.100000
44 w_1[5][3] = 0.100000
45 w_1[5][4] = 0.100000
46 w_1[5][5] = 0.100000
47 w_1[5][6] = 0.100000
48 w_1[5][7] = 0.100000
49 w_1[6][0] = 0.100000
50 w_1[6][1] = 0.100000
51 w_1[6][2] = 0.100000
52 w_1[6][3] = 0.100000
53 w_1[6][4] = 0.100000
54 w_1[6][5] = 0.100000
55 w_1[6][6] = 0.100000
56 w_1[6][7] = 0.100000
57 w_1[7][0] = 0.100000
58 w_1[7][1] = 0.100000
59 w_1[7][2] = 0.100000
60 w_1[7][3] = 0.100000
61 w_1[7][4] = 0.100000
62 w_1[7][5] = 0.100000
63 w_1[7][6] = 0.100000
```

```

64 w_1[7][7] = 0.100000
65 w_1[8][0] = 0.100000
66 w_1[8][1] = 0.100000
67 w_1[8][2] = 0.100000
68 w_1[8][3] = 0.100000
69 w_1[8][4] = 0.100000
70 w_1[8][5] = 0.100000
71 w_1[8][6] = 0.100000
72 w_1[8][7] = 0.100000
73 w_1[9][0] = 0.100000
74 w_1[9][1] = 0.100000
75 w_1[9][2] = 0.100000
76 w_1[9][3] = 0.100000
77 w_1[9][4] = 0.100000
78 w_1[9][5] = 0.100000
79 w_1[9][6] = 0.100000
80 w_1[9][7] = 0.100000
81 w_1[10][0] = 0.100000
82 w_1[10][1] = 0.100000
83 w_1[10][2] = 0.100000
84 w_1[10][3] = 0.100000
85 w_1[10][4] = 0.100000
86 w_1[10][5] = 0.100000
87 w_1[10][6] = 0.100000
88 w_1[10][7] = 0.100000
89 w_1[11][0] = 0.100000
90 w_1[11][1] = 0.100000
91 w_1[11][2] = 0.100000
92 w_1[11][3] = 0.100000
93 w_1[11][4] = 0.100000
94 w_1[11][5] = 0.100000
95 w_1[11][6] = 0.100000
96 w_1[11][7] = 0.100000
97 w_2[0][0] = 0.100000
98 w_2[0][1] = 0.100000
99 w_2[0][2] = 0.100000
100 w_2[0][3] = 0.100000
101 w_2[0][4] = 0.100000
102 w_2[0][5] = 0.100000
103 w_2[0][6] = 0.100000
104 w_2[0][7] = 0.100000
105 w_2[0][8] = 0.100000
106 w_2[0][9] = 0.100000
107 w_2[0][10] = 0.100000
108 w_2[0][11] = 0.100000
109 w_2[0][12] = 0.100000
110 w_2[1][0] = 0.100000
111 w_2[1][1] = 0.100000
112 w_2[1][2] = 0.100000
113 w_2[1][3] = 0.100000
114 w_2[1][4] = 0.100000
115 w_2[1][5] = 0.100000
116 w_2[1][6] = 0.100000
117 w_2[1][7] = 0.100000
118 w_2[1][8] = 0.100000
119 w_2[1][9] = 0.100000
120 w_2[1][10] = 0.100000
121 w_2[1][11] = 0.100000
122 w_2[1][12] = 0.100000
123 w_2[2][0] = 0.100000
124 w_2[2][1] = 0.100000
125 w_2[2][2] = 0.100000
126 w_2[2][3] = 0.100000

```

```

127 w_2[2][4] = 0.100000
128 w_2[2][5] = 0.100000
129 w_2[2][6] = 0.100000
130 w_2[2][7] = 0.100000
131 w_2[2][8] = 0.100000
132 w_2[2][9] = 0.100000
133 w_2[2][10] = 0.100000
134 w_2[2][11] = 0.100000
135 w_2[2][12] = 0.100000

```

## 前向传播

权重初始化完成后，就可以开始进行训练。首先第一步是前向传播，设 $a_j^i$ 为第 $i$ 层第 $j$ 个神经元的输出， $x_k$ 为输入向量的第 $k$ 个元素， $w_t^{ij}$ 为第 $i$ 层第 $j$ 个神经元第 $t$ 个权重值，计算出各个神经元的输出值：

$$a_j^1 = w_0^{1j} + \sum_t w_t^{1j} x_t$$

$$a_j^i = w_0^{ij} + \sum_t w_t^{ij} a_t^{i-1} \quad (i > 1)$$

代码片段如下：

```

1 double layer1_out[node1]; //第一层节点输出值
2 double layer2_out[node2]; //第二层节点输出值
3 double train[7] = {15.26,14.84,0.871,5.763,3.312,2.221,5.22};
4
5 for(j=0;j<node1;j++){
6     double sum = w_1[j][col-1];
7     for(k=0;k<col-1;k++){
8         sum += w_1[j][k]*train[k];
9     }
10    layer1_out[j] = sigmoid(sum);
11    printf("layer1[%d] = %f\n",j,layer1_out[j]);
12 }
13
14 for(j=0;j<node2;j++){
15     double sum = w_2[j][node1];
16     for(k=0;k<node1;k++){
17         sum += w_2[j][k]*layer1_out[k];
18     }
19    layer2_out[j] = sigmoid(sum);
20    printf("layer2[%d] = %f\n",j,layer2_out[j]);
21 }

```

经计算后，我们可以得到如下结果：

```

1 layer1[0] = 0.992222
2 layer1[1] = 0.992222
3 layer1[2] = 0.992222
4 layer1[3] = 0.992222
5 layer1[4] = 0.992222
6 layer1[5] = 0.992222
7 layer1[6] = 0.992222
8 layer1[7] = 0.992222
9 layer1[8] = 0.992222
10 layer1[9] = 0.992222
11 layer1[10] = 0.992222
12 layer1[11] = 0.992222
13 layer2[0] = 0.784260
14 layer2[1] = 0.784260
15 layer2[2] = 0.784260

```

## 反向传播

前向传播完成后，就可以计算预测值和真实值的误差，然后进行反向传播，为之后的权值更新做准备。设 $\Delta a_j^i$ 为第 $i$ 层第 $j$ 个神经元的误差， $I$ 为总层数，实际值的第 $j$ 个元素为 $y_j$ ， $\delta$ 为激活函数的导数，则其公式为：

$$\Delta a_j^I = (y_j - a_j^I) \delta(a_j^I)$$
$$\Delta a_t^i = \sum_j w_t^{(i+1)j} \delta(a_j^{i+1}) \quad (i < I)$$

本文选取的例子为多分类问题，需要把真实值（1,2,3）转换成one-hot形式（[1,0,0],[0,1,0],[0,0,1]）。定义转换函数如下：

```
1  #define class_num 3 //种类数量
2
3  double *transfer_to_one_hot(int y){
4      double *one_hot = (double *)malloc(class_num*sizeof(double));
5      int i;
6      for(i=0;i<class_num;i++){
7          one_hot[i] = 0;
8      }
9      one_hot[y-1] = 1;
10     return one_hot;
11 }
```

于是利用转换函数把原数据转换为one-hot形式，再计算误差进行反向传播，代码片段如下：

```
1  // 误差反向传播
2  int y = 1;
3  double *target = transfer_to_one_hot(y);
4  double layer2_delta[node2];
5  double layer1_delta[node1];
6  for(j=0;j<node2;j++){
7      double expected= (double) *(target + j);
8      layer2_delta[j] = (expected - layer2_out[j])*dsigmoid(layer2_out[j]);
9  }
10 for(j=0;j<node1;j++){
11     double error = 0.0;
12     for(k=0;k<node2;k++){
13         error += w_2[k][j]*layer2_delta[k];
14     }
15     layer1_delta[j] = error*dsigmoid(layer1_out[j]);
16 }
```

经过计算后，我们可以得到各个层的delta值，结果如下：

```
1  layer2_delta[0] = 0.036502
2  layer2_delta[1] = -0.132694
3  layer2_delta[2] = -0.132694
4  layer1_delta[0] = -0.000177
5  layer1_delta[1] = -0.000177
6  layer1_delta[2] = -0.000177
7  layer1_delta[3] = -0.000177
8  layer1_delta[4] = -0.000177
9  layer1_delta[5] = -0.000177
10 layer1_delta[6] = -0.000177
11 layer1_delta[7] = -0.000177
12 layer1_delta[8] = -0.000177
13 layer1_delta[9] = -0.000177
14 layer1_delta[10] = -0.000177
15 layer1_delta[11] = -0.000177
```

## 权重更新

在得到各个神经元的差值( $\Delta a$ )后, 就可以对原先的权重进行更新。设学习率为 $r$ , 于是可以得到公式:

$$w_t^{1j} = w_t^{1j} + rx_t \Delta a_j^i \quad (t > 0)$$

$$w_0^{1j} = w_0^{1j} + r \Delta a_j^i$$

$$w_t^{ij} = w_t^{ij} + ra_t^i \Delta a_j^i \quad (t > 0)$$

$$w_0^{ij} = w_0^{ij} + r \Delta a_j^i$$

写成代码片段如下:

```
1 // 更新权重
2 double l_rate = 0.01;
3 for(j=0;j<node1;j++){
4     for(k=0;k<col-1;k++){
5         w_1[j][k] += l_rate*layer1_delta[j]*train[k];
6         printf("w_1[%d][%d] = %f\n",j,k,w_1[j][k]);
7     }
8     w_1[j][col] += l_rate*layer1_delta[j];
9     printf("w_1[%d][%d] = %f\n",j,col,w_1[j][col]);
10 }
11 for(j=0;j<node2;j++){
12     for(k=0;k<node1+1;k++){
13         w_2[j][k] += l_rate*layer2_delta[j]*layer1_out[k];
14         printf("w_2[%d][%d] = %f\n",j,k,w_2[j][k]);
15     }
16     w_2[j][node1] += l_rate*layer2_delta[j];
17     printf("w_2[%d][%d] = %f\n",j,col,w_2[j][col]);
18 }
```

把更新后的权重结果打印如下:

```
1 w_1[0][0] = 0.099973
2 w_1[0][1] = 0.099974
3 w_1[0][2] = 0.099998
4 w_1[0][3] = 0.099990
5 w_1[0][4] = 0.099994
6 w_1[0][5] = 0.099996
7 w_1[0][6] = 0.099991
8 w_1[0][8] = 0.099998
9 w_1[1][0] = 0.099971
10 w_1[1][1] = 0.099974
11 w_1[1][2] = 0.099998
12 w_1[1][3] = 0.099990
13 w_1[1][4] = 0.099994
14 w_1[1][5] = 0.099996
15 w_1[1][6] = 0.099991
16 w_1[1][8] = 0.099998
17 w_1[2][0] = 0.099971
18 w_1[2][1] = 0.099974
19 w_1[2][2] = 0.099998
20 w_1[2][3] = 0.099990
21 w_1[2][4] = 0.099994
22 w_1[2][5] = 0.099996
23 w_1[2][6] = 0.099991
24 w_1[2][8] = 0.099998
25 w_1[3][0] = 0.099971
26 w_1[3][1] = 0.099974
27 w_1[3][2] = 0.099998
28 w_1[3][3] = 0.099990
```

```
29 w_1[3][4] = 0.099994
30 w_1[3][5] = 0.099996
31 w_1[3][6] = 0.099991
32 w_1[3][8] = 0.099998
33 w_1[4][0] = 0.099971
34 w_1[4][1] = 0.099974
35 w_1[4][2] = 0.099998
36 w_1[4][3] = 0.099990
37 w_1[4][4] = 0.099994
38 w_1[4][5] = 0.099996
39 w_1[4][6] = 0.099991
40 w_1[4][8] = 0.099998
41 w_1[5][0] = 0.099971
42 w_1[5][1] = 0.099974
43 w_1[5][2] = 0.099998
44 w_1[5][3] = 0.099990
45 w_1[5][4] = 0.099994
46 w_1[5][5] = 0.099996
47 w_1[5][6] = 0.099991
48 w_1[5][8] = 0.099998
49 w_1[6][0] = 0.099971
50 w_1[6][1] = 0.099974
51 w_1[6][2] = 0.099998
52 w_1[6][3] = 0.099990
53 w_1[6][4] = 0.099994
54 w_1[6][5] = 0.099996
55 w_1[6][6] = 0.099991
56 w_1[6][8] = 0.099998
57 w_1[7][0] = 0.099971
58 w_1[7][1] = 0.099974
59 w_1[7][2] = 0.099998
60 w_1[7][3] = 0.099990
61 w_1[7][4] = 0.099994
62 w_1[7][5] = 0.099996
63 w_1[7][6] = 0.099991
64 w_1[7][8] = 0.099998
65 w_1[8][0] = 0.099971
66 w_1[8][1] = 0.099974
67 w_1[8][2] = 0.099998
68 w_1[8][3] = 0.099990
69 w_1[8][4] = 0.099994
70 w_1[8][5] = 0.099996
71 w_1[8][6] = 0.099991
72 w_1[8][8] = 0.099998
73 w_1[9][0] = 0.099971
74 w_1[9][1] = 0.099974
75 w_1[9][2] = 0.099998
76 w_1[9][3] = 0.099990
77 w_1[9][4] = 0.099994
78 w_1[9][5] = 0.099996
79 w_1[9][6] = 0.099991
80 w_1[9][8] = 0.099998
81 w_1[10][0] = 0.099971
82 w_1[10][1] = 0.099974
83 w_1[10][2] = 0.099998
84 w_1[10][3] = 0.099990
85 w_1[10][4] = 0.099994
86 w_1[10][5] = 0.099996
87 w_1[10][6] = 0.099991
88 w_1[10][8] = 0.099998
89 w_1[11][0] = 0.099971
90 w_1[11][1] = 0.099974
91 w_1[11][2] = 0.099998
```

```

92 w_1[11][3] = 0.099990
93 w_1[11][4] = 0.099994
94 w_1[11][5] = 0.099996
95 w_1[11][6] = 0.099991
96 w_1[11][8] = -0.000002
97 w_2[0][0] = 0.100362
98 w_2[0][1] = 0.100362
99 w_2[0][2] = 0.100362
100 w_2[0][3] = 0.100362
101 w_2[0][4] = 0.100362
102 w_2[0][5] = 0.100362
103 w_2[0][6] = 0.100362
104 w_2[0][7] = 0.100362
105 w_2[0][8] = 0.100362
106 w_2[0][9] = 0.100362
107 w_2[0][10] = 0.100362
108 w_2[0][11] = 0.100362
109 w_2[0][12] = 0.100037
110 w_2[0][8] = 0.100362
111 w_2[1][0] = 0.098683
112 w_2[1][1] = 0.098683
113 w_2[1][2] = 0.098683
114 w_2[1][3] = 0.098683
115 w_2[1][4] = 0.098683
116 w_2[1][5] = 0.098683
117 w_2[1][6] = 0.098683
118 w_2[1][7] = 0.098683
119 w_2[1][8] = 0.098683
120 w_2[1][9] = 0.098683
121 w_2[1][10] = 0.098683
122 w_2[1][11] = 0.098683
123 w_2[1][12] = 0.099867
124 w_2[1][8] = 0.098683
125 w_2[2][0] = 0.098683
126 w_2[2][1] = 0.098683
127 w_2[2][2] = 0.098683
128 w_2[2][3] = 0.098683
129 w_2[2][4] = 0.098683
130 w_2[2][5] = 0.098683
131 w_2[2][6] = 0.098683
132 w_2[2][7] = 0.098683
133 w_2[2][8] = 0.098683
134 w_2[2][9] = 0.098683
135 w_2[2][10] = 0.098683
136 w_2[2][11] = 0.098683
137 w_2[2][12] = 0.099867
138 w_2[2][8] = 0.098683

```

## 预测

训练完成后，就可以利用训练好的权重矩阵进行预测。其过程和前向传播大致相同。代码如下：

```

1 // 预测
2 double *predictions = (double *)malloc(test_size*sizeof(double));
3 for(i=0;i<test_size;i++){
4     double out1[node1];
5     for(j=0;j<node1;j++){
6         out1[j] = w_1[j][col-1];
7         for(k=0;k<col-1;k++){
8             out1[j] += w_1[j][k]*test[i][k];
9         }
10        out1[j] = sigmoid(out1[j]);

```



```

11     }
12     double out2[node2];
13     for(j=0;j<node2;j++){
14         double max;
15         out2[j] = w_2[j][node1];
16         for(k=0;k<node1;k++){
17             out2[j] += w_2[j][k]*out1[k];
18         }
19         out2[j] = sigmoid(out2[j]);
20         if(j>0){
21             if(out2[j]>max){
22                 predictions[i] = j+1;
23                 max = out2[j];
24             }
25         }else{
26             predictions[i] = 1;
27             max = out2[j];
28         }
29     }
30 }

```

### 3.10.3 算法代码

我们现在知道了如何实现**BP算法**，那么我们把它应用到[小麦种子数据集 seeds\\_data.csv](https://aistudio.baidu.com/aistudio/datasetdetail/105756/0)

我们给出链接：<https://aistudio.baidu.com/aistudio/datasetdetail/105756/0>

### C语言细节讲解

本节假设您已下载数据集 `seeds_data.csv`，并且它在当前工作目录中可用。下面我们给出一个完整实例，使用C语言详细讲解每一处细节。我们给出每一个.c文件的所有代码：

#### 1) read\_csv.c

该文件代码与前面代码一致，不再重复给出。

#### 2) normalize.c

该文件代码与前面代码一致，不再重复给出。

#### 3) k\_fold.c

该文件代码与前面代码一致，不再重复给出。

#### 4) score.c

该文件代码与前面代码一致，不再重复给出。

#### 5) test\_prediction.c

```

1  #define randval(high) ((double)rand() / RAND_MAX * high)
2  #define uniform_plus_minus_one ((double)(2.0 * rand()) / ((double)RAND_MAX + 1.0) -
   1.0) //均匀随机分布
3
4  #define node1 12 //第一层节点数
5  #define node2 3 //第二层节点数
6  #define class_num 3 //种类数量
7
8  #include "math.h"
9  #include "stdlib.h"
10 #include "time.h"
11 #include "assert.h"
12 #include "string.h"
13 #include "stdio.h"

```

```

14
15 //激活函数
16 double sigmoid(double x)
17 {
18     return 1.0 / (1.0 + exp(-x));
19 }
20
21 //激活函数的导数, y为激活函数值
22 double dsigmoid(double y)
23 {
24     return y * (1.0 - y);
25 }
26
27 double *transfer_to_one_hot(int y)
28 {
29     double *one_hot = (double *)malloc(class_num * sizeof(double));
30     int i;
31     for (i = 0; i < class_num; i++)
32     {
33         one_hot[i] = 0;
34     }
35     one_hot[y - 1] = 1;
36     return one_hot;
37 }
38
39 //训练模型并获得预测值
40 double *get_test_prediction(double **train, double **test, double l_rate, int
n_epoch, int train_size, int test_size, int col)
41 {
42     int epoch, i, j, k;
43     // 初始化权重
44     double w_1[node1][col]; //第一层权重
45     double w_2[node2][node1 + 1]; //第二层权重
46     double out;
47     double predict;
48     double layer1_out[node1]; //第一层节点输出值
49     double layer2_out[node2]; //第二层节点输出值
50     // 初始化权重
51     for (j = 0; j < node1; j++)
52     {
53         for (k = 0; k < col; k++)
54         {
55             w_1[j][k] = uniform_plus_minus_one;
56         }
57     }
58     for (j = 0; j < node2; j++)
59     {
60         for (k = 0; k < node1 + 1; k++)
61         {
62             w_2[j][k] = uniform_plus_minus_one;
63         }
64     }
65     for (epoch = 0; epoch < n_epoch; epoch++)
66     {
67         for (i = 0; i < train_size; i++)
68         {
69             // 前向传播
70             for (j = 0; j < node1; j++)
71             {
72                 double sum = w_1[j][col - 1];
73                 for (k = 0; k < col - 1; k++)
74                 {
75                     sum += w_1[j][k] * train[i][k];

```

```

76         }
77         layer1_out[j] = sigmoid(sum);
78     }
79     for (j = 0; j < node2; j++)
80     {
81         double sum = w_2[j][node1];
82         for (k = 0; k < node1; k++)
83         {
84             sum += w_2[j][k] * layer1_out[k];
85         }
86         layer2_out[j] = sigmoid(sum);
87     }
88     // 误差反向传播
89     int y;
90     y = (int)train[i][col - 1];
91     double *target = transfer_to_one_hot(y);
92     double layer2_delta[node2];
93     double layer1_delta[node1];
94     for (j = 0; j < node2; j++)
95     {
96         double expected = (double)*(target + j);
97         layer2_delta[j] = (expected - layer2_out[j]) *
dsigmoid(layer2_out[j]);
98     }
99     for (j = 0; j < node1; j++)
100     {
101         double error = 0.0;
102         for (k = 0; k < node2; k++)
103         {
104             error += w_2[k][j] * layer2_delta[k];
105         }
106         layer1_delta[j] = error * dsigmoid(layer1_out[j]);
107     }
108     // 更新权重
109     for (j = 0; j < node1; j++)
110     {
111         for (k = 0; k < col - 1; k++)
112         {
113             w_1[j][k] += l_rate * layer1_delta[j] * train[i][k];
114         }
115         w_1[j][col] += l_rate * layer1_delta[j];
116     }
117     for (j = 0; j < node2; j++)
118     {
119         for (k = 0; k < node1 + 1; k++)
120         {
121             w_2[j][k] += l_rate * layer2_delta[j] * layer1_out[k];
122         }
123         w_2[j][node1] += l_rate * layer2_delta[j];
124     }
125 }
126 }
127 // 预测
128 double *predictions = (double *)malloc(test_size * sizeof(double));
129 for (i = 0; i < test_size; i++)
130 {
131     double out1[node1];
132     for (j = 0; j < node1; j++)
133     {
134         out1[j] = w_1[j][col - 1];
135         for (k = 0; k < col - 1; k++)
136         {
137             out1[j] += w_1[j][k] * test[i][k];

```

```

138     }
139     out1[j] = sigmoid(out1[j]);
140 }
141 double out2[node2];
142 for (j = 0; j < node2; j++)
143 {
144     double max;
145     out2[j] = w_2[j][node1];
146     for (k = 0; k < node1; k++)
147     {
148         out2[j] += w_2[j][k] * out1[k];
149     }
150     out2[j] = sigmoid(out2[j]);
151     if (j > 0)
152     {
153         if (out2[j] > max)
154         {
155             predictions[i] = j + 1;
156             max = out2[j];
157         }
158     }
159     else
160     {
161         predictions[i] = 1;
162         max = out2[j];
163     }
164 }
165 }
166 return predictions;
167 }

```

## 6) evaluate.c

```

1  #include <stdlib.h>
2  #include <stdio.h>
3
4  extern double ***cross_validation_split(double **dataset, int row, int n_folds, int
fold_size, int col);
5  extern double *get_test_prediction(double **train, double **test, double l_rate, int
n_epoch, int train_size, int test_size, int col);
6  extern double accuracy_metric(double *actual, double *predicted, int fold_size);
7
8  double *evaluate_algorithm(double **dataset, int n_folds, int fold_size, double
l_rate, int n_epoch, int col, int row)
9  {
10     double ***split = cross_validation_split(dataset, row, n_folds, fold_size, col);
11     int i, j, k, l;
12     int test_size = fold_size;
13     int train_size = fold_size * (n_folds - 1);
14     double *score = (double *)malloc(n_folds * sizeof(double));
15     for (i = 0; i < n_folds; i++)
16     {
17         double ***split_copy = (double ***)malloc(n_folds * sizeof(double **));
18         for (j = 0; j < n_folds; j++)
19         {
20             split_copy[j] = (double **)malloc(fold_size * sizeof(double **));
21             for (k = 0; k < fold_size; k++)
22             {
23                 split_copy[j][k] = (double *)malloc(col * sizeof(double));
24             }
25         }
26         for (j = 0; j < n_folds; j++)

```

```

27     {
28         for (k = 0; k < fold_size; k++)
29         {
30             for (l = 0; l < col; l++)
31             {
32                 split_copy[j][k][l] = split[j][k][l];
33             }
34         }
35     }
36     double **test_set = (double **)malloc(test_size * sizeof(double *));
37     for (j = 0; j < test_size; j++)
38     {
39         test_set[j] = (double *)malloc(col * sizeof(double));
40         for (k = 0; k < col; k++)
41         {
42             test_set[j][k] = split_copy[i][j][k];
43         }
44     }
45     for (j = i; j < n_folds - 1; j++)
46     {
47         split_copy[j] = split_copy[j + 1];
48     }
49     double **train_set = (double **)malloc(train_size * sizeof(double *));
50     for (k = 0; k < n_folds - 1; k++)
51     {
52         for (l = 0; l < fold_size; l++)
53         {
54             train_set[k * fold_size + l] = (double *)malloc(col *
sizeof(double));
55             train_set[k * fold_size + l] = split_copy[k][l];
56         }
57     }
58     double *predicted_2;
59     predicted_2 = get_test_prediction(train_set, test_set, l_rate, n_epoch,
train_size, test_size, col);
60     double predicted[test_size];
61     double *actual = (double *)malloc(test_size * sizeof(double));
62     for (l = 0; l < test_size; l++)
63     {
64         predicted[l] = (double)*(predicted_2 + l);
65         actual[l] = test_set[l][col - 1];
66     }
67     double accuracy = accuracy_metric(actual, predicted, test_size);
68     score[i] = accuracy;
69     printf("score[%d]=%f\n", i, score[i]);
70     free(split_copy);
71 }
72 double total = 0.0;
73 for (l = 0; l < n_folds; l++)
74 {
75     total += score[l];
76 }
77 printf("mean_accuracy=%f\n", total / n_folds);
78 return score;
79 }

```

## 7) main.c

```

1  #include <stdlib.h>
2  #include <math.h>
3  #include <stdio.h>
4  #include <time.h>

```

```

5
6 extern int get_row(char *filename);
7 extern int get_col(char *filename);
8 extern void get_two_dimension(char *line, double **dataset, char *filename);
9 extern double *evaluate_algorithm(double **dataset, int n_folds, int fold_size,
double l_rate, int n_epoch, int col, int row);
10
11 void main()
12 {
13     char filename[] = "seeds_data.csv";
14     char line[1024];
15     int row = get_row(filename);
16     int col = get_col(filename);
17     printf("row = %d\n", row);
18     printf("col = %d\n", col);
19     double **dataset = (double **)malloc(row * sizeof(int *));
20     int i;
21     for (i = 0; i < row; ++i)
22     {
23         dataset[i] = (double *)malloc(col * sizeof(double));
24     }
25     get_two_dimension(line, dataset, filename);
26     double l_rate = 0.1;
27     int n_epoch = 100;
28     int n_folds = 4;
29     int fold_size;
30     fold_size = (int)(row / n_folds);
31     evaluate_algorithm(dataset, n_folds, fold_size, l_rate, n_epoch, col, row);
32 }

```

## 8) compile.sh

```

1 gcc main.c read_csv.c normalize.c k_fold.c evaluate.c score.c test_prediction.c -o run
-lm && ./run

```

### 编译&运行:

```

1 bash compile.sh

```

运算后得到的结果如下:

```

1 score[0] = 95.918367%
2 score[1] = 83.673469%
3 score[2] = 79.591837%
4 score[3] = 57.142857%
5 mean_accuracy = 79.081633%

```

## Python语言实战

本节同样假设您已经下载数据集，我们使用著名机器学习开源库sklearn高效实现**BP算法**，以便您在实战中使用该算法：

```

1 import numpy as np
2 import pandas as pd
3
4 from sklearn.model_selection import KFold
5 from sklearn.metrics import accuracy_score
6 from sklearn.preprocessing import MinMaxScaler
7 from sklearn.neural_network import MLPClassifier
8

```

```

9  import warnings
10 warnings.filterwarnings('ignore')
11
12
13 if __name__ == '__main__':
14     dataset = np.array(pd.read_csv("seeds_data.csv", sep=',', header=None))
15     k_Cross = KFold(n_splits=5, random_state=8, shuffle=True)
16     index = 0
17     score = np.array([])
18     scaler = MinMaxScaler()
19     data, label = dataset[:, :-1], dataset[:, -1]
20     data = scaler.fit_transform(data)
21     for train_index, test_index in k_Cross.split(dataset):
22         train_data, train_label = data[train_index, :], label[train_index]
23         test_data, test_label = data[test_index, :], label[test_index]
24         model = MLPClassifier(max_iter=100)
25         model.fit(train_data, train_label)
26         pred = model.predict(test_data)
27         acc = accuracy_score(test_label, pred) * 100
28         score = np.append(score, acc)
29         print('score[{}] = {}'.format(index, acc))
30         index+=1
31     print('mean_accuracy = {}'.format(np.mean(score)))

```

输出结果如下：

```

1  score[0] = 92.5%
2  score[1] = 90.0%
3  score[2] = 85.0%
4  score[3] = 82.5%
5  score[4] = 100.0%
6  mean_accuracy = 90.0%

```

## 3.11 Bootstrap Aggregation

我们前面已经介绍了决策树算法的C语言实现。决策树是一种简单而强大的预测模型，可是决策树可能面临方差过大的问题，即输入数据的细微偏差会导致预测结果的较大不同，这使得训练出来的结果对于特定的训练数据而言并不具有良好的泛化能力。本节将介绍的算法可以使得决策树拥有更高的鲁棒性，进而在预测中取得更好的表现。这种算法称为bootstrap aggregation（引导聚合），简称bagging（袋装）。它的主要思想是，在原始数据集上，通过**有放回抽样**的方法，重新选择出S个新数据集来分别训练S个分类器的集成技术。也就是说，这些模型训练的数据中允许存在重复的数据。在学习完本节的内容后，您将理解：如何从数据集中创建一个自助样本（bootstrap sample），如何运用这样的自助模型（bootstrapped models）来做预测，以及如何在你自己的预测问题上运用这些技巧。

### 3.11.1 算法介绍

引导（bootstrap）指的是对原始数据集做某种替换形成的样本。这意味着根据现有数据集中随机抽取的样本创建一个新的数据集，这种抽取是有放回的。当可用的数据集比较有限时，这是一种有用的方法。Bootstrap Aggregation（下文简称bagging），其基本思想是给定一个弱学习算法和一个训练集，由于单个弱学习算法准确率不高，所以我们将该学习算法使用多次，得出预测函数序列进行投票，依据得票数的多少来决定预测结果。在本节中，我们将重点介绍以决策树（相关基本知识可参阅本书 **3.5 Classification and Regression Trees**）为基分类器的bagging算法的C语言实现，并以Sonar数据集为例进行演示。

单个决策树虽然原理简单，但其方差较大，这意味着数据的更改越大，算法的性能变化也越大。高方差机器学习算法的性能就可以通过训练许多模型（如许多棵决策树），并取其预测的平均值来提高。这样的结果通常比单个模型要好。

除了提高性能外，Bagging的另一个优势是多个模型的叠加不会产生过拟合，因此我们可以放心地继续添加基分类器，直至满足性能要求。

## 3.11.2 算法讲解

### 引导重采样

Bagging算法的重要内容就是对原始数据集进行有放回重采样，形成一系列子样本。显然，每个样本的抽取可以由随机数实现，子数据集中样本数占原数据集中样本数的比例可预先给定，由此可决定抽取样本的次数。

```
1 double subsample(double **dataset, double ratio)
2 {
3     int n_sample = (int)(row * ratio + 0.5);
4     double **sample=0;
5     sample = (double **)malloc(n_sample*sizeof(int *));
6     for (int i = 0; i < n_sample; i++){
7         sample[i] = (double *)malloc(col*sizeof(double));
8     }
9     int i,j,Index;
10    for(i=0;i<n_sample;i++)
11    {
12        Index = rand() % row;
13        for(j=0;j<col;j++)
14        {
15            sample[i][j] = dataset[Index][j];
16        }
17    }
18    return **sample;
19 }
```

下面我们通过一个具体的例子来体验引导重采样的威力。

首先我们随机生成一个较小的数据集（20个介于0~10之间的整数）。

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <time.h>
4 int main() {
5     int a[20];
6     srand((unsigned int)time(NULL));
7     for(int i=0; i<20; i++) {
8         a[i] = rand()%10;
9         printf("%d ",a[i]);
10    }
11    return 0;
12 }
```

输出：（此处输出仅作数据集生成函数效用的参考，以下演示用到的具体数据集在再次生成的时候必然不同于此）

```
1 // Output:
2 8 9 2 7 5 1 5 2 5 7 5 2 1 9 9 6 2 4 6 3
```

下面我们比较一下不同的采样率得到的子数据集，它们的均值：

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <time.h>
4
5 int main() {
6     int a[20];
7     srand((unsigned int)time(NULL));
8     for(int i=0; i<20; i++) {
```



```

9     a[i] = rand()%10;
10    printf("%d\n",a[i]);
11    }
12
13    double ratio=0.1;
14    double mean=0;
15    int Index;
16    int n_sample = (int)(20 * ratio + 0.5);
17
18    // 计算子采样的均值
19    for (int i = 0; i < n_sample; i++){
20        Index = rand() % 20;
21        mean+=a[Index];
22    }
23    mean/=n_sample;
24
25    // 计算原数据的均值
26    double sum=0;
27    double Tmean=0;
28    for(int i=0; i<20; i++) {
29        sum+=a[i];
30    }
31    Tmean=sum/20;
32
33    printf("\n子采样的均值为%.3f",mean);
34    printf("\n数据集的均值为%.3f",Tmean);
35    return 0;
36 }

```

当采样率为0.1时，可以发现子样本与原样本差别较大：

```

1 // Output:
2 0 8 6 9 3 3 8 3 6 1 2 9 9 9 4 3 4 8 0 2
3 子采样的均值为7.000
4 数据集的均值为4.850

```

随着采样率的提高，子样本越来越能够反映真实值。这与我们的常识也是相符的。

```

1 // ratio=0.2
2 1 8 4 8 4 7 1 7 8 4 8 6 3 3 9 0 2 8 9 1
3 子采样的均值为7.000
4 数据集的均值为5.050
5
6 // ratio=0.3
7 9 9 4 5 8 6 3 8 9 6 0 0 5 4 3 6 9 4 8 5
8 子采样的均值为4.667
9 数据集的均值为5.550
10
11 // ratio=0.4
12 9 5 9 2 1 8 6 0 8 3 1 0 6 7 6 9 8 4 2 8
13 子采样的均值为4.625
14 数据集的均值为5.100
15
16 // ratio=0.6
17 2 8 2 6 6 7 5 2 5 5 4 0 6 8 5 1 1 8 5 2
18 子采样的均值为4.667
19 数据集的均值为4.400
20
21 // ratio=0.8
22 0 9 9 3 0 6 7 3 6 7 3 4 8 9 9 7 8 4 4 6
23 子采样的均值为5.250

```

```

24 数据集的均值为5.600
25
26  // ratio=1.0
27  6 7 5 4 0 4 2 3 6 7 4 1 9 5 8 1 0 6 5 8
28 子采样的均值为4.650
29 数据集的均值为4.550

```

## 构建决策树

该步骤原理与代码可以参考3.5 Classification and Regression Trees.

### 3.11.1 算法代码

我们现在知道了如何实现Bootstrap Aggregation算法，那么我们把它应用到[声纳数据集 sonar.csv](#)

我们给出链接：<https://aistudio.baidu.com/aistudio/datasetdetail/105756/0>

## C语言细节讲解

本节假设您已下载数据集 `sonar.csv`，并且它在当前工作目录中可用。下面我们给出一个完整实例，使用C语言详细讲解每一处细节。我们给出每一个.c文件的所有代码：

### 1) read\_csv.c

该步骤代码与前面CART部分相似，不再重复给出。

### 2) k\_fold.c

该步骤代码与前面CART部分相似，不再重复给出。

### 3) BA.c

```

1  #include "BA.h"
2
3  // 切分函数，根据切分点将数据分为左右两组
4  struct dataset *test_split(int index, double value, int row, int col, double **data)
5  {
6      // 将切分结果作为结构体返回
7      struct dataset *split = (struct dataset *)malloc(sizeof(struct dataset));
8      int count1 = 0, count2 = 0;
9      double ***groups = (double ***)malloc(2 * sizeof(double **));
10     for (int i = 0; i < 2; i++)
11     {
12         groups[i] = (double **)malloc(row * sizeof(double *));
13         for (int j = 0; j < row; j++)
14         {
15             groups[i][j] = (double *)malloc(col * sizeof(double));
16         }
17     }
18     for (int i = 0; i < row; i++)
19     {
20         if (data[i][index] < value)
21         {
22             groups[0][count1] = data[i];
23             count1++;
24         }
25         else
26         {
27             groups[1][count2] = data[i];
28             count2++;
29         }
30     }
31     split->splitdata = groups;
32     split->row1 = count1;

```

```

33     split->row2 = count2;
34     return split;
35 }
36
37 // 计算Gini系数
38 double gini_index(int index, double value, int row, int col, double **dataset,
39 double *class, int classnum)
40 {
41     double *numcount1 = (double *)malloc(classnum * sizeof(double));
42     double *numcount2 = (double *)malloc(classnum * sizeof(double));
43     for (int i = 0; i < classnum; i++)
44         numcount1[i] = numcount2[i] = 0;
45
46     double count1 = 0, count2 = 0;
47     double gini1, gini2, gini;
48     gini1 = gini2 = gini = 0;
49     // 计算每一类的个数
50     for (int i = 0; i < row; i++)
51     {
52         if (dataset[i][index] < value)
53         {
54             count1++;
55             for (int j = 0; j < classnum; j++)
56                 if (dataset[i][col - 1] == class[j])
57                     numcount1[j] += 1;
58         }
59         else
60         {
61             count2++;
62             for (int j = 0; j < classnum; j++)
63                 if (dataset[i][col - 1] == class[j])
64                     numcount2[j]++;
65         }
66     }
67     // 判断分母是否为0, 防止运算错误
68     if (count1 == 0)
69     {
70         gini1 = 1;
71         for (int i = 0; i < classnum; i++)
72             gini2 += (numcount2[i] / count2) * (numcount2[i] / count2);
73     }
74     else if (count2 == 0)
75     {
76         gini2 = 1;
77         for (int i = 0; i < classnum; i++)
78             gini1 += (numcount1[i] / count1) * (numcount1[i] / count1);
79     }
80     else
81     {
82         for (int i = 0; i < classnum; i++)
83         {
84             gini1 += (numcount1[i] / count1) * (numcount1[i] / count1);
85             gini2 += (numcount2[i] / count2) * (numcount2[i] / count2);
86         }
87     }
88     // 计算Gini系数
89     gini1 = 1 - gini1;
90     gini2 = 1 - gini2;
91     gini = (count1 / row) * gini1 + (count2 / row) * gini2;
92     free(numcount1);
93     free(numcount2);
94     numcount1 = numcount2 = NULL;
95     return gini;

```

```

95 }
96
97 // 选取数据的最优切分点
98 struct treeBranch *get_split(int row, int col, double **dataset, double *class, int
classnum, int n_features)
99 {
100     struct treeBranch *tree = (struct treeBranch *)malloc(sizeof(struct
treeBranch));
101     int *featurelist = (int *)malloc(n_features * sizeof(int));
102     int count = 0, flag = 0, temp;
103     int b_index = 999;
104     double b_score = 999, b_value = 999, score;
105     // 随机选取n_features个特征
106     while (count < n_features)
107     {
108         featurelist[count] = count;
109         count++;
110     }
111     // 计算所有切分点Gini系数, 选出Gini系数最小的切分点
112     for (int i = 0; i < n_features; i++)
113     {
114         for (int j = 0; j < row; j++)
115         {
116             double value = dataset[j][featurelist[i]];
117             score = gini_index(featurelist[i], value, row, col, dataset, class,
classnum);
118             if (score < b_score)
119             {
120                 b_score = score;
121                 b_value = value;
122                 b_index = featurelist[i];
123             }
124         }
125     }
126     tree->index = b_index;
127     tree->value = b_value;
128     tree->flag = 0;
129     return tree;
130 }
131
132 // 计算叶节点结果
133 double to_terminal(int row, int col, double **data, double *class, int classnum)
134 {
135     int *num = (int *)malloc(classnum * sizeof(classnum));
136     double maxnum = 0;
137     int flag = 0;
138     // 计算所有样本中结果最多的一类
139     for (int i = 0; i < classnum; i++)
140         num[i] = 0;
141     for (int i = 0; i < row; i++)
142         for (int j = 0; j < classnum; j++)
143             if (data[i][col - 1] == class[j])
144                 num[j]++;
145     for (int j = 0; j < classnum; j++)
146     {
147         if (num[j] > flag)
148         {
149             flag = num[j];
150             maxnum = class[j];
151         }
152     }
153     free(num);
154     num = NULL;

```

```

155     return maxnum;
156 }
157
158 // 创建子树或生成叶节点
159 void split(struct treeBranch *tree, int row, int col, double **data, double *class,
160 int classnum, int depth, int min_size, int max_depth, int n_features)
161 {
162     // 判断是否已经达到最大层数
163     if (depth >= max_depth)
164     {
165         tree->flag = 1;
166         tree->output = to_terminal(row, col, data, class, classnum);
167         return;
168     }
169     struct dataset *childdata = test_split(tree->index, tree->value, row, col,
170 data);
171     // 判断样本是否已被分为一边
172     if (childdata->row1 == 0 || childdata->row2 == 0)
173     {
174         tree->flag = 1;
175         tree->output = to_terminal(row, col, data, class, classnum);
176         return;
177     }
178     // 左子树, 判断样本是否达到最小样本数, 如不是则继续迭代
179     if (childdata->row1 <= min_size)
180     {
181         struct treeBranch *leftchild = (struct treeBranch *)malloc(sizeof(struct
182 treeBranch));
183         leftchild->flag = 1;
184         leftchild->output = to_terminal(childdata->row1, col, childdata-
185 >splitdata[0], class, classnum);
186         tree->leftBranch = leftchild;
187     }
188     else
189     {
190         struct treeBranch *leftchild = get_split(childdata->row1, col, childdata-
191 >splitdata[0], class, classnum, n_features);
192         tree->leftBranch = leftchild;
193         split(leftchild, childdata->row1, col, childdata->splitdata[0], class,
194 classnum, depth + 1, min_size, max_depth, n_features);
195     }
196     // 右子树, 判断样本是否达到最小样本数, 如不是则继续迭代
197     if (childdata->row2 <= min_size)
198     {
199         struct treeBranch *rightchild = (struct treeBranch *)malloc(sizeof(struct
200 treeBranch));
201         rightchild->flag = 1;
202         rightchild->output = to_terminal(childdata->row2, col, childdata-
203 >splitdata[1], class, classnum);
204         tree->rightBranch = rightchild;
205     }
206     else
207     {
208         struct treeBranch *rightchild = get_split(childdata->row2, col, childdata-
209 >splitdata[1], class, classnum, n_features);
210         tree->rightBranch = rightchild;
211         split(rightchild, childdata->row2, col, childdata->splitdata[1], class,
212 classnum, depth + 1, min_size, max_depth, n_features);
213     }
214     free(childdata->splitdata);
215     childdata->splitdata = NULL;
216     free(childdata);
217     childdata = NULL;

```

```

208     return;
209 }
210
211 // 生成决策树
212 struct treeBranch *build_tree(int row, int col, double **data, int min_size, int
max_depth, int n_features)
213 {
214     int count1 = 0, flag1 = 0;
215     // 判断结果一共有多少类别, 此处classes[20]仅仅是取一个较大的数20, 默认类别不可能超过20类
216     double classes[20];
217     for (int i = 0; i < row; i++)
218     {
219         if (count1 == 0)
220         {
221             classes[0] = data[i][col - 1];
222             count1++;
223         }
224         else
225         {
226             flag1 = 0;
227             for (int j = 0; j < count1; j++)
228                 if (classes[j] == data[i][col - 1])
229                     flag1 = 1;
230             if (flag1 == 0)
231             {
232                 classes[count1] = data[i][col - 1];
233                 count1++;
234             }
235         }
236     }
237     // 生成切分点
238     struct treeBranch *result = get_split(row, col, data, classes, count1,
n_features);
239     // 进入迭代, 不断生成子树
240     split(result, row, col, data, classes, count1, 1, min_size, max_depth,
n_features);
241     return result;
242 }
243
244 // 随机森林算法, 将森林结果保存为结构体数组, 并返回结构体二重指针
245 struct treeBranch **random_forest(int row, int col, double **data, int min_size, int
max_depth, int n_features, int n_trees, double sample_size)
246 {
247     struct treeBranch **forest = (struct treeBranch **)malloc(n_trees *
sizeof(struct treeBranch *));
248     int samplenum = (int)(row * sample_size);
249     int temp;
250     // 生成随机训练集
251     double **subsample = (double **)malloc(samplenum * sizeof(double *));
252     for (int i = 0; i < samplenum; i++)
253     {
254         subsample[i] = (double *)malloc(col * sizeof(double));
255     }
256     // 生成所有决策树
257     for (int j = 0; j < n_trees; j++)
258     {
259         for (int i = 0; i < samplenum; i++)
260         {
261             temp = rand() % row;
262             subsample[i] = data[temp];
263         }
264         struct treeBranch *tree = build_tree(samplenum, col, subsample, min_size,
max_depth, n_features);

```

```

265     forest[j] = tree;
266 }
267 return forest;
268 }
269
270 // 决策树预测
271 double treepredict(double *test, struct treeBranch *tree)
272 {
273     double output;
274     // 判断是否达到叶节点, flag=1时为叶节点, flag=0时则继续判断
275     if (tree->flag == 1)
276     {
277         output = tree->output;
278         return output;
279     }
280     else
281     {
282         if (test[tree->index] < tree->value)
283         {
284             output = treepredict(test, tree->leftBranch);
285             return output;
286         }
287         else
288         {
289             output = treepredict(test, tree->rightBranch);
290             return output;
291         }
292     }
293 }
294
295 // 随机森林bagging预测
296 double predict(double *test, struct treeBranch **forest, int n_trees)
297 {
298     double output;
299     double *forest_result = (double *)malloc(n_trees * sizeof(double));
300     double *classes = (double *)malloc(n_trees * sizeof(double));
301     int *num = (int *)malloc(n_trees * sizeof(int));
302     for (int i = 0; i < n_trees; i++)
303         num[i] = 0;
304     int count = 0, flag, temp = 0;
305     // 将每棵树的判断结果保存
306     for (int i = 0; i < n_trees; i++)
307         forest_result[i] = treepredict(test, forest[i]);
308     // bagging选出最后结果
309     for (int i = 0; i < n_trees; i++)
310     {
311         flag = 0;
312         for (int j = 0; j < count; j++)
313         {
314             if (forest_result[i] == classes[j])
315             {
316                 flag = 1;
317                 num[j]++;
318             }
319         }
320         if (flag == 0)
321         {
322             classes[count] = forest_result[i];
323             num[count]++;
324             count++;
325         }
326     }
327     for (int i = 0; i < count; i++)

```

```

328     {
329         if (num[i] > temp)
330         {
331             temp = num[i];
332             output = classes[i];
333         }
334     }
335     return output;
336 }
337
338 // 从样本集中进行有放回子采样得到子样本集
339 double **subsample(double **dataset, double ratio)
340 {
341     int n_sample = (int)(row * ratio + 0.5);
342     double **sample = 0;
343     sample = (double **)malloc(n_sample * sizeof(int *));
344     for (int i = 0; i < n_sample; i++)
345     {
346         sample[i] = (double *)malloc(col * sizeof(double));
347     }
348     int i, j, Index;
349     for (i = 0; i < n_sample; i++)
350     {
351         Index = rand() % row;
352         for (j = 0; j < col; j++)
353         {
354             sample[i][j] = dataset[Index][j];
355         }
356     }
357     return sample;
358 }
359
360 // 用每个基模型进行预测并投票决定最终预测结果
361 double bagging_predict(struct treeBranch **trees, double *row, int n_trees)
362 {
363     int i;
364     int n_classes = 1000;
365     int class_count[1000] = {0};
366     double *predictions = 0;
367     predictions = (double *)malloc(n_trees * sizeof(double));
368     double counts = 0;
369     double result = 0;
370     for (i = 0; i < n_trees; i++)
371     {
372         predictions[i] = predict(row, &trees[i], n_trees);
373         class_count[(int)predictions[i]] += 1;
374     }
375     for (i = 0; i < n_classes; i++)
376     {
377         if (class_count[i] > counts)
378         {
379             counts = class_count[i];
380             result = i;
381         }
382     }
383     return result;
384 }
385
386 // Bagging 算法实现
387 double bagging(double **train, double **test, int max_depth, int min_size, double
sample_size, int n_trees, int n_features)
388 {
389     int i, j;

```



```

390     int n_classes = 1000;
391     int class_count[1000] = {0};
392     double **predictions = 0;
393     predictions = (double **)malloc(n_trees * sizeof(int));
394     for (int i = 0; i < n_trees; i++)
395     {
396         predictions[i] = (double *)malloc(sizeof(test) / sizeof(test[0]) *
sizeof(double));
397     }
398     double counts = 0;
399     double result = 0;
400     double *results = 0;
401     results = (double *)malloc(sizeof(test) / sizeof(test[0]) * sizeof(double *));
402     for (i = 0; i < n_trees; i++)
403     {
404         double **sample = subsample(train, sample_size); // sample_size表示ratio
405         struct treeBranch *tree = build_tree(row, col, sample, min_size, max_depth,
n_features);
406         for (j = 0; j < sizeof(test) / sizeof(test[0]); j++)
407         {
408             predictions[i][j] = predict(test[j], &tree, n_features);
409         }
410     }
411     for (j = 0; j < sizeof(test) / sizeof(test[0]); j++)
412     {
413         for (i = 0; i < n_trees; i++)
414         {
415             class_count[(int)predictions[i][j]] += 1;
416         }
417         for (i = 0; i < n_classes; i++)
418         {
419             if (class_count[i] > counts)
420             {
421                 counts = class_count[i];
422                 result = i;
423             }
424         }
425         results[j] = result;
426     }
427     return *results;
428 }

```

#### 4) BA.h

```

1  #ifndef BA
2  #define BA
3
4  #include <stdio.h>
5  #include <string.h>
6  #include <stdlib.h>
7
8  // 读取csv数据，全局变量
9  double **dataset;
10 int row, col;
11
12 // 树的结构体，flag判断是否为叶节点，index和value为切分点，Brance为对应子树
13 struct treeBranch
14 {
15     int flag;
16     int index;
17     double value;
18     double output;

```

```

19     struct treeBranch *leftBranch;
20     struct treeBranch *rightBranch;
21 };
22
23 // 切分数据，splitdata为切分成左右两组的三维数据，row1为左端数据行数，row2为右端
24 struct dataset
25 {
26     int row1;
27     int row2;
28     double ***splitdata;
29 };
30
31 int get_row(char *filename);
32 int get_col(char *filename);
33 void get_two_dimension(char *line, double **dataset, char *filename);
34 double ***cross_validation_split(double **dataset, int row, int n_folds, int
fold_size);
35 double *get_test_prediction(double **train, double **test, int column, int min_size,
int max_depth, int n_features, int n_trees, double sample_size, int fold_size, int
train_size);
36 struct dataset *test_split(int index, double value, int row, int col, double **data);
37 double gini_index(int index, double value, int row, int col, double **dataset, double
*class, int classnum);
38 struct treeBranch *get_split(int row, int col, double **dataset, double *class, int
classnum, int n_features);
39 double to_terminal(int row, int col, double **data, double *class, int classnum);
40 void split(struct treeBranch *tree, int row, int col, double **data, double *class,
int classnum, int depth, int min_size, int max_depth, int n_features);
41 struct treeBranch *build_tree(int row, int col, double **data, int min_size, int
max_depth, int n_features);
42 struct treeBranch **random_forest(int row, int col, double **data, int min_size, int
max_depth, int n_features, int n_trees, double sample_size);
43 double treepredict(double *test, struct treeBranch *tree);
44 double predict(double *test, struct treeBranch **tree, int n_trees);
45 double accuracy_metric(double *actual, double *predicted, int fold_size);
46 double *evaluate_algorithm(double **dataset, int column, int n_folds, int fold_size,
int min_size, int max_depth, int n_features, int n_trees, double sample_size);
47
48 #endif

```

## 5) score.c

该步骤代码与前面CART部分相似，不再重复给出。

## 6) test\_prediction.c

```

1  #include "BA.h"
2
3  double *get_test_prediction(double **train, double **test, int column, int min_size,
int max_depth, int n_features, int n_trees, double sample_size, int fold_size, int
train_size)
4  {
5      double *predictions = (double *)malloc(fold_size * sizeof(double)); //预测集的行数
就是数组prediction的长度
6      struct treeBranch **forest = random_forest(train_size, column, train, min_size,
max_depth, n_features, n_trees, sample_size);
7      for (int i = 0; i < fold_size; i++)
8      {
9          predictions[i] = predict(test[i], forest, n_trees);
10     }
11     return predictions; //返回对test的预测数组
12 }

```

## 7) evaluate.c

```

1  #include "BA.h"
2
3  double *evaluate_algorithm(double **dataset, int column, int n_folds, int fold_size,
4  int min_size, int max_depth, int n_features, int n_trees, double sample_size)
5  {
6      double ***split = cross_validation_split(dataset, row, n_folds, fold_size);
7      int i, j, k, l;
8      int test_size = fold_size;
9      int train_size = fold_size * (n_folds - 1); //train_size个一维数组
10     double *score = (double *)malloc(n_folds * sizeof(double));
11     for (i = 0; i < n_folds; i++)
12     { //因为要遍历删除, 所以拷贝一份split
13         double ***split_copy = (double ***)malloc(n_folds * sizeof(double **));
14         for (j = 0; j < n_folds; j++)
15         {
16             split_copy[j] = (double **)malloc(fold_size * sizeof(double *));
17             for (k = 0; k < fold_size; k++)
18             {
19                 split_copy[j][k] = (double *)malloc(column * sizeof(double));
20             }
21         }
22         for (j = 0; j < n_folds; j++)
23         {
24             for (k = 0; k < fold_size; k++)
25             {
26                 for (l = 0; l < column; l++)
27                 {
28                     split_copy[j][k][l] = split[j][k][l];
29                 }
30             }
31         }
32         double **test_set = (double **)malloc(test_size * sizeof(double *));
33         for (j = 0; j < test_size; j++)
34         { //对test_size中的每一行
35             test_set[j] = (double *)malloc(column * sizeof(double));
36             for (k = 0; k < column; k++)
37             {
38                 test_set[j][k] = split_copy[i][j][k];
39             }
40         }
41         for (j = i; j < n_folds - 1; j++)
42         {
43             split_copy[j] = split_copy[j + 1];
44         }
45         double **train_set = (double **)malloc(train_size * sizeof(double *));
46         for (k = 0; k < n_folds - 1; k++)
47         {
48             for (l = 0; l < fold_size; l++)
49             {
50                 train_set[k * fold_size + l] = (double *)malloc(column *
51                 sizeof(double));
52                 train_set[k * fold_size + l] = split_copy[k][l];
53             }
54         }
55         double *predicted = (double *)malloc(test_size * sizeof(double)); //predicted
56         有test_size个
57         predicted = get_test_prediction(train_set, test_set, column, min_size,
58         max_depth, n_features, n_trees, sample_size, fold_size, train_size);
59         double *actual = (double *)malloc(test_size * sizeof(double));
60         for (l = 0; l < test_size; l++)

```

```

57     {
58         actual[l] = test_set[l][column - 1];
59     }
60     double accuracy = accuracy_metric(actual, predicted, test_size);
61     score[i] = accuracy;
62     printf("score[%d]=%f%%\n", i, score[i]);
63     free(split_copy);
64 }
65 double total = 0.0;
66 for (l = 0; l < n_folds; l++)
67 {
68     total += score[l];
69 }
70 printf("mean_accuracy=%f%%\n", total / n_folds);
71 return score;
72 }

```

## 8) main.c

```

1  #include "BA.h"
2
3  int main()
4  {
5      char filename[] = "sonar.csv";
6      char line[1024];
7      row = get_row(filename);
8      col = get_col(filename);
9      int n_features = col - 1;
10     dataset = (double **)malloc(row * sizeof(int *));
11     for (int i = 0; i < row; ++i)
12     {
13         dataset[i] = (double *)malloc(col * sizeof(double));
14     } //动态申请二维数组
15     get_two_dimension(line, dataset, filename);
16
17     // 输入模型参数, 包括每个叶子最小样本数、最大层数、树木个数
18     int min_size = 2, max_depth = 10, n_trees = 20;
19     double sample_size = 1;
20     int n_folds = 8;
21     int fold_size = (int)(row / n_folds);
22
23     // Bagging算法, 返回交叉验证正确率
24     double *score = evaluate_algorithm(dataset, col, n_folds, fold_size, min_size,
25     max_depth, n_features, n_trees, sample_size);
26 }

```

## 9) compile.sh

```

1  gcc main.c read_csv.c BA.c k_fold.c evaluate.c score.c test_prediction.c -o run -lm &&
   ./run

```

**编译&运行:**

```

1  bash compile.sh

```

运算后得到的结果如下:

```

1 score[0] = 84.615385%
2 score[1] = 76.923077%
3 score[2] = 92.307692%
4 score[3] = 65.384615%
5 score[4] = 96.153846%
6 score[5] = 80.769231%
7 score[6] = 76.923077%
8 score[7] = 80.769231%
9 mean_accuracy = 81.730769%

```

## Python语言实战

本节同样假设您已经下载数据集，我们使用著名机器学习开源库sklearn高效实现**Bootstrap Aggregation**算法，以便您在实战中使用该算法：

```

1 import numpy as np
2 import pandas as pd
3
4 from sklearn.model_selection import KFold
5 from sklearn.tree import DecisionTreeClassifier
6 from sklearn.metrics import accuracy_score
7 from sklearn.ensemble import BaggingClassifier
8
9
10 if __name__ == '__main__':
11     dataset = np.array(pd.read_csv("sonar.csv", sep=',', header=None))
12     k_Cross = KFold(n_splits=8, random_state=0, shuffle=True)
13     index = 0
14     score = np.array([])
15     data, label = dataset[:, :-1], dataset[:, -1]
16     for train_index, test_index in k_Cross.split(dataset):
17         train_data, train_label = data[train_index, :], label[train_index]
18         test_data, test_label = data[test_index, :], label[test_index]
19         tree = DecisionTreeClassifier()
20         model = BaggingClassifier(base_estimator=tree, n_estimators=500,
21                                 max_samples=1.0, max_features=1.0, bootstrap=True, random_state=1)
22         model.fit(train_data, train_label)
23         pred = model.predict(test_data)
24         acc = accuracy_score(test_label, pred)
25         score = np.append(score, acc)
26         print('score[{}] = {}'.format(index, acc)) * 100
27         index+=1
28     print('mean_accuracy = {}'.format(np.mean(score)))

```

输出结果如下：

```

1 score[0] = 76.92307692307693%
2 score[1] = 84.61538461538461%
3 score[2] = 73.07692307692307%
4 score[3] = 69.23076923076923%
5 score[4] = 80.76923076923077%
6 score[5] = 80.76923076923077%
7 score[6] = 84.61538461538461%
8 score[7] = 80.76923076923077%
9 mean_accuracy = 78.84615384615384%

```

## 3.12 Random Forest

我们前面已经介绍了决策树和Bagging算法的C语言实现。随机森林是基于Bagging的一个扩展，它在Bagging对原始训练数据进行有放回抽样形成子数据集的基础上，对构建决策树的特征也进行随机的选择。不同的树之间可能长得非常不同，这种做法可以提高算法的性能。在本节中，您将了解Bagging和随机森林的区别，如何由高方差的决策树构建随机森林，如何将随机森林算法应用于实际的预测问题中。

### 3.12.1 算法介绍

Bagging算法虽然已在一定程度上提高了单棵决策树的性能，但它也有一些限制。它对每一棵构建的树，应用相同的贪婪算法，这将导致每棵树中都会选中相似的分割点，也就是说每一棵树将长得非常相似。相似的树做出的预测也是相似的，这样让它们做投票就达不到预想的效果。为了解决这一问题，我们可以通过限制每一棵树分割点评估的特征来强制决策树不同，这种算法就叫做随机森林。和bagging相同的是，都需要获取原始数据集的多个子样本，并用之分别训练不同的决策树。与bagging不同的是，在每个点上对数据进行拆分并将其添加到树中时，只能考虑固定的特征子集，也就是说对每一棵树所关注的特征，也进行随机选择。

一般来说，对于分类问题，每一棵树考虑的分割特征数量限制为输入特征数量的平方根。即

$$num\_features\_for\_split = \sqrt{total\_input\_features}$$

如此一来，树与树之间的差异会更大，从而导致预测结果多元化，组合预测的性能将比单棵决策树或者只进行套袋要好很多。

### 3.12.2 算法讲解

#### 2.1 构建决策树

与构建单纯决策树的不同之处，主要体现在寻找最佳分点的方式。需要随机选取`n_features`（预先设定）个特征值用于树的构建。

```
1 struct treeBranch *get_split(int row, int col, double **dataset, double *class, int
  classnum, int n_features)
2 {
3     struct treeBranch *tree=(struct treeBranch *)malloc(sizeof(struct treeBranch));
4     int *featurelist=(int *)malloc(n_features * sizeof(int));
5     int count=0,flag=0,temp;
6     int b_index=999;
7     double b_score = 999, b_value = 999,score;
8     // 随机选取n_features个特征
9     while (count<n_features)
10    {
11        flag=0;
12        temp=rand()%(col-1);
13        for (int i = 0; i < count; i++)
14            if (temp==featurelist[i])
15                flag=1;
16        if (flag==0)
17        {
18            featurelist[count]=temp;
19            count++;
20        }
21    }
22    // 计算所有切分点Gini系数，选出Gini系数最小的切分点
23    for (int i = 0; i < n_features; i++)
24    {
25        for (int j = 0; j < row; j++)
26        {
27            double value=dataset[j][featurelist[i]];
28            score = gini_index(featurelist[i], value, row, col, dataset, class,
classnum);
29            if (score<b_score)
30            {
31                b_score = score;
32                b_value = value;
```

```

33         b_index = featurelist[i];
34     }
35 }
36 }
37 tree->index=b_index;tree->value=b_value;tree->flag=0;
38 return tree;
39 }

```

该步骤其他代码与CART的代码一致，不再重复给出。

## 2.2 随机森林算法

由2.1的步骤得到单棵决策树，将森林结果保存为结构体数组，并返回结构体二重指针。其中每棵树是用随机采样的训练集训练出来的。

```

1  struct treeBranch **random_forest(int row, int col, double **data, int min_size, int
max_depth, int n_features, int n_trees, float sample_size)
2  {
3      struct treeBranch ** forest = (struct treeBranch **)malloc(n_trees *
sizeof(struct treeBranch*));
4      int samplenum = (int)(row*sample_size);
5      int temp;
6      // 生成随机训练集
7      double ** subsample=(double **)malloc(samplenum * sizeof(double *));
8      for (int i = 0; i < samplenum; i++)
9      {
10         subsample[i]=(double *)malloc(col * sizeof(double));
11     }
12     // 生成所有决策树
13     for (int j = 0; j < n_trees; j++)
14     {
15         for (int i = 0; i < samplenum; i++)
16         {
17             temp = rand() % row;
18             subsample[i] = data[temp];
19         }
20         struct treeBranch *tree = build_tree(samplenum, col, subsample, min_size,
max_depth, n_features);
21         forest[j]=tree;
22     }
23     return forest;
24 }

```

### 3.12.3 算法代码

我们现在知道了如何实现**随机森林算法**，那么我们把它应用到[声纳数据集 sonar.csv](https://aistudio.baidu.com/aistudio/datasetdetail/105756/0)

我们给出链接：<https://aistudio.baidu.com/aistudio/datasetdetail/105756/0>

## C语言细节讲解

本节假设您已下载数据集 `sonar.csv`，并且它在当前工作目录中可用。下面我们给出一个完整实例，使用C语言详细讲解每一处细节。我们给出每一个.c文件的所有代码：

### 1) read\_csv.c

该步骤代码与前面CART部分相似，不再重复给出。

## 2) k\_fold.c

该步骤代码与前面CART部分相似，不再重复给出。

## 3) RF.c

```
1  #include "RF.h"
2
3  // 切分函数，根据切分点将数据分为左右两组
4  struct dataset *test_split(int index, double value, int row, int col, double **data)
5  {
6      // 将切分结果作为结构体返回
7      struct dataset *split = (struct dataset *)malloc(sizeof(struct dataset));
8      int count1 = 0, count2 = 0;
9      double ***groups = (double ***)malloc(2 * sizeof(double **));
10     for (int i = 0; i < 2; i++)
11     {
12         groups[i] = (double **)malloc(row * sizeof(double *));
13         for (int j = 0; j < row; j++)
14         {
15             groups[i][j] = (double *)malloc(col * sizeof(double));
16         }
17     }
18     for (int i = 0; i < row; i++)
19     {
20         if (data[i][index] < value)
21         {
22             groups[0][count1] = data[i];
23             count1++;
24         }
25         else
26         {
27             groups[1][count2] = data[i];
28             count2++;
29         }
30     }
31     split->splitdata = groups;
32     split->row1 = count1;
33     split->row2 = count2;
34     return split;
35 }
36
37 // 计算Gini系数
38 double gini_index(int index, double value, int row, int col, double **dataset,
39 double *class, int classnum)
40 {
41     double *numcount1 = (double *)malloc(classnum * sizeof(double));
42     double *numcount2 = (double *)malloc(classnum * sizeof(double));
43     for (int i = 0; i < classnum; i++)
44         numcount1[i] = numcount2[i] = 0;
45
46     double count1 = 0, count2 = 0;
47     double gini1, gini2, gini;
48     gini1 = gini2 = gini = 0;
49     // 计算每一类的个数
50     for (int i = 0; i < row; i++)
51     {
52         if (dataset[i][index] < value)
53         {
54             count1++;
55             for (int j = 0; j < classnum; j++)
56                 if (dataset[i][col - 1] == class[j])
57                     numcount1[j] += 1;
```



```

57     }
58     else
59     {
60         count2++;
61         for (int j = 0; j < classnum; j++)
62             if (dataset[i][col - 1] == class[j])
63                 numcount2[j]++;
64     }
65 }
66 // 判断分母是否为0, 防止运算错误
67 if (count1 == 0)
68 {
69     gini1 = 1;
70     for (int i = 0; i < classnum; i++)
71         gini2 += (numcount2[i] / count2) * (numcount2[i] / count2);
72 }
73 else if (count2 == 0)
74 {
75     gini2 = 1;
76     for (int i = 0; i < classnum; i++)
77         gini1 += (numcount1[i] / count1) * (numcount1[i] / count1);
78 }
79 else
80 {
81     for (int i = 0; i < classnum; i++)
82     {
83         gini1 += (numcount1[i] / count1) * (numcount1[i] / count1);
84         gini2 += (numcount2[i] / count2) * (numcount2[i] / count2);
85     }
86 }
87 // 计算Gini系数
88 gini1 = 1 - gini1;
89 gini2 = 1 - gini2;
90 gini = (count1 / row) * gini1 + (count2 / row) * gini2;
91 free(numcount1);
92 free(numcount2);
93 numcount1 = numcount2 = NULL;
94 return gini;
95 }
96
97 // 选取数据的最优切分点
98 struct treeBranch *get_split(int row, int col, double **dataset, double *class, int
classnum, int n_features)
99 {
100     struct treeBranch *tree = (struct treeBranch *)malloc(sizeof(struct
treeBranch));
101     int *featurelist = (int *)malloc(n_features * sizeof(int));
102     int count = 0, flag = 0, temp;
103     int b_index = 999;
104     double b_score = 999, b_value = 999, score;
105     // 随机选取n_features个特征
106     while (count < n_features)
107     {
108         flag = 0;
109         temp = rand() % (col - 1);
110         for (int i = 0; i < count; i++)
111             if (temp == featurelist[i])
112                 flag = 1;
113         if (flag == 0)
114         {
115             featurelist[count] = temp;
116             count++;
117         }

```

```

118     }
119     // 计算所有切分点Gini系数，选出Gini系数最小的切分点
120     for (int i = 0; i < n_features; i++)
121     {
122         for (int j = 0; j < row; j++)
123         {
124             double value = dataset[j][featurelist[i]];
125             score = gini_index(featurelist[i], value, row, col, dataset, class,
classnum);
126             if (score < b_score)
127             {
128                 b_score = score;
129                 b_value = value;
130                 b_index = featurelist[i];
131             }
132         }
133     }
134     tree->index = b_index;
135     tree->value = b_value;
136     tree->flag = 0;
137     return tree;
138 }
139
140 // 计算叶节点结果
141 double to_terminal(int row, int col, double **data, double *class, int classnum)
142 {
143     int *num = (int *)malloc(classnum * sizeof(classnum));
144     double maxnum = 0;
145     int flag = 0;
146     // 计算所有样本中结果最多的一类
147     for (int i = 0; i < classnum; i++)
148         num[i] = 0;
149     for (int i = 0; i < row; i++)
150         for (int j = 0; j < classnum; j++)
151             if (data[i][col - 1] == class[j])
152                 num[j]++;
153     for (int j = 0; j < classnum; j++)
154     {
155         if (num[j] > flag)
156         {
157             flag = num[j];
158             maxnum = class[j];
159         }
160     }
161     free(num);
162     num = NULL;
163     return maxnum;
164 }
165
166 // 创建子树或生成叶节点
167 void split(struct treeBranch *tree, int row, int col, double **data, double *class,
int classnum, int depth, int min_size, int max_depth, int n_features)
168 {
169     // 判断是否已经达到最大层数
170     if (depth >= max_depth)
171     {
172         tree->flag = 1;
173         tree->output = to_terminal(row, col, data, class, classnum);
174         return;
175     }
176     struct dataset *childdata = test_split(tree->index, tree->value, row, col,
data);
177     // 判断样本是否已被分为一边

```

```

178     if (childdata->row1 == 0 || childdata->row2 == 0)
179     {
180         tree->flag = 1;
181         tree->output = to_terminal(row, col, data, class, classnum);
182         return;
183     }
184     // 左子树, 判断样本是否达到最小样本数, 如不是则继续迭代
185     if (childdata->row1 <= min_size)
186     {
187         struct treeBranch *leftchild = (struct treeBranch *)malloc(sizeof(struct
179 treeBranch));
188         leftchild->flag = 1;
189         leftchild->output = to_terminal(childdata->row1, col, childdata-
180 >splitdata[0], class, classnum);
190         tree->leftBranch = leftchild;
191     }
192     else
193     {
194         struct treeBranch *leftchild = get_split(childdata->row1, col, childdata-
181 >splitdata[0], class, classnum, n_features);
195         tree->leftBranch = leftchild;
196         split(leftchild, childdata->row1, col, childdata->splitdata[0], class,
182 classnum, depth + 1, min_size, max_depth, n_features);
197     }
198     // 右子树, 判断样本是否达到最小样本数, 如不是则继续迭代
199     if (childdata->row2 <= min_size)
200     {
201         struct treeBranch *rightchild = (struct treeBranch *)malloc(sizeof(struct
183 treeBranch));
202         rightchild->flag = 1;
203         rightchild->output = to_terminal(childdata->row2, col, childdata-
184 >splitdata[1], class, classnum);
204         tree->rightBranch = rightchild;
205     }
206     else
207     {
208         struct treeBranch *rightchild = get_split(childdata->row2, col, childdata-
185 >splitdata[1], class, classnum, n_features);
209         tree->rightBranch = rightchild;
210         split(rightchild, childdata->row2, col, childdata->splitdata[1], class,
186 classnum, depth + 1, min_size, max_depth, n_features);
211     }
212     free(childdata->splitdata);
213     childdata->splitdata = NULL;
214     free(childdata);
215     childdata = NULL;
216     return;
217 }
218
219 // 生成决策树
220 struct treeBranch *build_tree(int row, int col, double **data, int min_size, int
221 max_depth, int n_features)
222 {
223     int count1 = 0, flag1 = 0;
224     // 判断结果一共有多少类别, 此处classes[20]仅仅是取一个较大的数20, 默认类别不可能超过20类
225     double classes[20];
226     for (int i = 0; i < row; i++)
227     {
228         if (count1 == 0)
229         {
230             classes[0] = data[i][col - 1];
231             count1++;
232         }
233     }

```

```

232         else
233         {
234             flag1 = 0;
235             for (int j = 0; j < count1; j++)
236                 if (classes[j] == data[i][col - 1])
237                     flag1 = 1;
238             if (flag1 == 0)
239             {
240                 classes[count1] = data[i][col - 1];
241                 count1++;
242             }
243         }
244     }
245     // 生成切分点
246     struct treeBranch *result = get_split(row, col, data, classes, count1,
n_features);
247     // 进入迭代, 不断生成子树
248     split(result, row, col, data, classes, count1, 1, min_size, max_depth,
n_features);
249     return result;
250 }
251
252 // 随机森林算法, 将森林结果保存为结构体数组, 并返回结构体二重指针
253 struct treeBranch **random_forest(int row, int col, double **data, int min_size, int
max_depth, int n_features, int n_trees, double sample_size)
254 {
255     struct treeBranch **forest = (struct treeBranch **)malloc(n_trees *
sizeof(struct treeBranch *));
256     int samplenum = (int)(row * sample_size);
257     int temp;
258     // 生成随机训练集
259     double **subsample = (double **)malloc(samplenum * sizeof(double *));
260     for (int i = 0; i < samplenum; i++)
261     {
262         subsample[i] = (double *)malloc(col * sizeof(double));
263     }
264     // 生成所有决策树
265     for (int j = 0; j < n_trees; j++)
266     {
267         for (int i = 0; i < samplenum; i++)
268         {
269             temp = rand() % row;
270             subsample[i] = data[temp];
271         }
272         struct treeBranch *tree = build_tree(samplenum, col, subsample, min_size,
max_depth, n_features);
273         forest[j] = tree;
274     }
275     return forest;
276 }
277
278 // 决策树预测
279 double treepredict(double *test, struct treeBranch *tree)
280 {
281     double output;
282     // 判断是否达到叶节点, flag=1时为叶节点, flag=0时则继续判断
283     if (tree->flag == 1)
284     {
285         output = tree->output;
286         return output;
287     }
288     else
289     {

```

```

290     if (test[tree->index] < tree->value)
291     {
292         output = treepredict(test, tree->leftBranch);
293         return output;
294     }
295     else
296     {
297         output = treepredict(test, tree->rightBranch);
298         return output;
299     }
300 }
301 }
302
303 // 随机森林bagging预测
304 double predict(double *test, struct treeBranch **forest, int n_trees)
305 {
306     double output;
307     double *forest_result = (double *)malloc(n_trees * sizeof(double));
308     double *classes = (double *)malloc(n_trees * sizeof(double));
309     int *num = (int *)malloc(n_trees * sizeof(int));
310     for (int i = 0; i < n_trees; i++)
311         num[i] = 0;
312     int count = 0, flag, temp = 0;
313     // 将每棵树的判断结果保存
314     for (int i = 0; i < n_trees; i++)
315         forest_result[i] = treepredict(test, forest[i]);
316     // bagging选出最后结果
317     for (int i = 0; i < n_trees; i++)
318     {
319         flag = 0;
320         for (int j = 0; j < count; j++)
321         {
322             if (forest_result[i] == classes[j])
323             {
324                 flag = 1;
325                 num[j]++;
326             }
327         }
328         if (flag == 0)
329         {
330             classes[count] = forest_result[i];
331             num[count]++;
332             count++;
333         }
334     }
335     for (int i = 0; i < count; i++)
336     {
337         if (num[i] > temp)
338         {
339             temp = num[i];
340             output = classes[i];
341         }
342     }
343     return output;
344 }

```

#### 4) RF.h

```

1  #ifndef RF
2  #define RF
3
4  #include <stdio.h>

```

```

5  #include <string.h>
6  #include <stdlib.h>
7
8  // 读取csv数据, 全局变量
9  double **dataset;
10 int row, col;
11
12 // 树的结构体, flag判断是否为叶节点, index和value为切分点, Brance为对应子树
13 struct treeBranch
14 {
15     int flag;
16     int index;
17     double value;
18     double output;
19     struct treeBranch *leftBranch;
20     struct treeBranch *rightBranch;
21 };
22
23 // 切分数据, splitdata为切分成左右两组的三维数据, row1为左端数据行数, row2为右端
24 struct dataset
25 {
26     int row1;
27     int row2;
28     double ***splitdata;
29 };
30
31 int get_row(char *filename);
32 int get_col(char *filename);
33 void get_two_dimension(char *line, double **dataset, char *filename);
34 double ***cross_validation_split(double **dataset, int row, int n_folds, int
fold_size);
35 double *get_test_prediction(double **train, double **test, int column, int min_size,
int max_depth, int n_features, int n_trees, double sample_size, int fold_size, int
train_size);
36
37 struct dataset *test_split(int index, double value, int row, int col, double **data);
38 double gini_index(int index, double value, int row, int col, double **dataset, double
*class, int classnum);
39 struct treeBranch *get_split(int row, int col, double **dataset, double *class, int
classnum, int n_features);
40 double to_terminal(int row, int col, double **data, double *class, int classnum);
41 void split(struct treeBranch *tree, int row, int col, double **data, double *class,
int classnum, int depth, int min_size, int max_depth, int n_features);
42 struct treeBranch *build_tree(int row, int col, double **data, int min_size, int
max_depth, int n_features);
43 struct treeBranch **random_forest(int row, int col, double **data, int min_size, int
max_depth, int n_features, int n_trees, double sample_size);
44 double treepredict(double *test, struct treeBranch *tree);
45 double predict(double *test, struct treeBranch **tree, int n_trees);
46 double accuracy_metric(double *actual, double *predicted, int fold_size);
47 double *evaluate_algorithm(double **dataset, int column, int n_folds, int fold_size,
int min_size, int max_depth, int n_features, int n_trees, double sample_size);
48
49 #endif

```

## 5) score.c

该步骤代码与前面CART部分相似, 不再重复给出。

## 6) test\_prediction.c

```
1 #include "RF.h"
2
3 double *get_test_prediction(double **train, double **test, int column, int min_size,
4 int max_depth, int n_features, int n_trees, double sample_size, int fold_size, int
5 train_size)
6 {
7     double *predictions = (double *)malloc(fold_size * sizeof(double)); //预测集的行数
8     //就是数组prediction的长度
9     struct treeBranch **forest = random_forest(train_size, column, train, min_size,
10 max_depth, n_features, n_trees, sample_size);
11     for (int i = 0; i < fold_size; i++)
12     {
13         predictions[i] = predict(test[i], forest, n_trees);
14     }
15     return predictions; //返回对test的预测数组
16 }
```

## 7) evaluate.c

```
1 #include "RF.h"
2
3 double *evaluate_algorithm(double **dataset, int column, int n_folds, int fold_size,
4 int min_size, int max_depth, int n_features, int n_trees, double sample_size)
5 {
6     double ***split = cross_validation_split(dataset, row, n_folds, fold_size);
7     int i, j, k, l;
8     int test_size = fold_size;
9     int train_size = fold_size * (n_folds - 1); //train_size个一维数组
10    double *score = (double *)malloc(n_folds * sizeof(double));
11    for (i = 0; i < n_folds; i++)
12    { //因为要遍历删除，所以拷贝一份split
13        double ***split_copy = (double ***)malloc(n_folds * sizeof(double **));
14        for (j = 0; j < n_folds; j++)
15        {
16            split_copy[j] = (double **)malloc(fold_size * sizeof(double *));
17            for (k = 0; k < fold_size; k++)
18            {
19                split_copy[j][k] = (double *)malloc(column * sizeof(double));
20            }
21        }
22        for (j = 0; j < n_folds; j++)
23        {
24            for (k = 0; k < fold_size; k++)
25            {
26                for (l = 0; l < column; l++)
27                {
28                    split_copy[j][k][l] = split[j][k][l];
29                }
30            }
31        }
32        double **test_set = (double **)malloc(test_size * sizeof(double *));
33        for (j = 0; j < test_size; j++)
34        { //对test_size中的每一行
35            test_set[j] = (double *)malloc(column * sizeof(double));
36            for (k = 0; k < column; k++)
37            {
38                test_set[j][k] = split_copy[i][j][k];
39            }
40        }
41        for (j = i; j < n_folds - 1; j++)
```

```

41     {
42         split_copy[j] = split_copy[j + 1];
43     }
44     double **train_set = (double **)malloc(train_size * sizeof(double *));
45     for (k = 0; k < n_folds - 1; k++)
46     {
47         for (l = 0; l < fold_size; l++)
48         {
49             train_set[k * fold_size + l] = (double *)malloc(column *
50             sizeof(double));
51             train_set[k * fold_size + l] = split_copy[k][l];
52         }
53     }
54     double *predicted = (double *)malloc(test_size * sizeof(double)); //predicted
有test_size个
55     predicted = get_test_prediction(train_set, test_set, column, min_size,
max_depth, n_features, n_trees, sample_size, fold_size, train_size);
56     double *actual = (double *)malloc(test_size * sizeof(double));
57     for (l = 0; l < test_size; l++)
58     {
59         actual[l] = test_set[l][column - 1];
60     }
61     double accuracy = accuracy_metric(actual, predicted, test_size);
62     score[i] = accuracy;
63     printf("score[%d] = %f%%\n", i, score[i]);
64     free(split_copy);
65 }
66 double total = 0.0;
67 for (l = 0; l < n_folds; l++)
68 {
69     total += score[l];
70 }
71 printf("mean_accuracy = %f%%\n", total / n_folds);
72 return score;
73 }

```

## 8) main.c

```

1  #include "RF.h"
2
3  int main()
4  {
5      char filename[] = "sonar.csv";
6      char line[1024];
7      row = get_row(filename);
8      col = get_col(filename);
9      dataset = (double **)malloc(row * sizeof(int *));
10     for (int i = 0; i < row; ++i)
11     {
12         dataset[i] = (double *)malloc(col * sizeof(double));
13     } //动态申请二维数组
14     get_two_dimension(line, dataset, filename);
15
16     // 输入模型参数, 包括每个叶子最小样本数、最大层数、特征值选取个数、树木个数
17     int min_size = 2, max_depth = 10, n_features = 7, n_trees = 100;
18     double sample_size = 1;
19     int n_folds = 5;
20     int fold_size = (int)(row / n_folds);
21
22     // 随机森林算法, 返回交叉验证正确率
23     double *score = evaluate_algorithm(dataset, col, n_folds, fold_size, min_size,
max_depth, n_features, n_trees, sample_size);

```



## 9) compile.sh

```
1 gcc main.c read_csv.c BA.c k_fold.c evaluate.c score.c test_prediction.c -o run -lm &&
  ./run
```

编译&运行:

```
1 bash compile.sh
```

运算后得到的结果如下:

```
1 score[0] = 73.170732%
2 score[1] = 68.292683%
3 score[2] = 58.536585%
4 score[3] = 68.292683%
5 score[4] = 65.853659%
6 mean_accuracy = 66.829268%
```

## Python语言实战

本节同样假设您已经下载数据集，我们使用著名机器学习开源库sklearn高效实现**随机森林算法**，以便您在实战中使用该算法：

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import KFold
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.metrics import accuracy_score
6
7
8 if __name__ == '__main__':
9     dataset = np.array(pd.read_csv("sonar.csv", sep=',', header=None))
10    k_Cross = KFold(n_splits=5, random_state=0, shuffle=True)
11    index = 0
12    score = np.array([])
13    data,label = dataset[:, :-1], dataset[:, -1]
14    for train_index, test_index in k_Cross.split(dataset):
15        train_data, train_label = data[train_index, :], label[train_index]
16        test_data, test_label = data[test_index, :], label[test_index]
17        model = RandomForestClassifier()
18        model.fit(train_data, train_label)
19        pred = model.predict(test_data)
20        acc = accuracy_score(test_label, pred)
21        score = np.append(score, acc)
22        print('score[{}] = {}'.format(index, acc))
23        index+=1
24    print('mean_accuracy = {}'.format(np.mean(score)))
```

输出结果如下:

```
1 score[0] = 0.8333333333333334%
2 score[1] = 0.6666666666666666%
3 score[2] = 0.8095238095238095%
4 score[3] = 0.8048780487804879%
5 score[4] = 0.8780487804878049%
6 mean_accuracy = 0.7984901277584203%
```

## 3.13 Stacked Generalization

David H. Wolpert在1992年发表了SG算法论文，正式提出SG算法。该算法属于集成学习，集成一系列子模型，提高了机器学习的准确率。

### 3.13.1 算法介绍

SG算法，即堆栈泛化。堆栈泛化算法使用一系列的子模型

$$(m_1, m_2, \dots, m_n)$$

每一个子模型分别对数据集进行分类/回归预测，得到结果：

$$(r_1, r_2, \dots, r_n)$$

将每一个样本的  $n$  个预测值组成新的 *stack\_row*，从而生成一个新的数据集，再使用一个新的模型 (*Aggregator Model*) 对新的数据集进行分类/回归预测，完成算法。

### 3.13.2 算法讲解

本节以子模型分别为 *KNN*、*Perceptron*，集成模型为 *Logistic Regression*

#### KNN子模型

该部分代码主体与前面章节介绍的代码一致。

【注意】

- 函数名需要相应的改变
- 此时KNN用于分类任务，故预测值为近邻类别的众数

```
1  #include<stdlib.h>
2  #include<string.h>
3  #include<stdio.h>
4  #include<math.h>
5
6  void QuickSort(double **arr, int L, int R) {
7      int i = L;
8      int j = R;
9      //支点
10     int kk = (L + R) / 2;
11     double pivot = arr[kk][0];
12     //左右两端进行扫描，只要两端还没有交替，就一直扫描
13     while (i <= j) {
14         //寻找直到比支点大的数
15         while (pivot > arr[i][0])
16             {
17                 i++;
18             }
19         //寻找直到比支点小的数
20         while (pivot < arr[j][0])
21             {
22                 j--;
23             }
24         //此时已经分别找到了比支点小的数(右边)、比支点大的数(左边)，它们进行交换
25         if (i <= j) {
26             double *temp = arr[i];
27             arr[i] = arr[j];
28             arr[j] = temp;
29             //double temp2 = arr[i][1];
30             //arr[i][1] = arr[j][1];
31             //arr[j][1] = temp2;
32             i++;
33         }
34     }
```

```

33         j--;
34     }
35 }//上面一个while保证了第一趟排序支点的左边比支点小，支点的右边比支点大了。
36 //“左边”再做排序，直到左边剩下一个数(递归出口)
37 if (L < j)
38 {
39     QuickSort(arr, L, j);
40 }
41 //“右边”再做排序，直到右边剩下一个数(递归出口)
42 if (i < R)
43 {
44     QuickSort(arr, i, R);
45 }
46 }
47 double euclidean_distance(double *row1, double *row2, int col) {
48     double distance = 0;
49     for (int i = 0; i < col - 1; i++) {
50         distance += pow((row1[i] - row2[i]), 2);
51     }
52     return distance;
53 }
54 double* get_neighbors(double **train_data, int train_row, int col, double *test_row,
55 int num_neighbors) {
56     double *neighbors = (double *)malloc(num_neighbors * sizeof(double));
57     double **distances = (double **)malloc(train_row * sizeof(double *));
58     for (int i = 0; i < train_row; i++) {
59         distances[i] = (double *)malloc(2 * sizeof(double));
60         distances[i][0] = euclidean_distance(train_data[i], test_row, col);
61         distances[i][1] = train_data[i][col - 1];
62     }
63     QuickSort(distances, 0, train_row - 1);
64     for (int i = 0; i < num_neighbors; i++) {
65         neighbors[i] = distances[i][1];
66     }
67     return neighbors;
68 }
69 double knn_single_predict(double **train_data, int train_row, int col, double
70 *test_row, int num_neighbors) {
71     double* neighbors = get_neighbors(train_data, train_row, col, test_row,
72 num_neighbors);
73     double result = 0;
74     for (int i = 0; i < num_neighbors; i++) {
75         result += neighbors[i];
76     }
77     return round(result / num_neighbors);
78 }
79 double* knn_predict(double **train, int train_size, double **test, int test_size, int
80 col, int num_neighbors)
81 {
82     double* predictions = (double*)malloc(test_size * sizeof(double));
83     for (int i = 0; i < test_size; i++)
84     {
85         predictions[i] = knn_single_predict(train, train_size, col, test[i],
86 num_neighbors);
87     }
88     return predictions;
89 }

```

## Perceptron子模型

该部分代码主体与前面章节介绍的代码一致。

### 【注意】

- 函数名需要相应的改变

```
1  #include<stdlib.h>
2  #include<stdio.h>
3
4  double perceptron_single_predict(int col, double *array, double *weights) {
5      double activation = weights[0];
6      for (int i = 0; i < col - 1; i++)
7      {
8          activation += weights[i + 1] * array[i];
9      }
10     double output = 0.0;
11     if (activation >= 0.0)
12     {
13         output = 1.0;
14     }
15     else
16     {
17         output = 0.0;
18     }
19     return output;
20 }
21
22 void train_weights(double **data, int col, double *weights, double l_rate, int
n_epoch, int train_size) {
23     for (int i = 0; i < n_epoch; i++)
24     {
25         for (int j = 0; j < train_size; j++)
26         {
27             double yhat = perceptron_single_predict(col, data[j], weights);
28             double err = data[j][col - 1] - yhat;
29             weights[0] += l_rate * err;
30             for (int k = 0; k < col - 1; k++)
31             {
32                 weights[k + 1] += l_rate * err * data[j][k];
33             }
34         }
35     }
36 }
37
38 double* perceptron_predict(double** train, int train_size, double** test, int
test_size, int col, double l_rate, int n_epoch) {
39     double* weights = (double*)malloc(col * sizeof(double));
40     int i;
41     for (i = 0; i < col; i++) {
42         weights[i] = 0.0;
43     }
44     train_weights(train, col, weights, l_rate, n_epoch, train_size); //核心算法执行部分
45     double* predictions = (double*)malloc(test_size * sizeof(double)); //因为test_size
和fold_size一样大
46     for (i = 0; i < test_size; i++) { //因为test_size和fold_size一样大
47         predictions[i] = perceptron_single_predict(col, test[i], weights);
48     }
49     return predictions; //返回对test的预测数组
50 }
```

## 2.3 Logistic Regression集成模型

该部分代码主体与前面章节介绍的代码一致。

【注意】

- 函数名需要相应的改变

```
1 #include<stdlib.h>
2 #include<stdio.h>
3 #include<math.h>
4
5 double stack_predict(int col, double *array, double *coefficients)
6 {
7     double yhat = coefficients[0];
8     int i;
9     for (i = 0; i < col - 1; i++){
10         yhat += coefficients[i + 1] * array[i];
11     }
12     return 1 / (1 + exp(-yhat));
13 }
14
15 void coefficients_sgd(double ** dataset, int col, double *coef, double l_rate, int
n_epoch, int train_size)
16 {
17     for (int i = 0; i < n_epoch; i++)
18     {
19         for (int j = 0; j < train_size; j++)
20         {
21             double yhat = stack_predict(col, dataset[j], coef);
22             double err = dataset[j][col - 1] - yhat;
23             coef[0] += l_rate * err * yhat * (1 - yhat);
24
25             for (int k = 0; k < col - 1; k++)
26             {
27                 coef[k + 1] += l_rate * err * yhat * (1 - yhat) * dataset[j][k];
28             }
29         }
30     }
31 }
```

## 2.4 利用KNN与Perceptron的结果，预测结果

```
1 #include<stdlib.h>
2 #include<stdio.h>
3 #include<math.h>
4
5 extern double* knn_predict(double **train, int train_size, double **test, int
test_size, int col, int num_neighbors);
6 extern double* perceptron_predict(double** train, int train_size, double** test, int
test_size, int col, double l_rate, int n_epoch);
7 extern void coefficients_sgd(double ** dataset, int col, double *coef, double l_rate,
int n_epoch, int train_size);
8 extern double stack_predict(int col, double *array, double *coefficients);
9
10 double* get_test_prediction(double **train, int train_size, double **test, int
test_size, int col, int num_neighbors, double l_rate, int n_epoch)
11 {
12     double* coef = (double*)malloc(3 * sizeof(double));
13     for (int i = 0; i < 3; i++)
14     {
15         coef[i] = 0.0;
```

```

16     }
17     // 训练
18     double* train_knn_predictions = knn_predict(train, train_size, train, train_size,
col, num_neighbors);
19     double* train_perceptron_predictions = perceptron_predict(train, train_size,
train, train_size, col, l_rate, n_epoch);
20     double** new_train = (double **)malloc(train_size * sizeof(double *));
21     for (int i = 0; i < train_size; i++)
22     {
23         new_train[i] = (double *)malloc(3 * sizeof(double));
24     }
25     // 生成新的数据集
26     for (int i = 0; i < train_size; i++)
27     {
28         new_train[i][0] = train_knn_predictions[i];
29         new_train[i][1] = train_perceptron_predictions[i];
30         new_train[i][2] = train[i][col - 1];
31     }
32
33     coefficients_sgd(new_train, 3, coef, l_rate, n_epoch, train_size);
34
35     // 预测
36     double* knn_predictions = knn_predict(train, train_size, test, test_size, col,
num_neighbors);
37     double* perceptron_predictions = perceptron_predict(train, train_size, test,
test_size, col, l_rate, n_epoch);
38
39     double** new_test = (double **)malloc(test_size * sizeof(double *));
40     for (int i = 0; i < test_size; ++i)
41     {
42         new_test[i] = (double *)malloc(2 * sizeof(double));
43     }
44     // 生成新的数据集
45     for (int i = 0; i < test_size; i++)
46     {
47         new_test[i][0] = knn_predictions[i];
48         new_test[i][1] = perceptron_predictions[i];
49     }
50     double* predictions = (double*)malloc(test_size * sizeof(double));
51     for (int i = 0; i < test_size; i++)
52     {
53         predictions[i] = round(stack_predict(3, new_test[i], coef));
54     }
55
56     return predictions; //返回对test的预测数组
57 }

```

### 3.13.3 算法代码

我们现在知道了如何实现**堆栈泛化算法**，那么我们把它应用到[声纳数据集 sonar.csv](https://aistudio.baidu.com/aistudio/datasetdetail/105756/0)

我们给出链接: <https://aistudio.baidu.com/aistudio/datasetdetail/105756/0>

### C语言细节讲解

本节假设您已下载数据集 `sonar.csv`，并且它在当前工作目录中可用。下面我们给出一个完整实例，使用C语言详细讲解每一处细节。我们给出每一个.c文件的所有代码：

### 1) read\_csv.c

该文件代码与前面代码一致，不再重复给出。

### 2) normalize.c

该文件代码与前面代码一致，不再重复给出。

### 3) k\_fold.c

该文件代码与前面代码一致，不再重复给出。

### 4) score.c

该文件代码与前面代码一致，不再重复给出。

### 5) knn\_model.c

```
1  #include <stdlib.h>
2  #include <string.h>
3  #include <stdio.h>
4  #include <math.h>
5
6  void QuickSort(double **arr, int L, int R)
7  {
8      int i = L;
9      int j = R;
10     //支点
11     int kk = (L + R) / 2;
12     double pivot = arr[kk][0];
13     //左右两端进行扫描，只要两端还没有交替，就一直扫描
14     while (i <= j)
15     {
16         //寻找直到比支点大的数
17         while (pivot > arr[i][0])
18         {
19             i++;
20         }
21         //寻找直到比支点小的数
22         while (pivot < arr[j][0])
23         {
24             j--;
25         }
26         //此时已经分别找到了比支点小的数(右边)、比支点大的数(左边)，它们进行交换
27         if (i <= j)
28         {
29             double *temp = arr[i];
30             arr[i] = arr[j];
31             arr[j] = temp;
32             //double temp2 = arr[i][1];
33             //arr[i][1] = arr[j][1];
34             //arr[j][1] = temp2;
35             i++;
36             j--;
37         }
38     } //上面一个while保证了第一趟排序支点的左边比支点小，支点的右边比支点大了。
39     //“左边”再做排序，直到左边剩下一个数(递归出口)
40     if (L < j)
41     {
42         QuickSort(arr, L, j);
43     }
44     //“右边”再做排序，直到右边剩下一个数(递归出口)
45     if (i < R)
46     {
```

```

47     QuickSort(arr, i, R);
48 }
49 }
50 double euclidean_distance(double *row1, double *row2, int col)
51 {
52     double distance = 0;
53     for (int i = 0; i < col - 1; i++)
54     {
55         distance += pow((row1[i] - row2[i]), 2);
56     }
57     return distance;
58 }
59 double *get_neighbors(double **train_data, int train_row, int col, double *test_row,
60 int num_neighbors)
61 {
62     double *neighbors = (double *)malloc(num_neighbors * sizeof(double));
63     double **distances = (double **)malloc(train_row * sizeof(double *));
64     for (int i = 0; i < train_row; i++)
65     {
66         distances[i] = (double *)malloc(2 * sizeof(double));
67         distances[i][0] = euclidean_distance(train_data[i], test_row, col);
68         distances[i][1] = train_data[i][col - 1];
69     }
70     QuickSort(distances, 0, train_row - 1);
71     for (int i = 0; i < num_neighbors; i++)
72     {
73         neighbors[i] = distances[i][1];
74     }
75     return neighbors;
76 }
77 double knn_single_predict(double **train_data, int train_row, int col, double
78 *test_row, int num_neighbors)
79 {
80     double *neighbors = get_neighbors(train_data, train_row, col, test_row,
81 num_neighbors);
82     double result = 0;
83     for (int i = 0; i < num_neighbors; i++)
84     {
85         result += neighbors[i];
86     }
87     return round(result / num_neighbors);
88 }
89 }
90 double *knn_predict(double **train, int train_size, double **test, int test_size, int
91 col, int num_neighbors)
92 {
93     double *predictions = (double *)malloc(test_size * sizeof(double));
94     for (int i = 0; i < test_size; i++)
95     {
96         predictions[i] = knn_single_predict(train, train_size, col, test[i],
97 num_neighbors);
98     }
99     return predictions; //返回对test的预测数组
100 }

```

## 6) perceptron\_model.c

```

1 #include <stdlib.h>
2 #include <stdio.h>
3
4 double perceptron_single_predict(int col, double *array, double *weights)
5 {

```



```

6     double activation = weights[0];
7     for (int i = 0; i < col - 1; i++)
8     {
9         activation += weights[i + 1] * array[i];
10    }
11    double output = 0.0;
12    if (activation >= 0.0)
13    {
14        output = 1.0;
15    }
16    else
17    {
18        output = 0.0;
19    }
20    return output;
21 }
22
23 void train_weights(double **data, int col, double *weights, double l_rate, int
n_epoch, int train_size)
24 {
25     for (int i = 0; i < n_epoch; i++)
26     {
27         for (int j = 0; j < train_size; j++)
28         {
29             double yhat = perceptron_single_predict(col, data[j], weights);
30             double err = data[j][col - 1] - yhat;
31             weights[0] += l_rate * err;
32             for (int k = 0; k < col - 1; k++)
33             {
34                 weights[k + 1] += l_rate * err * data[j][k];
35             }
36         }
37     }
38 }
39
40 double *perceptron_predict(double **train, int train_size, double **test, int
test_size, int col, double l_rate, int n_epoch)
41 {
42     double *weights = (double *)malloc(col * sizeof(double));
43     int i;
44     for (i = 0; i < col; i++)
45     {
46         weights[i] = 0.0;
47     }
48     train_weights(train, col, weights, l_rate, n_epoch, train_size);
49     double *predictions = (double *)malloc(test_size * sizeof(double));
50     for (i = 0; i < test_size; i++)
51     {
52         predictions[i] = perceptron_single_predict(col, test[i], weights);
53     }
54     return predictions;
55 }

```

## 7) stacking\_model.c

```

1  #include <stdlib.h>
2  #include <stdio.h>
3  #include <math.h>
4
5  double stack_predict(int col, double *array, double *coefficients)
6  {
7      double yhat = coefficients[0];

```

```

8     int i;
9     for (i = 0; i < col - 1; i++)
10    {
11        yhat += coefficients[i + 1] * array[i];
12    }
13    return 1 / (1 + exp(-yhat));
14 }
15
16 void coefficients_sgd(double **dataset, int col, double *coef, double l_rate, int
n_epoch, int train_size)
17 {
18     for (int i = 0; i < n_epoch; i++)
19     {
20         for (int j = 0; j < train_size; j++)
21         {
22             double yhat = stack_predict(col, dataset[j], coef);
23             double err = dataset[j][col - 1] - yhat;
24             coef[0] += l_rate * err * yhat * (1 - yhat);
25
26             for (int k = 0; k < col - 1; k++)
27             {
28                 coef[k + 1] += l_rate * err * yhat * (1 - yhat) * dataset[j][k];
29             }
30         }
31     }
32 }

```

## 8) test\_prediction.c

```

1  #include <stdlib.h>
2  #include <stdio.h>
3  #include <math.h>
4
5  extern double *knn_predict(double **train, int train_size, double **test, int
test_size, int col, int num_neighbors);
6  extern double *perceptron_predict(double **train, int train_size, double **test, int
test_size, int col, double l_rate, int n_epoch);
7  extern void coefficients_sgd(double **dataset, int col, double *coef, double l_rate,
int n_epoch, int train_size);
8  extern double stack_predict(int col, double *array, double *coefficients);
9
10 double *get_test_prediction(double **train, int train_size, double **test, int
test_size, int col, int num_neighbors, double l_rate, int n_epoch)
11 {
12     double *coef = (double *)malloc(3 * sizeof(double));
13     for (int i = 0; i < 3; i++)
14     {
15         coef[i] = 0.0;
16     }
17     // 训练
18     double *train_knn_predictions = knn_predict(train, train_size, train, train_size,
col, num_neighbors);
19     double *train_perceptron_predictions = perceptron_predict(train, train_size,
train, train_size, col, l_rate, n_epoch);
20     double **new_train = (double **)malloc(train_size * sizeof(double *));
21     for (int i = 0; i < train_size; i++)
22     {
23         new_train[i] = (double *)malloc(3 * sizeof(double));
24     }
25     // 生成新的数据集
26     for (int i = 0; i < train_size; i++)
27     {

```

```

28     new_train[i][0] = train_knn_predictions[i];
29     new_train[i][1] = train_perceptron_predictions[i];
30     new_train[i][2] = train[i][col - 1];
31 }
32
33 coefficients_sgd(new_train, 3, coef, l_rate, n_epoch, train_size);
34
35 // 预测
36 double *knn_predictions = knn_predict(train, train_size, test, test_size, col,
num_neighbors);
37 double *perceptron_predictions = perceptron_predict(train, train_size, test,
test_size, col, l_rate, n_epoch);
38
39 double **new_test = (double **)malloc(test_size * sizeof(double *));
40 for (int i = 0; i < test_size; ++i)
41 {
42     new_test[i] = (double *)malloc(2 * sizeof(double));
43 }
44 // 生成新的数据集
45 for (int i = 0; i < test_size; i++)
46 {
47     new_test[i][0] = knn_predictions[i];
48     new_test[i][1] = perceptron_predictions[i];
49 }
50 double *predictions = (double *)malloc(test_size * sizeof(double)); //因为
test_size和fold_size一样大
51 for (int i = 0; i < test_size; i++)
52 {
53     predictions[i] = round(stack_predict(3, new_test[i], coef));
54 }
55 return predictions; //返回对test的预测数组
56 }

```

## 9) evaluate.c

```

1  #include <stdlib.h>
2  #include <stdio.h>
3
4  extern double *get_test_prediction(double **train, int train_size, double **test, int
test_size, int col, int num_neighbors, double l_rate, int n_epoch);
5  extern double accuracy_metric(double *actual, double *predicted, int fold_size);
6  extern double ***cross_validation_split(double **dataset, int row, int col, int
n_folds, int fold_size);
7
8  void evaluate_algorithm(double **dataset, int row, int col, int n_folds, int
num_neighbors, double l_rate, int n_epoch)
9  {
10     int fold_size = (int)row / n_folds;
11     double ***split = cross_validation_split(dataset, row, n_folds, fold_size, col);
12     int i, j, k, l;
13     int test_size = fold_size;
14     int train_size = fold_size * (n_folds - 1);
15     double *score = (double *)malloc(n_folds * sizeof(double));
16
17     for (i = 0; i < n_folds; i++)
18     {
19         double ***split_copy = (double ***)malloc(n_folds * sizeof(double **));
20         for (j = 0; j < n_folds; j++)
21         {
22             split_copy[j] = (double **)malloc(fold_size * sizeof(double *));
23             for (k = 0; k < fold_size; k++)
24             {

```

```

25         split_copy[j][k] = (double *)malloc(col * sizeof(double));
26     }
27 }
28 for (j = 0; j < n_folds; j++)
29 {
30     for (k = 0; k < fold_size; k++)
31     {
32         for (l = 0; l < col; l++)
33         {
34             split_copy[j][k][l] = split[j][k][l];
35         }
36     }
37 }
38 double **test_set = (double **)malloc(test_size * sizeof(double *));
39 for (j = 0; j < test_size; j++)
40 {
41     test_set[j] = (double *)malloc(col * sizeof(double));
42     for (k = 0; k < col; k++)
43     {
44         test_set[j][k] = split_copy[i][j][k];
45     }
46 }
47 for (j = i; j < n_folds - 1; j++)
48 {
49     split_copy[j] = split_copy[j + 1];
50 }
51 double **train_set = (double **)malloc(train_size * sizeof(double *));
52 for (k = 0; k < n_folds - 1; k++)
53 {
54     for (l = 0; l < fold_size; l++)
55     {
56         train_set[k * fold_size + l] = (double *)malloc(col *
sizeof(double));
57         train_set[k * fold_size + l] = split_copy[k][l];
58     }
59 }
60 double *predicted = (double *)malloc(test_size * sizeof(double));
61 predicted = get_test_prediction(train_set, train_size, test_set, test_size,
col, num_neighbors, l_rate, n_epoch);
62 double *actual = (double *)malloc(test_size * sizeof(double));
63 for (l = 0; l < test_size; l++)
64 {
65     actual[l] = test_set[l][col - 1];
66 }
67 double acc = accuracy_metric(actual, predicted, test_size);
68 score[i] = acc;
69 printf("scores[%d] = %f%%\n", i, score[i]);
70 free(split_copy);
71 }
72 double total = 0;
73 for (l = 0; l < n_folds; l++)
74 {
75     total += score[l];
76 }
77 printf("mean_accuracy = %f%%\n", total / n_folds);
78 }

```

## 10) main.c

```
1  #include <stdlib.h>
2  #include <string.h>
3  #include <stdio.h>
4  #include <math.h>
5
6  extern int get_row(char *filename);
7  extern int get_col(char *filename);
8  extern void get_two_dimension(char *line, double **dataset, char *filename);
9  extern void normalize_dataset(double **dataset, int row, int col);
10 extern void evaluate_algorithm(double **dataset, int row, int col, int n_folds, int
    num_neighbors, double l_rate, int n_epoch);
11
12 void main()
13 {
14     char filename[] = "sonar.csv";
15     char line[1024];
16     int row = get_row(filename);
17     int col = get_col(filename);
18     double **dataset;
19     dataset = (double **)malloc(row * sizeof(double *));
20     for (int i = 0; i < row; ++i)
21     {
22         dataset[i] = (double *)malloc(col * sizeof(double));
23     }
24     get_two_dimension(line, dataset, filename);
25     normalize_dataset(dataset, row, col);
26     int k_fold = 3;
27     int num_neighbours = 2;
28     double l_rate = 0.01;
29     int n_epoch = 5000;
30     evaluate_algorithm(dataset, row, col, k_fold, num_neighbours, l_rate, n_epoch);
31 }
```

## 11) compile.sh

```
1  gcc main.c normalize.c score.c test_prediction.c stacking_model.c k_fold.c knn_model.c
    evaluate.c read_csv.c perceptron_model.c -o test -lm && ./run
```

**编译&运行:**

```
1  bash compile.sh
```

运算后得到的结果如下:

```
1  Scores[0] = 82.608696%
2  Scores[1] = 79.710145%
3  Scores[2] = 73.913043%
4  mean_accuracy = 78.743961%
```

## Python语言实战

本节同样假设您已经下载数据集，我们使用著名机器学习开源库sklearn高效实现**堆栈泛化算法**，以便您在实战中使用该算法：

```
1  import pandas as pd
2  import numpy as np
3  from sklearn.model_selection import KFold
4  from sklearn.metrics import accuracy_score
```

```

5 from sklearn.linear_model import LogisticRegression
6 from sklearn.neighbors import KNeighborsClassifier
7 from sklearn.linear_model import Perceptron
8 from mlxtend.classifier import StackingClassifier
9
10
11 if __name__ == '__main__':
12     dataset = np.array(pd.read_csv("sonar.csv", sep=',', header=None))
13     k_Cross = KFold(n_splits=3, random_state=0, shuffle=True)
14     index = 0
15     score = np.array([])
16     data, label = dataset[:, :-1], dataset[:, -1]
17     for train_index, test_index in k_Cross.split(dataset):
18         train_data, train_label = data[train_index, :], label[train_index]
19         test_data, test_label = data[test_index, :], label[test_index]
20         model = StackingClassifier(
21             classifiers=[
22                 KNeighborsClassifier(n_neighbors=2),
23                 Perceptron(eta0=0.01, max_iter=5000)
24             ],
25             meta_classifier=LogisticRegression()
26         )
27         model.fit(train_data, train_label)
28         pred = model.predict(test_data)
29         acc = accuracy_score(test_label, pred)
30         score = np.append(score, acc)
31         print('score[{}] = {}'.format(index, acc))
32         index += 1
33     print('mean_accuracy = {}'.format(np.mean(score)))

```

输出结果如下:

```

1 score[0] = 0.8857142857142857%
2 score[1] = 0.7971014492753623%
3 score[2] = 0.8405797101449275%
4 mean_accuracy = 0.8411318150448585%

```

# 四、简易C语言教程

## 一、C语言的初步使用

我们先写一个简单的C语言程序main.c。

```
1  #include <stdio.h>
2
3  int main(void)
4  {
5      int num;
6      num = 0;
7      printf("Hello world!\n");
8      printf("num = %d\n",num);
9      return 0;
10 }
```

```
1  gcc main.c -o run && ./run
```

这样屏幕上打印出了一行“Hello World!”。

我们可以看出这个C程序的基本构成：

1) 指令和头文件：对应上述代码中的 `#include <stdio.h>`，它的作用相当于把stdio.h文件中的所有内容都输入到该行所在的位置。

2) 函数：对应上述代码中的 `int main(void)` 及其大括号下的部分。C程序一定要有一个main函数（特殊情况暂不考虑），它是C程序的核心。你可以创建多个函数，但main函数必须是开始的函数。main前面的int代表的是main函数的返回类型，说明main返回的值是整数。main小括号（()）内的内容代表的是传入函数的信息，其中main(void)代表的是函数不传入任何信息。main大括号（{}）内的内容代表的是函数执行的命令。我们可以在函数里写合适的代码，来完成我们想要的指令。

对于一般的C程序而言，除了main函数之外，可能还会有其他多个函数。这些函数是C程序的构造块，它们相互协调相互作用，共同构成了一个完整的C程序。

接下来，我们来详细看一下main函数里面的内容。它主要由两部分组成：声明和语句。

1) 声明：对应程序里的 `int num;` 这行代码完成了两件事，一是说明在函数中有一个名为num的变量，二是表明num这个变量是一个整数。int是一种数据类型，它代表的是整数；除了int之外，还有其他的数据类型。下表整理了一些比较常用的数据类型：

| 名称     | 数据类型   |
|--------|--------|
| int    | 整型     |
| float  | 浮点型    |
| double | 双精度浮点型 |
| char   | 字符型    |

同时int也是C语言中的关键字。在C语言中，关键字是语言定义的单词，不能用作其它用途，因此不能把int作为变量名；num是一个标识符，是变量的名称。需要注意的是，在C语言中，所有变量都必须先声明才能使用。

2) 语句：语句是C程序的基本构件块。一条语句就相当于一条完整的计算机指令。在上述的例子中，main函数里除了变量声明外，其他的部分都是语句。语句又分为简单语句和复合语句。

这里面比较特殊的是return语句，对应例子中的return 0。它代表的是函数的返回值（执行函数中的代码所得到的结果）。对于没有返回值的函数（void类型），return语句可以不加，而对于其他类型的函数，return语句是需要加的。

return语句可以有多个，可以出现在函数里的任意位置，但每次调用函数都只能有一个return语句被执行，执行之后，return后面的语句就都不会执行了。

需要注意的是，语句和声明的末尾，都需加上“;”，以表明语句\变量的结束。

下面将详细讲解语句的概念和使用规则。

## 二、运算符、表达式、语句

为了理解语句，首先需要了解运算符和表达式是什么。

### 1.运算符

#### 1.1 基本运算符和其他运算符

运算符的作用是用来表示算术运算。数学上的加减乘除符号(+,-,\*,/)在C语言里都属于运算符，并且是基本运算符。除此之外，还有其他的基本运算符。总结后表格如下：

| 基本运算符 | 作用                                 |
|-------|------------------------------------|
| =     | 赋值运算符，把等号右边的值赋给左边变量                |
| +     | 做加法运算，把两侧的值相加                      |
| -     | 做减法运算，把左侧的值减掉右侧的值。除此之外，它也可以让值取相反数。 |
| *     | 做乘法运算                              |
| /     | 做除法运算                              |
| ()    | 括号，和数学中的用法一样                       |

需要特别注意的是，在C语言中，除法运算的结果会因为数据类型的不同而不同！比如，浮点数除法的结果是浮点数，而整数除法的结果是整数。比如我们运行以下代码：

```
1  #include "stdio.h"
2
3  int main(void) {
4      printf("5/4 = %d\n",5/4);
5      return 0;
6  }
```

正确的结果应该是 $5/4 = 1.25$ ，然而程序运行后得到的却是 $5/4 = 1$ ：

这种情况被称为“截断”。在做除法运算的时候，要特别牢记这一点。

除了基本运算符，还有一些其他比较常用的运算符，整理表格如下：



| 运算符       | 作用                                                                    |
|-----------|-----------------------------------------------------------------------|
| %         | 求模运算符，用以求出左侧整数除以右侧整数后得到的余数。                                           |
| ++        | 递增运算符，将其运算对象递增1。该运算符可以放在变量前（前缀模式），也可放在变量后（后缀模式）；两种方式的区别在于递增行为发生的时间不同。 |
| --        | 递减运算符，将其运算对象递减1。该运算符可以放在变量前（前缀模式），也可放在变量后（后缀模式）；两种方式的区别在于递增行为发生的时间不同。 |
| <<br>(<=) | 关系运算符中的小于（小于等于），如果左侧变量的值小于（小于等于）右侧变量的值，则取值为1，否则为0。                    |
| ><br>(>=) | 关系运算符中的大于（大于等于），如果左侧变量的值大于（大于等于）右侧变量的值，则取值为1，否则为0。                    |
| ==        | 关系运算符中的等于，如果左侧变量的值等于右侧变量的值，则取值为1，否则为0。                                |
| !=        | 关系运算符中的不等于，如果左侧变量的值不等于右侧变量的值，则取值为1，否则为0。                              |

对于递增、递减运算符，前缀后缀模式会对代码产生不同的影响，比如我们运行下方代码：

```

1  #include "stdio.h"
2
3  int main(void) {
4      int i = 0;
5      int j = 0;
6      printf("i = %d\n", i++);
7      printf("j = %d\n", ++j);
8      return 0;
9  }
```

得到的结果为：

```

1  i = 0
2  j = 1
```

可以看到打印出来的i，j两个值并不相同。因此我们需要特别注意递增递减运算符前后缀形式的不同。

## 1.2 运算符优先级

不同运算符的优先级是不同的。优先级高的运算符首先被程序执行，优先级相同的从左到右执行（=运算符除外）。因此我们需要了解各个运算符的优先级。整理后得到表格如下：

| 优先级 | 运算符     | 名称          |
|-----|---------|-------------|
| 1   | ()      | 括号          |
| 1   | ++      | 递增运算符（后缀形式） |
| 1   | --      | 递减运算符（后缀形式） |
| 2   | -       | 负号运算符       |
| 2   | ++      | 递增运算符（前缀形式） |
| 2   | --      | 递减运算符（前缀形式） |
| 3   | /       | 除法运算符       |
| 3   | *       | 乘法运算符       |
| 3   | %       | 求模运算符       |
| 4   | +       | 加法运算符       |
| 4   | -       | 减法运算符       |
| 5   | << (>>) | 左移（右移）      |
| 6   | > (>=)  | 大于（大于等于）    |
| 6   | < (<=)  | 小于（小于等于）    |
| 7   | ==      | 等于          |
| 7   | !=      | 不等于         |

## 2.表达式

表达式由运算符和运算对象组成。以下的例子都是表达式：

```
1 4
2 23
3 a*b+c-d/3+100
4 p = 3*a
5 x = q++
6 x > 3
```

C语言中，表达式的一个最重要的特性就是每个表达式都有一个值。这个值是按照运算符优先级的顺序执行获得的。在上述的例子中，前三个都比较清晰。对于有赋值运算符（=）的表达式，其值和赋值运算符左侧变量的值相同。对于有不等式的表达式，如果不等式成立，则值为1，否则为0。

## 3.语句

### 3.1 简单语句

上面我们简单介绍了运算符和表达式，接下来我们来正式介绍语句。

语句是C程序的基本构建块，它有很多类型，其中较为常见的赋值表达式语句——它是由赋值表达式构成的，比如：

```
1 a = 3*5 + 2;
```

函数表达式语句会引起函数调用。比如，调用 `printf()` 函数可以打印结果。

这里需要注意的是：赋值和函数调用都是表达式，这些语句本质上都是表达式语句。

## 3.2 复合语句

以上介绍的基本上都是简单语句，下面将介绍相对复杂一点的复合语句。

复合语句使用花括号括起来的一条或者多条语句，也被称之为块。比较常见的复合语句类型有：循环语句和分支语句。

## 3.3 循环语句

循环语句，顾名思义，就是让程序反复执行一段或者多段命令。循环语句一般都是复合语句，其中比较常见的两种类型有while语句和for语句。

### 3.3.1 while语句

while语句的通用格式如下：

```
1 while(expression){
2     statement
3 }
```

其中expression代表的是表达式，它的作用是判断是否继续循环执行while语句内（花括号内的内容）的语句。如果表达式的值为1，则循环继续（从头执行花括号内的语句）；如果表达式的值为0，则循环结束，往下执行语句（花括号外的语句）。statement代表的是while语句内包含的语句。

通常，我们一般用关系表达式（含关系运算符）作为while循环的判断条件。如果我们想要终止while循环，则必须让while循环内的表达式的值有所变化。比如下面的代码就是错误的，它会导致程序陷入无限循环：

```
1 int index = 2;
2 while(index < 3){
3     printf("hello world!\n");
4 }
```

如果不想让程序想入无限循环，则我们可以写成如下形式：

```
1 int index = 2;
2 while(index < 3){
3     printf("hello world!\n");
4     index++;
5 }
```

### 3.3.2 for语句

for语句的通用格式如下：

```
1 for(exp1;exp2;exp3){
2     statement;
3 }
```

可以看到，for语句小括号里面含有三个表达式。第一个表达式exp1是初始化，它的作用在于初始化计数器的值；第二个表达式是测试条件，如果表达式为假（值为0），则结束循环；第三个表达式是执行更新，在每次循环结束后求值，更新计数。statement代表的是for语句里面包含的语句。

### 3.3.3 嵌套循环

嵌套循环指的是在一个循环内包含另一个循环。嵌套循环常用于按行和列显示数据。下面是一个嵌套循环的例子：

```

1  for(i=0;i<10;i++){
2      for(j=0;j<10;j++){
3          printf("i=%d,j=%d\n",i,j);
4      }
5  }

```

### 3.4 if语句

基本的if语句的通用格式如下：

```

1  if(expression){
2      statement;
3  }

```

if小括号内的表达式expression用于判断。如果表达式为真（1），则执行花括号内的statement（语句）；如果为假，则不执行花括号内的语句。

简单形式的if语句可以让程序选择执行或者不执行一条或者多个语句。除此之外，我们还可以用if和else，让程序在两个语句块中选择其一执行。其格式如下：

```

1  if(expression){
2      statement1;
3  }else{
4      statement2;
5  }

```

这段代码的含义是：如果expression为真，则执行statement1，否则执行statement2。

另外，我们还可以把if语句进行多层嵌套，比如可以写成如下格式：

```

1  if(expression1){
2      statement1;
3  }else if(expression2){
4      statement2;
5  }else{
6      statement3;
7  }

```

这段代码的含义是：如果expression1为真，则执行statement1，否则继续判断expression2的真假——如果expression2为真，则执行statement2，否则执行statement3。

### 3.5 逻辑运算符

在讲完循环语句和if语句后，我们可以继续深入了解一种特殊的运算符——逻辑运算符。顾名思义，它是用来表明逻辑的。

逻辑运算符共有三种：与、或、非。具体如下：

| 逻辑运算符 | 含义 |
|-------|----|
| &&    | 与  |
|       | 或  |
| !     | 非  |

假设exp1和exp2是两个简单的关系表达式，那么：

- 当且仅当exp1和exp2都为真时，exp1&&exp2才为真
- 如果exp1和exp2二者有其一为真，则exp1||exp2为真
- 如果exp1为真，则!exp1为假；如果exp1为假，则!exp1为真

## 三、数组和指针

### 1.数组

数组是由数据类型相同的一系列元素组成的。需要使用数组时，我们首先需要声明数组，告诉编译器数组内有多少元素以及其元素的类型。比如 `int group[10]` 或者 `float group2[7]`，中括号内的数字表示的是数组的长度。除此之外，我们还可以用逗号分隔的列表（用花括号括起来）来初始化数组，比如 `int group[3] = {1,2,3}`。

在初始化数组后，我们可以对数组的某一位置的元素进行赋值。比如这样：

```
1 int a[3];
2 a[0] = 1;
```

其中 `a[0] = 1;` 里的中括号内的数字表示的是数组的下标，指明对应数组下标位置的数组元素。需要注意的是，数组下标是从0开始的，而且存在上界。因此，我们必须保证数组下标在有效范围内。对于 `int a[3];` 这一长度为3的数组来说，它的下标最大值为2。

### 2.多维数组

有时我们想存储矩阵或者表格形式的数据，这时我们就需要用到多维数组。

我们先拿二维数组举例。如果我们想要初始化二维数组，则首先要初始化一维数组。比如我们想要初始化一个二维数组，则可以这样表示：

```
1 int group[2][3] = {
2     {1,2,3},{4,5,6}
3 };
```

可以看到最外层的花括号内含有2个长度为3的数值列表。可以理解为group这个二维数组含有两个长度为3的一维数组。如果花括号内的数值列表的长度小于对应的数组下标，则程序会默认把其他的元素初始化为0。比如：

```
1 int group[2][3] = {
2     {1,2},{4,5,6}
3 };
```

可以看到，最外层花括号内的第一个数值列表的长度只有2，小于对应的数组下标3，因此程序会默认在花括号内第一个数组里数组下标大于1的元素设置为0。也就是等价为：

```
1 int group[2][3] = {
2     {1,2,0},{4,5,6}
3 };
```

初始化时也可以省略内部的花括号，只保留最外层的花括号。只要保证初始化的数值个数正确，初始化的效果与上面相同。但如果初始化的数值不够，则按照先后顺序逐一进行初始化，直到用完所有值。后面没有值初始化的元素会被统一设置为0。比如：

```
1 int group[2][3] = {
2     5,6,7,8
3 };
```

实际上等同于：

```
1 int group[2][3] = {
2     {5,5,7},{8,0,0}
3 };
```

对于二维数组的讨论同样也适用于三维数组以及更多维的数组。比如我们可以声明一个三维数组：

```
1 | int box[10][20][30];
```

可以理解为box这个三维数组中，含有10个大小为20×30的二维数组。

### 3.指针

指针是一个值为内存地址的变量。正如char类型的变量的值为字符，int类型的变量的值为整数，指针变量的值为地址。

假如一个指针变量名为ptr，可以编写如下语句：

```
1 | ptr = &a;
```

其中“&”为地址运算符，表示取右侧向量的内存地址。对于这条语句，我们可以说ptr指向了a。ptr和&a的区别在于ptr是变量，而&a是常量。

要创建指针变量，先要声明指针变量的类型。假设想要把ptr声明为存储int类型变量地址的指针，就需要用到间接运算符“\*”。

假设已知ptr指向b，如：

```
1 | ptr = &b;
```

然后使用间接运算符\*找出存储b中的值，如：

```
1 | val = *ptr;
```

二者相结合相当于下面的语句：

```
1 | val = b;
```

声明指针时，必须指定指针所指向变量的类型。比如：

```
1 | int * pi; //pi是指向int类型变量的指针
2 | char * ch; //ch是指向char类型变量的指针
```

类型说明符表明了指针所指向对象的类型，星号（\*）表明声明的变量是一个指针。