

## 2.1 读取CSV

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。

- 功能——读取csv文件
- 最终返回行数、列数及数据的二维数组
- 字符串作为参数

```
1  double **dataset;
2  int row,col;
3  int get_row(char *filename)//获取行数
4  {
5      char line[1024];
6      int i = 0;
7      FILE* stream = fopen(filename, "r");
8      while(fgets(line, 1024, stream)){
9          i++;
10     }
11     fclose(stream);
12     return i;
13 }
14
15 int get_col(char *filename)//获取列数
16 {
17     char line[1024];
18     int i = 0;
19     FILE* stream = fopen(filename, "r");
20     fgets(line, 1024, stream);
21     char* token = strtok(line, ",");
22     while(token){
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,
40 ",\n")){
41             data[i][j] = atof(tok);//转换成浮点数
42             //字符串拆分操作
43             i++;
44         }
```

```

43     free(tmp);
44 }
45 fclose(stream); //文件打开后要进行关闭操作
46 }
47

```

我们使用data.csv文件，内容如下：

```

1  1  12.2  12.5  11.1
2  2.5 555.2 121.4  2.1

```

调用函数：

```

1  int main()
2  {
3      char filename[] = "data.csv";
4      char line[1024];
5      double **data;
6      int row, col;
7      row = get_row(filename);
8      col = get_col(filename);
9      data = (double **)malloc(row * sizeof(int *));
10     for (int i = 0; i < row; ++i){
11         data[i] = (double *)malloc(col * sizeof(double));
12     } //动态申请二维数组
13     get_two_dimension(line, data, filename);
14     printf("row = %d\n", row);
15     printf("col = %d\n", col);
16
17     int i, j;
18     for(i=0; i<row; i++){
19         for(j=0; j<col; j++){
20             printf("%f\t", data[i][j]);
21         }
22         printf("\n");
23     }
24 }

```

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 计算RMSE

衡量观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。

- 功能——计算按规定划分的测试集的RMSE
- 最终输出RMSE值

```

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

```

我们利用如下数据：

```

1 float actual[] = {1,2,3,4};
2 float pre[] = {4,3,2,1};
3 printf("%f",rmse_metric(actual,pre,3));

```

得到结果如下：

```

1 3.316625

```

## 2.3 score.c

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。

功能——计算准确率

输入真实值一维数组，预测值一维数组，交叉验证fold的长度

输出准确率值

```

1 float accuracy_metric(float *actual, float *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 / (float)len)*100.0;
12 }

```

我们利用如下数据集：

```

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

```

得到如下结果:

```

1 75.0000

```

## 2.4 划分数据为k折

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。

- 功能——划分数据为k折
- 输入数据集的二维数组, 行数, 交叉验证折数, 交叉验证fold的长度
- 输出划分后的三维数组

```

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

```

```

34         split[i]=fold;
35     }
36     return split;
37 }
38

```

我们运行以下代码：

```

1  int main()
2  {
3      double data[6][2];
4      double *data_ptr[6];
5      double *** split;
6      for(int i=0;i<6;i++)
7      {
8
9          for(int j=0;j<2;j++)
10         {
11             data[i][j]=i+j;
12             //printf("%f",actual[i][j]);
13         }
14         data_ptr[i] = data[i];
15     };
16
17     split = cross_validation_split(data_ptr,6,3,2);
18     printf("%f",split[0][0][1]);
19 }
20

```

结果如下：

```

1 | 6.0000

```

## 2.5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。

```

1  void normalize_dataset(float **dataset,int row, int col)
2  {
3      // 先 对列循环
4      float 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 }

```

使用以下数据

1	0.000000	1.000000	2.000000	3.000000	4.000000	5.000000
	10.000000	11.000000	12.000000	13.000000	14.000000	15.000000
	20.000000	21.000000	22.000000	23.000000	24.000000	25.000000
	30.000000	31.000000	32.000000	33.000000	34.000000	35.000000
	40.000000	41.000000	42.000000	43.000000	44.000000	45.000000
	50.000000	51.000000	52.000000	53.000000	54.000000	55.000000
	60.000000	61.000000	62.000000	63.000000	64.000000	65.000000
	70.000000	71.000000	72.000000	73.000000	74.000000	75.000000
	80.000000	81.000000	82.000000	83.000000	84.000000	85.000000
	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
	0.111111	0.111111	0.111111	0.111111	0.111111	0.111111
	0.222222	0.222222	0.222222	0.222222	0.222222	0.222222
	0.333333	0.333333	0.333333	0.333333	0.333333	0.333333
	0.444444	0.444444	0.444444	0.444444	0.444444	0.444444
	0.555556	0.555556	0.555556	0.555556	0.555556	0.555556
	0.666667	0.666667	0.666667	0.666667	0.666667	0.666667
	0.777778	0.777778	0.777778	0.777778	0.777778	0.777778
	0.888889	0.888889	0.888889	0.888889	0.888889	0.888889
	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000

## 2.6测试预测函数

将训练集、测试集、学习率，epoch数，交叉验证fold的长度输入函数，并用函数中的模型框架进行对训练集训练后得到模型后，对预测集进行预测分类，并返回预测结果。

函数如下：

```

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

```

## 2.7利用k-fold交叉验证计算预测准确率

输入训练集、测试集，学习率，epoch数，交叉验证折数，交叉验证fold的长度等，并用函数中的模型进行预测，输出最终准确率。

函数实例：