

First set the environment & load csv files, I load two data sets, one is sample from the slide and one is the iris flower data set.

Code: `mydata<-read.csv("sample.csv")`

Code: `myflower<-read.csv("raw data.csv")`

## 1. Implement chi-merge

I use the function `chiSq` from package “discretization” to calculate the chi-square value between two clusters. So we need to install and import this package.

Code: `install.packages("discretization")`

Code: `library(discretization)`

Below is my `chiMerge` algorithm, I have attached a text file in my homework report for readability:

```
1 chi_merge <- function(ATTR, FREQU, THRESHOLD){
2   contingency_table <- table(ATTR, FREQU)
3   contingency_table <- as.data.frame.matrix(contingency_table)
4
5   print(contingency_table)
6   #calculate chi value and create a table
7   n = nrow(contingency_table)
8   for (i in 1:(n-1)) {
9     chi_val = chiSq(contingency_table[i:(i+1), 1:ncol(contingency_table)])
10    if(grepl("-", rownames(contingency_table)[i])){
11      lower_bound = as.numeric(strsplit(rownames(contingency_table)[i], "-")[[1]][1])
12    }else{
13      lower_bound = as.numeric(rownames(contingency_table)[i])
14    }
15    if(grepl("-", rownames(contingency_table)[(i+1)])){
16      upper_bound = as.numeric(strsplit(rownames(contingency_table)[(i+1)], "-")[[1]][2])
17    }else{
18      upper_bound = as.numeric(rownames(contingency_table)[(i+1)])
19    }
20    if (!exists("chi_table")) chi_table <- c(lower_bound, upper_bound, chi_val)
21    else chi_table <- rbind(chi_table, c(lower_bound, upper_bound, chi_val))
22  }
23  colnames(chi_table) <- c("lower_bound", "upper_bound", "chi_val")
24  chi_table = as.data.frame(chi_table)
25  min_chi = min(chi_table$chi_val)
26  p_value = qchisq(1 - threshold, ncol(contingency_table)-1)
27  print(chi_table)
28  #merge the min chi_val until the threshold is reached
29  while(min_chi < p_value){
30    n = nrow(chi_table)
31    # generate new contingency_table
32    for(i in 1:n){
33      if(chi_table[i,3] == min_chi){
34        if((i==n) || chi_table[(i+1),3] != min_chi){
35          if(!exists("index_lower")) index_lower = i
36          index_upper = (i+1)
37          interval = paste(chi_table[index_lower,1], chi_table[i,2], sep="-")
38          rownames(contingency_table)[index_lower] <- interval
39          col_num = ncol(contingency_table)
40          while(index_upper != index_lower){
41            for (j in 1:col_num){
42              contingency_table[(index_upper-1),j] = contingency_table[(index_upper-1),j] + contingency_table[index_upper,j]
43            }
44            if(!exists("row_index_to_delete")){
45              row_index_to_delete <- index_upper
46            }
47            else{
48              row_index_to_delete <- c(row_index_to_delete, index_upper)
49            }
50            index_upper = index_upper + 1
51          }
52        }
53        # we need to do cumulative merging if two adjacent chi_val are min value
54        else{
55          if(!exists("index_lower")) index_lower = i
56          next
57        }
58      }
59      if(exists("index_lower")) rm(index_lower)
60    }
61    contingency_table <- contingency_table[-row_index_to_delete,]
62    rm(row_index_to_delete)
63    print(contingency_table)
64    # calculate chi value and generate chi_table
65    n = nrow(contingency_table)
66    rm(chi_table)
67    for (i in 1:(n-1)){
68      chi_val = chiSq(contingency_table[i:(i+1), 1:ncol(contingency_table)])
69      if(grepl("-", rownames(contingency_table)[i])){
70        lower_bound = as.numeric(strsplit(rownames(contingency_table)[i], "-")[[1]][1])
71      }else{
72        lower_bound = as.numeric(rownames(contingency_table)[i])
73      }
74      if(grepl("-", rownames(contingency_table)[(i+1)])){
75        upper_bound = as.numeric(strsplit(rownames(contingency_table)[(i+1)], "-")[[1]][2])
76      }else{
```

```

77     upper_bound = as.numeric(rownames(contingency_table)[(i+1)])
78   }
79   if (exists("chi_table")) chi_table <- c(lower_bound,upper_bound,chi_val)
80   else chi_table <- rbind(chi_table,c(lower_bound,upper_bound,chi_val))
81 }
82 colnames(chi_table) <- c("lower_bound","upper_bound","chi_val")
83 chi_table = as.data.frame(chi_table)
84 print(chi_table)
85 min_chi = min(chi_table$chi_val)
86 }
87 return(contingency_table)
88 }
89
90

```

The `chi_merge` function takes two vectors as input , one is **attr** , the data we want for discretization, one is **freque** , the class frequency of the corresponding data . The function also need a **threshold** parameter , which is needed to calculate the corresponding p values.

The output of this function is a contingency table(I print the contingency table and `chi_val` table every calculation in my function for debugging and better illustration) , the row names specify the cluster interval (inclusively) , the columns are the class frequency of this specific data interval. For example:

```

      Iris-setosa Iris-versicolor Iris-virginica
4.3-4.8          16              0              0
4.9              4              1              1
5-5.4           25              5              0
5.5-5.7          4             15              2
5.8-6.2          1             15             10
6.3-7            0             14             25
7.1-7.9          0              0             12
> |

```

We take the first line for illustration , it says , in 4.3-4.8 interval (inclusively) , there are 16 instances of class Iris-setosa , 0 instance of Iris-versicolor and 0 instance of class Iris-virginica.

Basically the function `chi_merge` do two things , first , generate the **contingency table** of **attr** and **freque** . Then generate the corresponding **chi\_value** table . which is needed to update our original contingency table. I will use the example in the slides for demonstration:

Step 1 : Generate the contingency table of **attr** F and **freque** K

F	K=1	K=2
1	1	0
3	0	1
7	1	0
8	1	0
9	1	0
11	0	1
23	0	1
37	1	0
39	0	1

45	1	0
46	1	0
59	1	0

Step 2: Generate the chi\_val table between clusters in the contingency table. We need to calculate the chi\_square value for clusters pairs in the contingency table.

Chi_table	Lower_bound	Upper_bound	Chi_val
	1	3	1.9996
	3	7	1.9996
	7	8	0.0000
	8	9	0.0000
	9	11	1.9996
	11	23	0.0000
	23	37	1.9996
	37	39	1.9996
	39	45	1.9996
	45	46	0.0000
	46	59	0.0000

In this table **lower bound** means the lower bound of the first cluster , **upper bound** means the upper bound of the second cluster. **Chi\_val** is the chi value of these two clusters. We will then find the min value of chi-square.

Step 3:If the min value of chi-square is below corresponding p\_value, merge these two clusters into one cluster. And modify the original contingency table.

Before:

F	K=1	K=2
1	1	0
3	0	1
7-9	1+1+1	0+0+0
8	1	0
9	1	0
11-23	0+0	1+1
23	0	1
37	1	0
39	0	1
45-59	1+1+1	0+0+0
46	1	0
59	1	0

After:

F	K=1	K=2
1	1	0
3	0	1
7-9	3	0
11-23	0	2

37	1	0
39	0	1
45-59	3	0

Step 4: Generate new chi\_table for our updated contingency table.

Chi_table	Lower_bound	Upper_bound	Chi_val
	1	3	1.9996
	3	9	3.999333
	7	23	4.999567
	11	37	2.999500
	37	39	1.999600
	39	59	3.999333

Step 5: Then we will do Step 3 and Step 4 again and again until the min chi\_val is above our p\_value , which means the adjacent clusters are different enough and we do not need to merge them into one cluster.

So much for the illustration of my chi\_merge algorithm.

## 2. Discretize sepal length in cm using 5 different thresholds using chi-merge (0.10 0.05 0.025 0.02 0.01)

Code: `chi_merge(myflower$sepal.length,myflower$class,0.1)`

Output:

```

      Iris-setosa Iris-versicolor Iris-virginica
4.3-4.8          16              0              0
4.9              4              1              1
5-5.4           25              5              0
5.5-5.7          4             15              2
5.8-6.2          1             15             10
6.3-7            0             14             25
7.1-7.9          0              0             12
> |

```

Code: `chi_merge(myflower$sepal.length,myflower$class,0.05)`

Output:

```

      Iris-setosa Iris-versicolor Iris-virginica
4.3-5.4          45              6              1
5.5-5.7          4             15              2
5.8-7            1             29             35
7.1-7.9          0              0             12
> |

```

Code: `chi_merge(myflower$sepal.length,myflower$class,0.025)`

Output:

```
      Iris-setosa Iris-versicolor Iris-virginica
4.3-5.4          45              6              1
5.5-5.7           4             15              2
5.8-7             1             29             35
7.1-7.9           0              0             12
> |
```

Code: `chi_merge(myflower$sepal.length,myflower$class,0.02)`

Output:

```
      Iris-setosa Iris-versicolor Iris-virginica
4.3-5.4          45              6              1
5.5-5.7           4             15              2
5.8-7             1             29             35
7.1-7.9           0              0             12
> |
```

---

Code: `chi_merge(myflower$sepal.length,myflower$class,0.01)`

Output:

```
      Iris-setosa Iris-versicolor Iris-virginica
4.3-5.4          45              6              1
5.5-5.7           4             15              2
5.8-7.9           1             29             47
> |
```

### 3. Discretize petal length in cm using 5 different thresholds using chi-merge (0.10 0.05 0.025 0.02 0.01)

Code: `chi_merge(myflower$petal.length,myflower$class,0.1)`

Output:

```
      Iris-setosa Iris-versicolor Iris-virginica
1-1.9           50              0              0
3-4.7            0             44              1
4.8-5.1           0              6             15
5.2-6.9           0              0             34
>
```

Code: `chi_merge(myflower$petal.length,myflower$class,0.05)`

Output:

```
      Iris-setosa Iris-versicolor Iris-virginica
1-1.9           50              0              0
3-4.7            0             44              1
4.8-5.1           0              6             15
5.2-6.9           0              0             34
>
```

Code: `chi_merge(myflower$petal.length,myflower$class,0.025)`

Output:

```
      Iris-setosa Iris-versicolor Iris-virginica
1-1.9           50              0              0
3-4.7            0             44              1
4.8-5.1          0              6             15
5.2-6.9          0              0             34
> |
```

Code: `chi_merge(myflower$petal.length,myflower$class,0.02)`

Output:

```
      Iris-setosa Iris-versicolor Iris-virginica
1-1.9           50              0              0
3-4.7            0             44              1
4.8-5.1          0              6             15
5.2-6.9          0              0             34
> |
```

Code: `chi_merge(myflower$petal.length,myflower$class,0.01)`

Output:

```
      Iris-setosa Iris-versicolor Iris-virginica
1-1.9           50              0              0
3-4.7            0             44              1
4.8-5.1          0              6             15
5.2-6.9          0              0             34
> |
```