

Group-4-Project-2

Yi Qian; Kailing Huang; Tarek Hatata; Shusen Wu; Bohui NI

R function:

1. `as.data.frame(x, row.names = NULL, optional = FALSE, ...)`
S3 method for character
`as.data.frame(x, ..., stringsAsFactors = default.stringsAsFactors())`
S3 method for list
`as.data.frame(x, row.names = NULL, optional = FALSE, ..., cut.names = FALSE, col.names = names(x), fix.empty.names = TRUE, stringsAsFactors = default.stringsAsFactors())`
S3 method for matrix
`as.data.frame(x, row.names = NULL, optional = FALSE, ..., stringsAsFactors = default.stringsAsFactors())`
Functions to check if an object is a data frame, or coerce it if possible.
2. `Supply(X, FUN, ..., simplify = TRUE, USE.NAMES = TRUE)`
Supply is equivalent to `sapply`, except that it preserves the dimension and dimension names of the argument X. It also preserves the dimension of results of the function FUN. It is intended for application to results e.g. of a call to `by`. Lapply is an analog to lapply insofar as it does not try to simplify the resulting list of results of FUN.
3. `as.character(x, ...)`
Create or test for objects of type "character".
4. `kmeans(x, centers, iter.max = 10, nstart = 1, algorithm = c("Hartigan-Wong", "Lloyd", "Forgy", "MacQueen"), trace=FALSE)`
The data given by x are clustered by the k-means method, which aims to partition the points into k groups such that the sum of squares from points to the assigned cluster centres is minimized. At the minimum, all cluster centres are at the mean of their Voronoi sets (the set of data points which are nearest to the cluster centre).
5. `knn(train, test, cl, k = 1, l = 0, prob = FALSE, use.all = TRUE)`
k-nearest neighbour classification for test set from training set. For each row of the test set, the k nearest (in Euclidean distance) training set vectors are found, and the classification is decided by majority vote, with ties broken at random. If there are ties for the kth nearest vector, all candidates are included in the vote.
6. `CrossTable(x, y, digits=3, max.width = 5, expected=FALSE, prop.r=TRUE, prop.c=TRUE, prop.t=TRUE, prop.chisq=TRUE, chisq = FALSE, fisher=FALSE, mcnemar=FALSE, resid=FALSE, sresid=FALSE, asresid=FALSE, missing.include=FALSE, format=c("SAS","SPSS"), dnn = NULL, ...)`
An implementation of a cross-tabulation function with output similar to S-Plus `crosstabs()` and SAS Proc Freq (or SPSS format) with Chi-square, Fisher and McNemar tests of the independence of all table factors.
7. `vegdist(x, method="bray", binary=FALSE, diag=FALSE, upper=FALSE, na.rm = FALSE, ...)`
The function computes dissimilarity indices that are useful for or popular with community ecologists. All indices use quantitative data, although they would be named by the corresponding binary index, but you can calculate the binary index using an appropriate argument. If you do not find your favourite index here, you can see if it can be implemented using `designdist`.

[illegible]

Knn :

RStudio interface showing R code execution and environment. The code defines a k-NN model and evaluates its performance. The console output shows the cell contents and a confusion matrix.

```
# RStudio
# Source on Save
# Run
# Source

# knn
# get train label
knn_train_label = mrftrain.df[,1]
# remove label
knn_train_data = mrftrain.df[,2:23]
# return prediction of test set
knn_pred = knn(train=knn_train_data, test=knn_test_data, cl=knn_train_label, k=3)
# evaluating knn
knn_test_label = mrfest.df[,1]
# cross table
crosstable(x=knn_test_label, y=knn_pred, prop.chisq=FALSE)
# hierarchical clustering
# (no levels)
```

cell contents

	N / Row Total	N / Col Total	N / Table Total

Total observations in Table: 2438

knn_test_label	knn_pred	0	1	Row Total
0	0	1216	0	1216
	1	1.000	0.000	0.499
	1.000	0.000		
	0.499	0.000		
1	0	0	1222	1222
	1	0.000	1.000	0.501
	0.000	1.000		
	0.000	0.501		
column Total		1216	1222	2438
	0.499	0.501		

RStudio interface showing R code execution and environment. The code defines a k-NN model and evaluates its performance. The console output shows the cell contents and a confusion matrix.

```
# RStudio
# Source on Save
# Run
# Source

# knn
# get train label
knn_train_label = mrftrain.df[,1]
# remove label
knn_train_data = mrftrain.df[,2:23]
# return prediction of test set
knn_pred = knn(train=knn_train_data, test=knn_test_data, cl=knn_train_label, k=5)
# evaluating knn
knn_test_label = mrfest.df[,1]
# cross table
crosstable(x=knn_test_label, y=knn_pred, prop.chisq=FALSE)
# hierarchical clustering
# (no levels)
```

cell contents

	N / Row Total	N / Col Total	N / Table Total

Total observations in Table: 2438

knn_test_label	knn_pred	0	1	Row Total
0	0	1216	0	1216
	1	1.000	0.000	0.499
	1.000	0.000		
	0.499	0.000		
1	0	0	1222	1222
	1	0.000	1.000	0.501
	0.000	1.000		
	0.000	0.501		
column Total		1216	1222	2438
	0.499	0.501		

RStudio interface showing R code execution and environment. The code defines a k-NN model and evaluates its performance. The console output shows the cell contents and a confusion matrix.

```
# RStudio
# Source on Save
# Run
# Source

# knn
# get train label
knn_train_label = mrftrain.df[,1]
# remove label
knn_train_data = mrftrain.df[,2:23]
# return prediction of test set
knn_pred = knn(train=knn_train_data, test=knn_test_data, cl=knn_train_label, k=7)
# evaluating knn
knn_test_label = mrfest.df[,1]
# cross table
crosstable(x=knn_test_label, y=knn_pred, prop.chisq=FALSE)
# hierarchical clustering
library(vegan)
# (no levels)
```

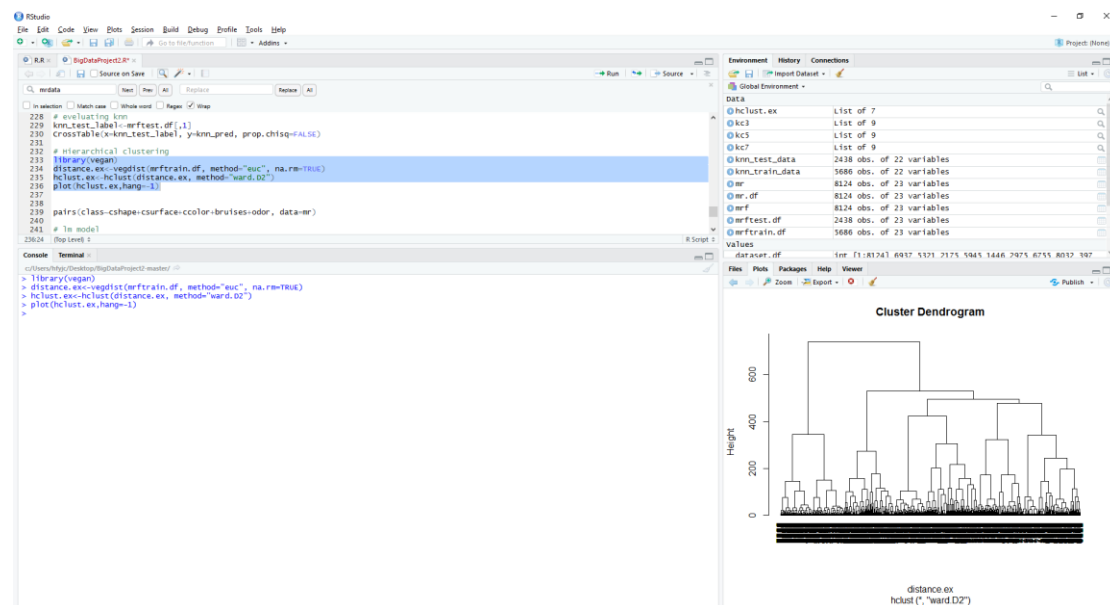
cell contents

	N / Row Total	N / Col Total	N / Table Total

Total observations in Table: 2438

knn_test_label	knn_pred	0	1	Row Total
0	0	1212	4	1216
	1	0.997	0.003	0.499
	1.000	0.003		
	0.497	0.002		
1	0	0	1222	1222
	1	0.000	1.000	0.501
	0.000	0.997		
	0.000	0.501		
column Total		1212	1226	2438
	0.497	0.503		

Hierarchical clustering

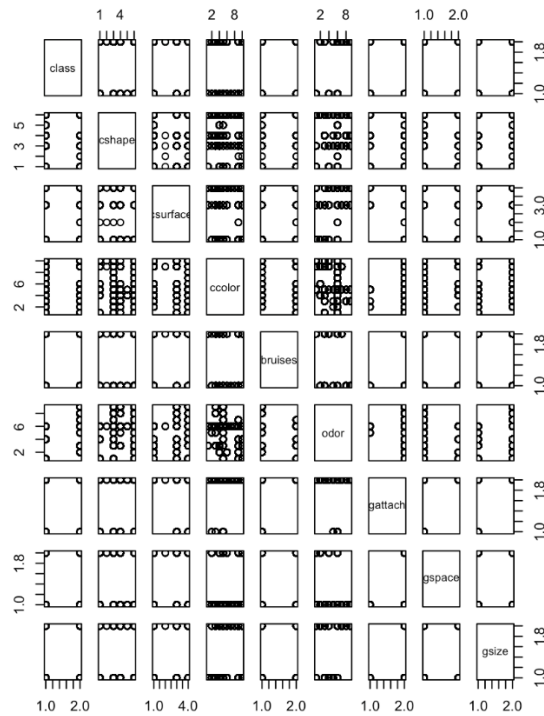


Discussion of results from the experiments that you run in part 3b above. Which clustering version gives the best results? What seems to be the best number of clusters for each method?

A clustering of the samples into 3 and 5 gives the best results, because k-nearest neighbor classification for test set from training set have a good accuracy, they all have 100% accuracy rate.

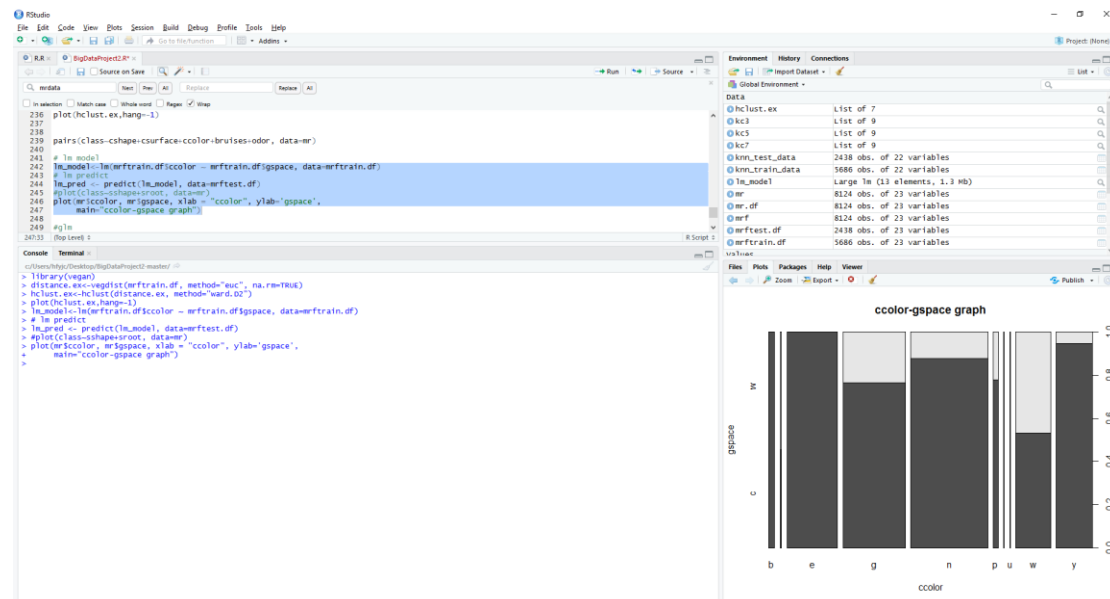
In our opinion, the KNN give us best result. Undefinedly, the other method also can give us a good result. But the KNN accuracy rate close to 100%.

Table from 4 above

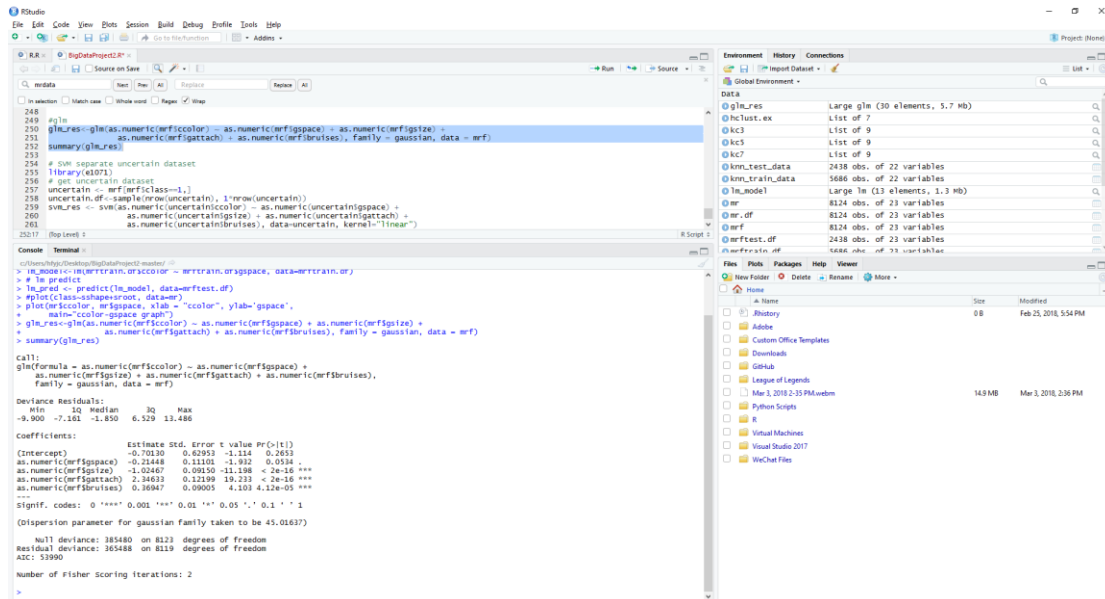


From above pair picture, we get some class could be compare and utilized.

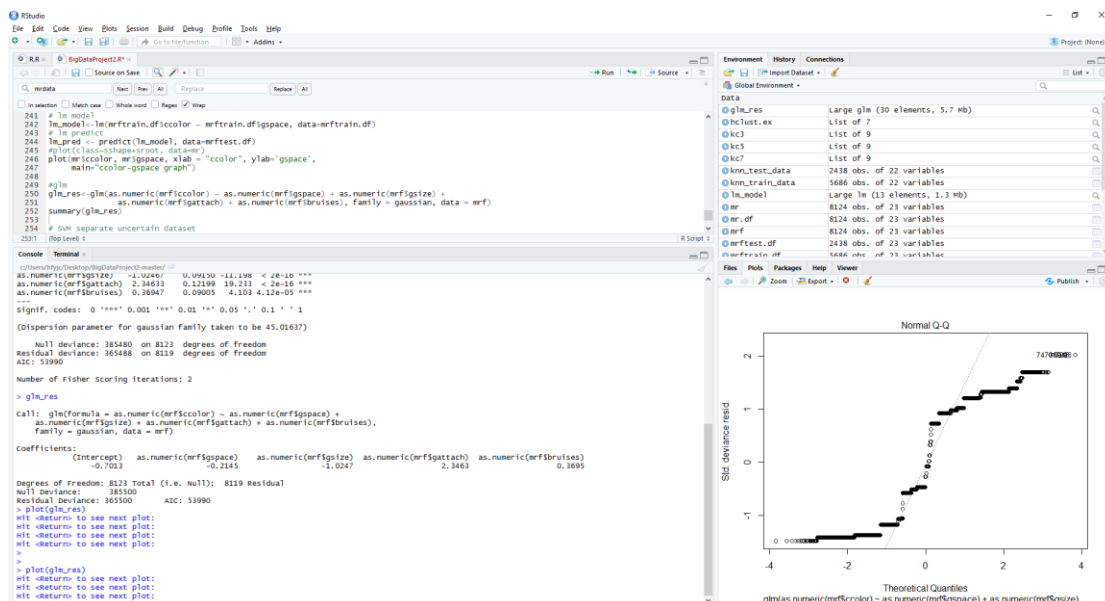
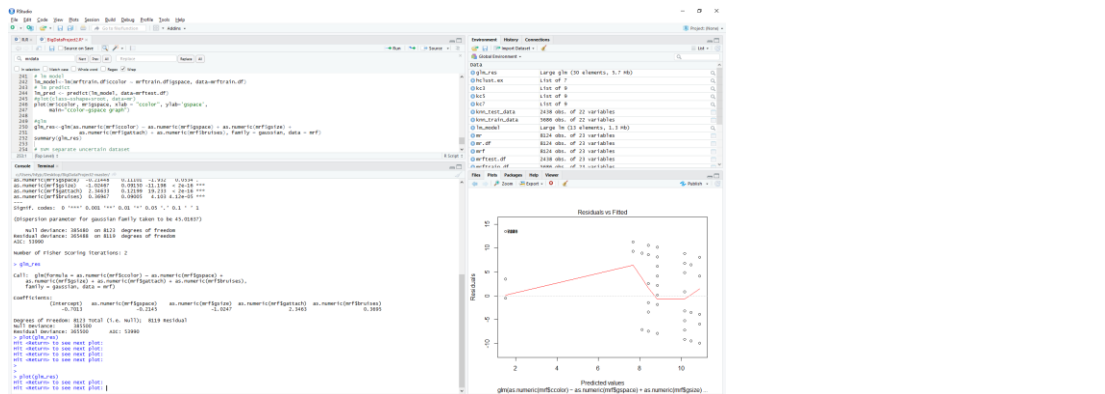
In lm model, we pick color and space:

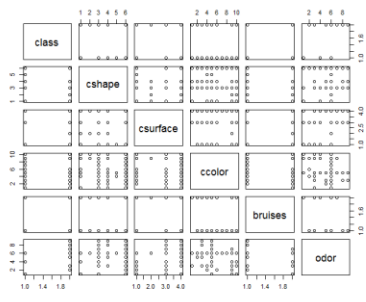
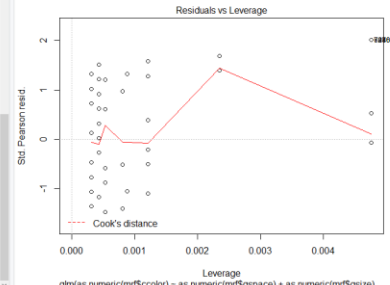
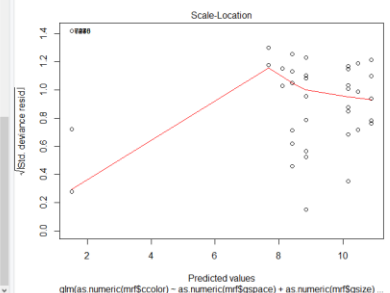


In glm model, we can pick more class: color, space, size, attach, bruises:



Plots with discussion of results: which plot helped you understand the data best and why?





The above picture is best plot for us to understand the data. Because, from this picture, they show all relationship between different class. From this picture, we can find which class can help us to distinguish the kind of mushroom. We compare the picture, and try to find which two class have linear relationship, and then use those data further.

Analysis of what this project helped you learn about data science, e.g., the exploration of data which is what you have been doing – 1 point. You must argue persuasively.

First, we learn some new function and how to use it. Based on those function, we can get some graph and data what we need for analyses. Then the most important things is know how to analyses those graph, and data. We have those data, then convert to different table or graph, then we make assumption from that information. And then we try to prove it or find some new viewpoint. We work this as a loop, so that improve our ability of integrating and analyzing data.