

STAT 501 Exam 2

Yifan Zhu

April 28, 2018

1. We check multivariate normality for each of the 3 cultivars using `testnormality` and `mvnorm.etest`.

```
1 # read data
2 wine <- read.table("wine.dat", head = F, sep = ",",
3                   col.names = c("Cultivar", "Alcohol", "Malic acid", "Ash", "Alkalinity",
4                                "Magnesium", "Total phenols", "Flavanoids",
5                                "Nonflavanoid phenols", "Proanthocyanins", "Color intensity", "Hue",
6                                "OD280/OD315", "Proline"))
7 wine$Cultivar <- as.factor(wine$Cultivar)
8 # check multivariate normality for each cultivar
9 source("testnormality.R")
10 library(energy)
11 library(dplyr)
12 wine %>% group_by(Cultivar) %>% do(data.frame(testnormality = testnormality(., -1),
13                                              energytest = mvnorm.etest(., -1, R = 999)$p.value))
```

The result is shown in a table, each row is a group and the `testnormality` and `energytest` are p-values of the tests.

	Cultivar	testnormality	energytest
	<fct>	<dbl>	<dbl>
1	1	0.656	0.0280
2	2	0.000000712	0.
3	3	0.794	0.456

From the result we can see, the p-values from both tests for Cultivar 3 are large, thus we conclude the multivariate normality holds for Cultivar 3. The p-values from both tests are small for Cultivar 2, and we reject the null hypothesis and conclude that the multivariate normality does not hold for Cultivar 2. For Cultivar 1, the p-values from `testnormality` is large, while the p-value from `energytest` is small. But since the p-value from `testnormality` is pretty big, we still conclude that multivariate normality holds for Cultivar 1.

2. Before doing the clustering, we first standardize the data.

```
1 # cultivar information
2 cultivar <- wine$Cultivar
3
4 # standardize data
5 wine.sc <- scale(wine[, -1])
```

- (a) Hierarchical clustering with average linkage:

The hierarchical tree is shown in Figure 1. Then we take number of group to be 3 and display the clustering result using `ggandrews` in Figure 2.

```

1 source("ggandrews.R")
2 # Hierarchical clustering with average linkage
3 hc <- hclust(dist(wine.sc), method = "average")
4 plot(hc, main = "")
5
6 # display using ggandrews
7 ggandrews(data.frame(cutree(hc, k = 3), wine.sc), clr = 1, return_value = F)

```

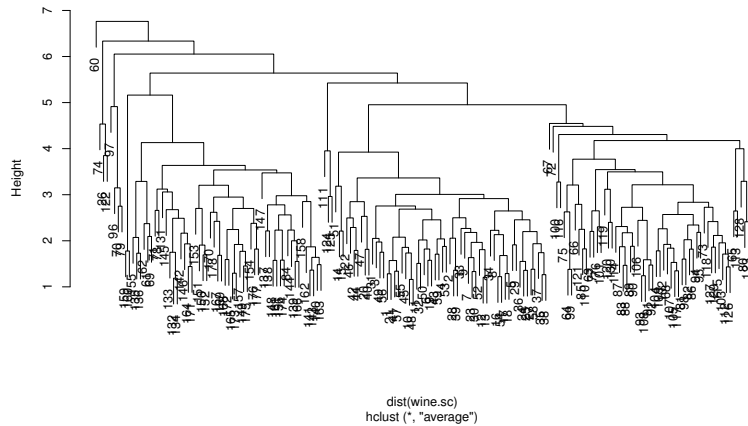


Figure 1: Hierarchical clustering tree

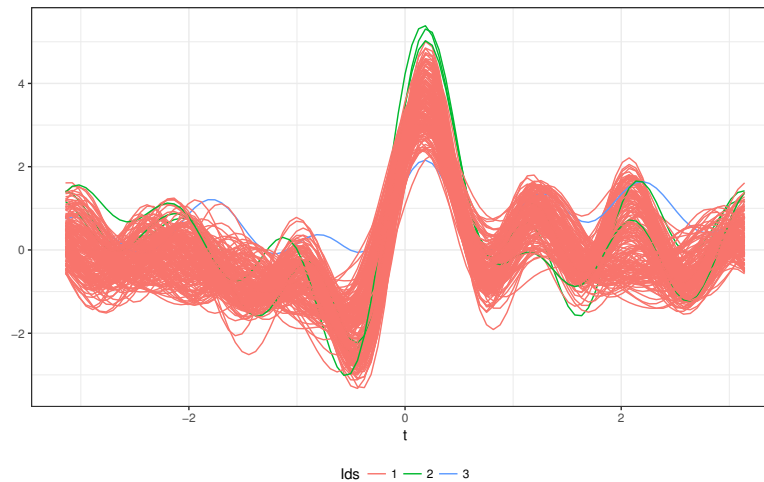


Figure 2: Hierarchical clustering with average linkage result, $K = 3$

(b) K-means clustering:

Here we use 2 different initializations. One is using the result from hierarchical clustering with average linkage in (a), another is random initialization. The results are displayed using ggandrews in Figure 3 and Figure 4.

```

1 # k-means initialized with hc
2 kmnsinithcl <- function(x.data, nclus, ncut = nclus, hcl.tree)
3 {

```

```

4  x.hcl <- hcl.tree
5  x.cl <- cutree(x.hcl, k = ncut)
6  data.x <- data.frame(x.data, cl = x.cl)
7  means <- aggregate(. ~ cl, data = data.x, FUN = mean)
8  return(kmeans(x.data, centers = means[, -1]))
9  }
10
11 km <- kmnsinithcl(wine.sc, nclus = 3, ncut = 3, hcl.tree = hc)
12
13 # display using ggandrews
14 ggandrews(data.frame(km$cluster, wine.sc), clr = 1, return_value = F)
15
16 # k-means with random initialization
17 km.r <- kmeans(wine.sc, centers = 3, nstart = 10000)
18
19 # display using ggandrews
20 ggandrews(data.frame(km.r$cluster, wine.sc), clr = 1, return_value = F)

```

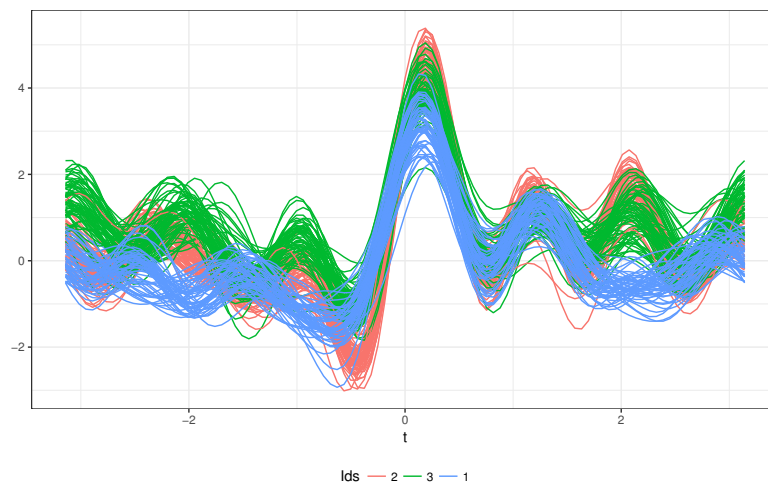


Figure 3: K-means clustering result initialized with hierarchical clustering result, $K = 3$

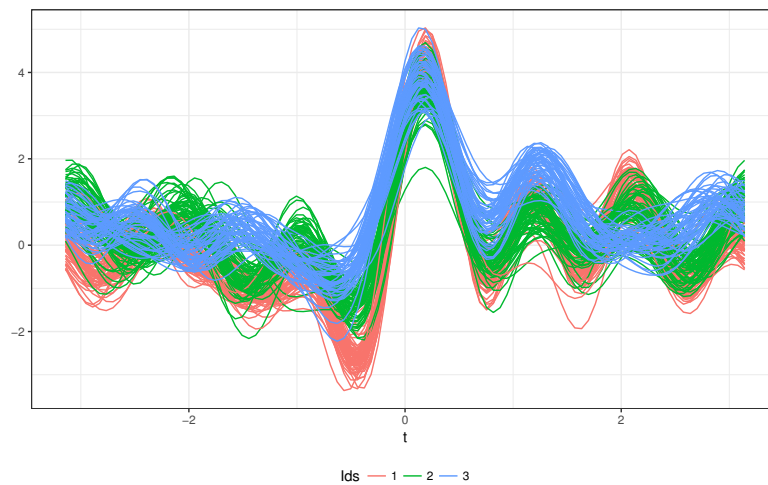


Figure 4: K-means clustering result with random initialization, $K = 3$

(c) Model based clustering:

We use BIC to pick the model and number of clusters. In Figure 5 we can see, model VVE and 3 clusters is the best. The display of clustering is shown in Figure 6. It is not easy to see, so we also displayed the result using ggandrews as in (a) and (b) in Figure .

```

1 # model based clustering
2 library(mclust)
3 mcl <- Mclust(wine.sc)
4 plot(mcl$BIC)
5 plot.Mclust(mcl, what = "classification")
6
7 # display using ggandrews
8 ggandrews(data.frame(mcl$classification, wine.sc), clr = 1, return_value = F)

```

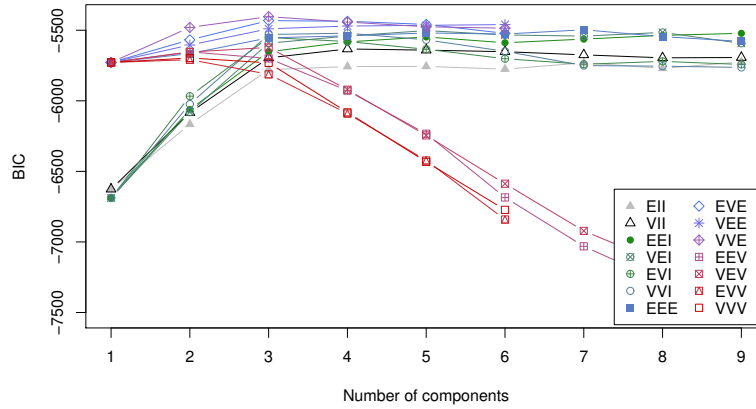


Figure 5: Model selection with BIC

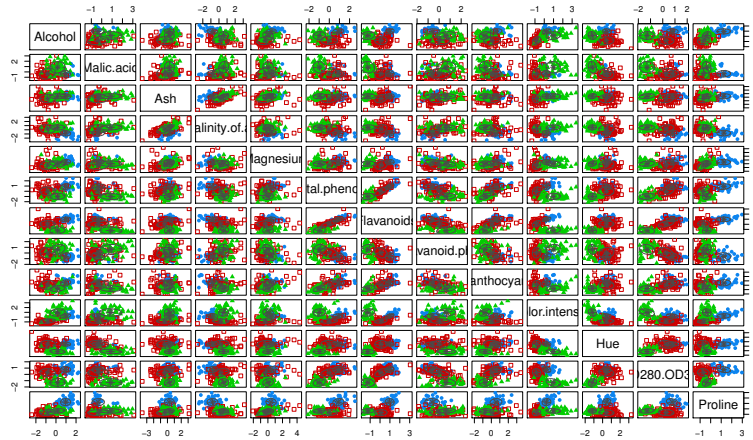


Figure 6: Model based clustering result, $K = 3$

From the display, we can see the result of K-means (both initializations) and model based clustering are similar, and $K = 3$ is quite reasonable from the result. The result of hierarchical clustering is different and it seems like there is only one cluster in the final result.

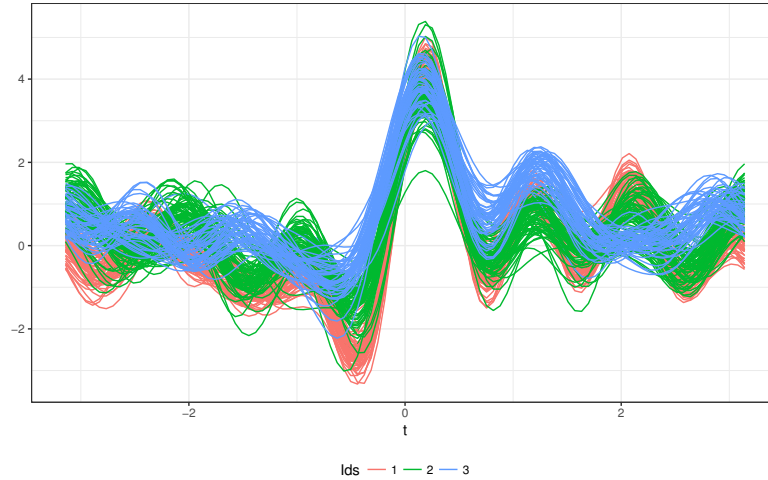


Figure 7: Model based clustering result, $K = 3$

We then compare the clustering result with cultivar information with cross table of frequencies.

```

1  # compare the results with cultivar information
2  # hierarchical clustering with average linkage
3  ftable(table(cultivar, cutree(hc, k = 3)))
4
5      1  2  3
6  cultivar
7  1      58  1  0
8  2      68  2  1
9  3      48  0  0
10
11 # k-means innitialized with hc
12 ftable(table(cultivar, mapvalues(km$cluster, from = c(2,3,1), c(1,2,3))))
13
14      1  2  3
15 cultivar
16 1      59  0  0
17 2       3 65  3
18 3       0  0 48
19
20 # k-means with random initialization
21 ftable(table(cultivar, km.r$cluster))
22
23      1  2  3
24 cultivar
25 1      59  0  0
26 2       3 65  3
27 3       0  0 48
28
29 # model based clustering
30 ftable(table(cultivar, mcl$classification))
31
32      1  2  3
33 cultivar
34 1      56  3  0
35 2       0 70  1
36 3       0  0 48

```

We can see K-means and model based clustering gives the clusters that match the cultivar information very well (in K-means with hierarchical clustering initialization, we renamed the clusters to

match the cultivar information). And the result of hierarchical clustering is not good. It clustered almost all observations to one cultivar.

3. We use the standardized data to perform PCA.

```

1 # PCA
2 wine.pc <- prcomp(wine.sc)
3 # compute proportion of total variance explained by
4 # each component
5
6 s <- wine.pc$sdev^2
7
8 pvar<-s/sum(s)
9 cat("proportion of variance: ", pvar, fill=T)
10
11 # cumulative proportion of total variance explained
12 # by each component
13
14 cpvar <- cumsum(s)/sum(s)
15 cat("cumulative proportion of variance: ", cpvar, fill=T)

```

The proportion is variance and cumulative proportion of variance are:

```

1 proportion of variance:  0.3619885 0.1920749 0.1112363 0.0706903 0.06563294 0.04935823
  ↳ 0.04238679 0.02680749 0.02222153
2 0.01930019 0.01736836 0.01298233 0.007952149
3
4 cumulative proportion of variance: 0.3619885 0.5540634 0.6652997 0.73599 0.8016229
  ↳ 0.8509812 0.893368 0.9201754
5 0.942397 0.9616972 0.9790655 0.9920479 1

```

Based on the cumulative proportion of variance, we then test for 4 and 5 PCs for explaining 80% of total variance.

```

1 # test 5 is enough while 4 is not
2 source("PCs.proportion.variation.enuff.R")
3 PCs.proportion.variation.enuff(lambda = s, q = 4, nobs = nrow(wine.sc), propn = 0.8)
4 PCs.proportion.variation.enuff(lambda = s, q = 5, nobs = nrow(wine.sc), propn = 0.8)

```

The p-values are 0 and 1 respectively. So we conclude that the first 5 PCs are enough to explain 80% of variance in the data.

The first 5 PCs are:

```

1 # first 5 PCs
2 wine.pc$rotation[,1:5]
3
4           PC1          PC2          PC3          PC4          PC5
5 Alcohol    -0.144329395  0.483651548 -0.20738262  0.01785630 -0.26566365
6 Malic.acid  0.245187580  0.224930935  0.08901289 -0.53689028  0.03521363
7 Ash        0.002051061  0.316068814  0.62622390  0.21417556 -0.14302547
8 Alkalinity.of.ash 0.239320405 -0.010590502  0.61208035 -0.06085941  0.06610294
9 Magnesium  -0.141992042  0.299634003  0.13075693  0.35179658  0.72704851
10 Total.phenols -0.394660845  0.065039512  0.14617896 -0.19806835 -0.14931841
11 Flavanoids -0.422934297 -0.003359812  0.15068190 -0.15229479 -0.10902584
12 Nonflavanoid.phenols 0.298533103  0.028779488  0.17036816  0.20330102 -0.50070298
13 Proanthocyanins -0.313429488  0.039301722  0.14945431 -0.39905653  0.13685982
14 Color.intensity 0.088616705  0.529995672 -0.13730621 -0.06592568 -0.07643678
15 Hue        -0.296714564 -0.279235148  0.08522192  0.42777141 -0.17361452
16 OD280.OD315 -0.376167411 -0.164496193  0.16600459 -0.18412074 -0.10116099
17 Proline    -0.286752227  0.364902832 -0.12674592  0.23207086 -0.15786880

```

Each of the 5 PCs seems like contrast of these 13 measures. For example, the first PC is the difference of weighted mean of malic acid, ash, alkalinity of ash, nonflavanoid phenols and color intensity between the weighted mean of alcohol, magnesium, total phenols, flavanoids, proanthocyanins, hue, OD280/OD315 and proline. The rest can also be interpreted in this way.

4. We first split the data into training set and test set.

```

1  # split data
2  set.seed(8413)
3  train.idx <- sample(1:nrow(wine), size = 128, replace = F)
4  wine.train <- wine[train.idx,]
5  wine.test <- wine[-train.idx,]

```

Then for each method, we calculate the AER and LOOCV misclassification rate in training set and also the misclassification rate in test set with the classification rule trained with training set.

- (a) QDA:

```

1  # QDA
2  library(MASS)
3  wine.qda <- qda(Cultivar ~ ., data = wine.train, CV = F)
4
5  # AER
6  mean(wine.train$Cultivar != predict(wine.qda)$class)
7
8  # CV
9  wine.qda.cv <- qda(Cultivar ~ ., data = wine.train, CV = T)
10 mean(wine.train$Cultivar != wine.qda.cv$class)
11
12 # misclassification on test set
13 mean(wine.test$Cultivar != predict(wine.qda, newdata = wine.test[, -1])$class)

```

The result of QDA:

AER = 0
LOOCV = 0.0078125
test = 0

- (b) k-NN:

We first find the optimal k among $\{1, 2, \dots, 10\}$ with leave-one-out cross-validation (function `knn.cv` is used). The cross validation error rate for each k is shown in Figure 8. From Figure 8 we can see $k = 5$ is the optimal one, and we use $k = 5$ to calculate our misclassification rates.

```

1  #kNN
2  library(class)
3  # using cross-validation to pick k
4  # using scaled data
5  wine.train.sc <- scale(wine.train[, -1])
6  wine.test.sc <- scale(wine.test[, -1])
7  # try k = 1, ..., 10
8  knn.cv.err <- NULL
9  knn.cv.sd <- NULL
10 for (i in 1:10) {
11   temp <- NULL
12   for (j in 1:10000)
13     temp <- c(temp, mean(knn.cv(wine.train.sc,
14                               cl = wine.train$Cultivar, k = i) !=
15                               wine.train$Cultivar))
16   knn.cv.err <- c(knn.cv.err, mean(temp))

```

```

16 knn.cv.sd<-c(knn.cv.sd,sd(temp))
17 cat("\n Done i= ",i)
18 }
19
20
21 plot(knn.cv.err, xlim = c(1, 10),
22      ylim=c(min(knn.cv.err - 1.96 * knn.cv.sd),
23             max(knn.cv.err + 1.96 * knn.cv.sd)), type = "n")
24 lines(knn.cv.err + 1.96 * knn.cv.sd, lty = 2, col = "blue")
25 lines(knn.cv.err - 1.96 * knn.cv.sd, lty = 2, col = "green")
26 lines(knn.cv.err, col = "red")
27
28 # use k = 5
29
30 wine.knn.train <- knn(train = wine.train.sc, test = wine.train.sc, cl =
31   ↪ wine.train$Cultivar, k = 5)
32
33 #AER
34 mean(wine.knn.train != wine.train$Cultivar)
35
36 #CV
37 knn.cv.err[5]
38
39 # misclassification on test set
40 wine.knn <- knn(train = wine.train.sc, test = wine.test.sc, cl =
41   ↪ wine.train$Cultivar, k = 5)
42 mean(wine.knn != wine.test$Cultivar)

```

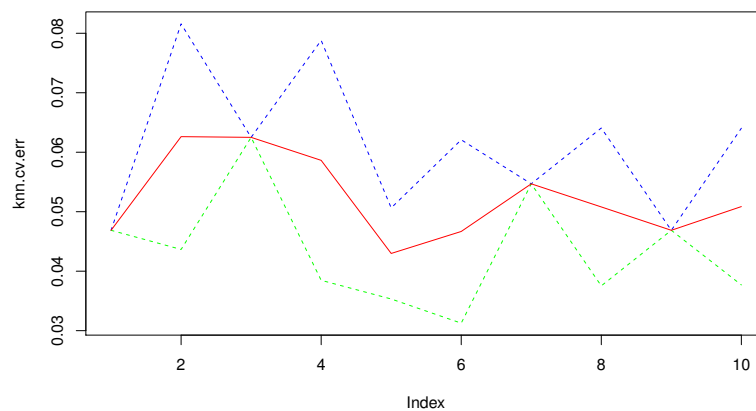


Figure 8: Finding the optimal k for k-NN

The result of k-NN:

$AER = 0.03125$
 $LOOCV = 0.04298281$
 $test = 0.06$

(c) CART:

We first find the optimal tree with leave-one-out cross-validation (function `cv.tree` is used). From Figure 9 we can see when number of nodes is 4, we got the optimal tree (Figure!10). Then we calculated misclassification rates with this tree.


```

1  # CART
2  library(tree)
3  # getting optimal tree using cross-validation
4  wine.tree <- tree(formula = Cultivar ~ ., data = wine.train)
5  wine.tree.cv <- cv.tree(wine.tree, K = nrow(wine.train))
6  plot(wine.tree.cv)
7  # the best one is 4. Plot the best one.
8  wine.tree.opt <- prune.tree(wine.tree, k = 4)
9  plot(wine.tree.opt)
10 text(wine.tree.opt)
11
12 #AER
13 mean(apply(predict(wine.tree.opt), 1, which.max) != wine.train$Cultivar)
14
15 #CV
16 mean(sapply(1:nrow(wine.train), function(x) mean(mean(apply(predict(tree(Cultivar ~
  ↪ ., data = wine.train[-x,]), newdata = wine.train[, -1]), 1,
  ↪ which.max) != wine.train$Cultivar))))
17
18 # misclassification on test set
19 mean(apply(predict(wine.tree.opt, newdata = wine.test[, -1]), 1,
  ↪ which.max) != wine.test$Cultivar)

```

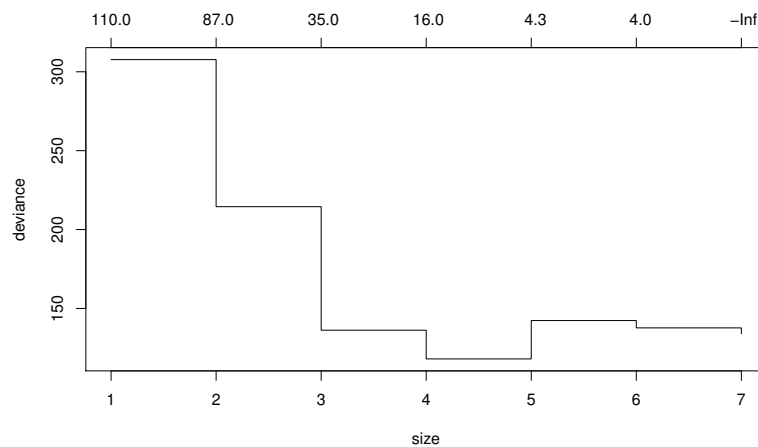


Figure 9: Finding the optimal tree

The result of CART:

AER = 0.03125
 LOOCV = 0.03167725
 test = 0.02

Now we summarise the missclassification rates in the Table 1 below. From the table, we can see that QDA gives the smallest missclassification rates in AER, LOOCV and test. Thus for this data set, we can use QDA as a classification rule as it out-performs others.

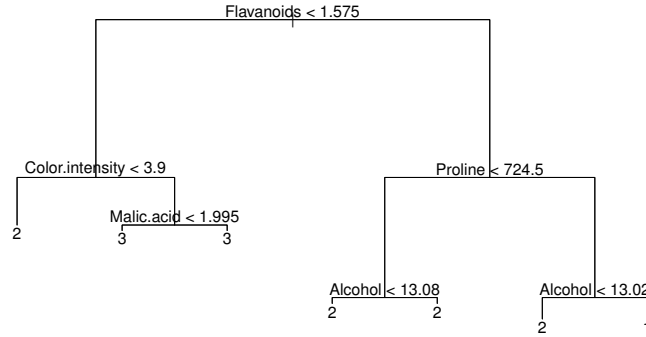


Figure 10: Optimal tree, $k = 4$

Table 1: Summary of missclassification rates

	QDA	k-NN	CART
AER	0	0.03125	0.03125
LOOVA	0.0078125	0.04298281	0.03167725
test	0	0.06	0.02

5. We display the first 5 PCs using 2D and 3D radial visualization in Figure 11 and Figure 12. The cultivars are distinguishable in both 2D and 3D radial visualization, but we can see less overlapping with the 3D radial visualization (a specific angle of view was picked here in Figure 12 to show the non-overlapping of 3 cultivars).

```

1 # display the first 5 PCs using 3d radviz
2 source("radviz3d-3.R")
3 radialvis3d(data = wine.pc$x[,1:5], cl = wine$Cultivar)
4 rgl.snapshot("radviz3d.jpg")
5 # using 2d radviz
6 source("radviz2d.R")
7 source("mmnorm.R")
8 source("circledraw.R")
9 radviz2d(dataset = cbind(wine.pc$x[,1:5], wine$Cultivar), name = "wine")

```

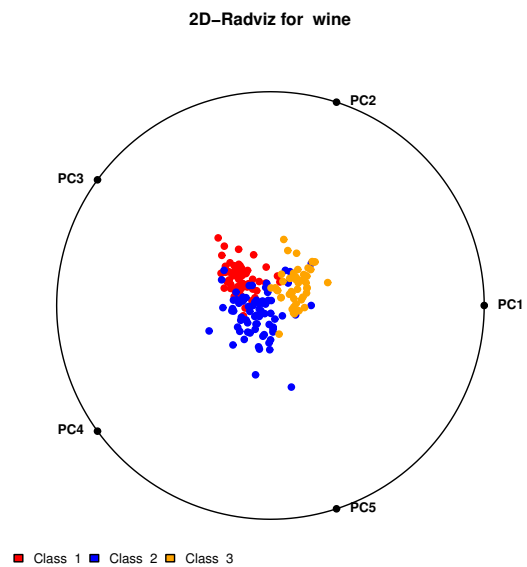


Figure 11: 2D radial visualization

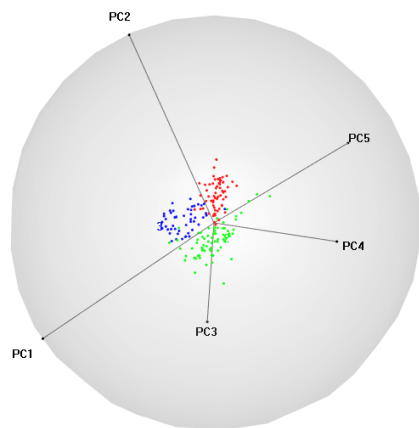


Figure 12: 3D radial visualization