# CS-GY 6923 Machine Learning

*Professor: Dr. Raman Kannan*

# Homework 2:
# Individual Classifier

Hansheng Li

# Contents

# 1. Import Models

Here we import the libraries and data. Then set training and valid set.

```
1. > library(caret)
2. > library(rpart)
3. > library(rpart.plot)
4. > library(nnet)
5. > library(e1071)
6. > library(class)
7. > library(pROC)
8. >
```

```
9. > loadData = function(csvfile) { read.csv(csvfile,head=T,sep=',',stringsAsFa
   ctors=F)} # Function to load the data
10. > data = loadData('Dry_Bean_Dataset.csv') # Load the data
11. >
12. > head(data)[,1:6]
13.    Area Perimeter MajorAxisLength MinorAxisLength AspectRation Eccentricity
14. 1 28395   610.291        208.1781        173.8887     1.197191    0.5498122
15. 2 28734   638.018        200.5248        182.7344     1.097356    0.4117853
16. 3 29380   624.110        212.8261        175.9311     1.209713    0.5627273
17. 4 30008   645.884        210.5580        182.5165     1.153638    0.4986160
18. 5 30140   620.134        201.8479        190.2793     1.060798    0.3336797
19. 6 30279   634.927        212.5606        181.5102     1.171067    0.5204007
20. >
21. > dim(data) # Check the number of rows and columns of the data
22. [1] 13611     17
23. > names(data) # Check the names of each column (features and class variable)
24.  [1] "Area"           "Perimeter"       "MajorAxisLength" "MinorAxisLength"
25.  [5] "AspectRation"   "Eccentricity"    "ConvexArea"      "EquivDiameter"
26.  [9] "Extent"         "Solidity"        "roundness"       "Compactness"
27. [13] "ShapeFactor1"   "ShapeFactor2"    "ShapeFactor3"    "ShapeFactor4"
28. [17] "Class"
29. > length(names(data)) # Double check that the length of "names" matches the
   number of columns
30. [1] 17
31. > which(names(data)=='Class') # Position of the class variable column
32. [1] 17
33. > table(data$Class)
34.
35. BARBUNYA    BOMBAY    CALI DERMASON    HOROZ    SEKER    SIRA
36.     1322       522    1630     3546     1928     2027    2636
37. >
38. > data$Class <- as.factor(data$Class)
39. > set.seed(18)
40. > train.idx <- sample(1:nrow(data), 0.7*nrow(data),replace = F)
41. > data.train <- data[train.idx,]
42. > data.valid <- data[-train.idx,]
```

# 2. Multinomial Logistic Regression

Build Multinomial Logistic Regression Model

```
1. > log.fit <- multinom(Class ~ ., data = data.train)
2. # weights:  126 (102 variable)
```

```
3.  initial  value 18538.685990
4.  iter  10 value 13417.578057
5.  iter  20 value 9406.201803
6.  iter  30 value 7131.638986
7.  iter  40 value 3444.545535
8.  iter  50 value 2079.206940
9.  iter  60 value 1982.080845
10. iter  70 value 1951.091573
11. iter  80 value 1925.082531
12. iter  90 value 1917.739223
13. iter 100 value 1909.562134
14. final  value 1909.562134
15. stopped after 100 iterations
```

# 3. Support Vector Machine Classification

Build Support Vector Machine Model

```
1.  > svm.fit <- svm(Class ~ ., data = data.train)
2.  > svm.fit
3.
4.  Call:
5.  svm(formula = Class ~ ., data = data.train)
6.
7.
8.  Parameters:
9.     SVM-Type:  C-classification
10.  SVM-Kernel:  radial
11.        cost:  1
12.
13. Number of Support Vectors:  2064
```
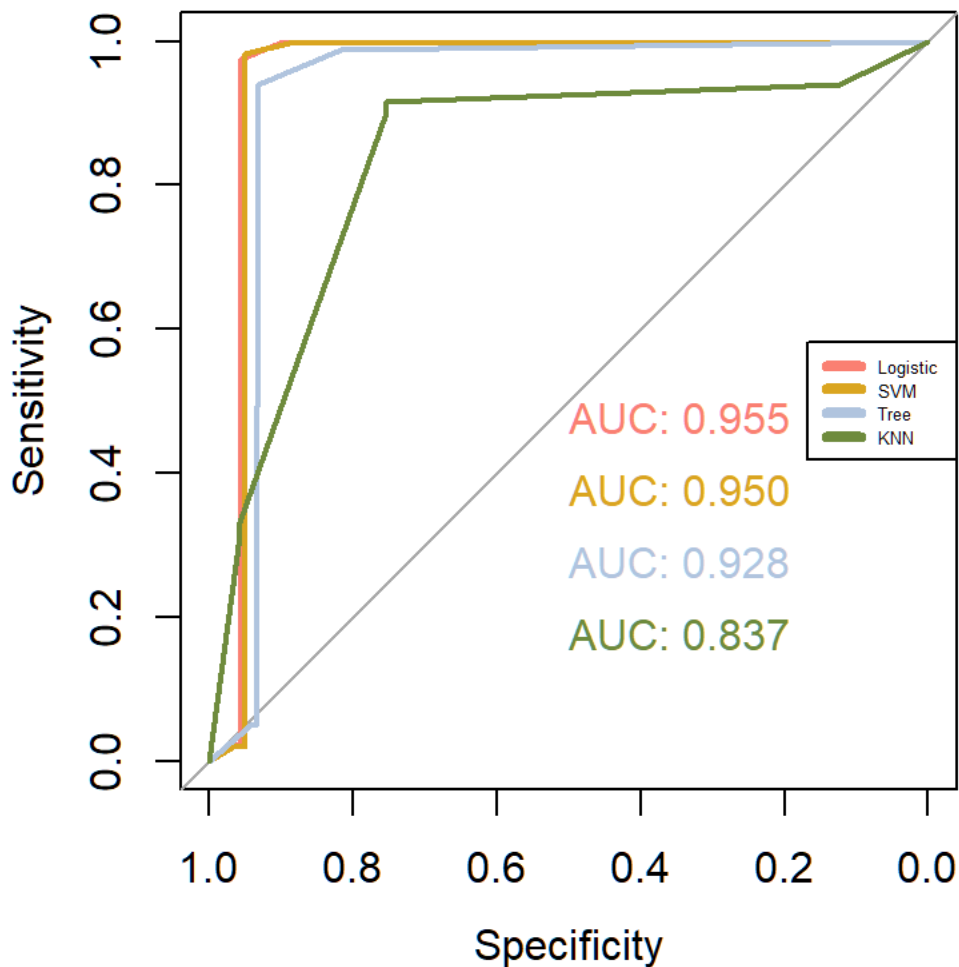
# 4. Decision Tree Classification

Build Decision Tree Model

```
1. > tree.fit <- rpart(Class ~ ., data = data.train)
2. > prp(tree.fit, extra=104)
```

```
                                    MajorAxi >= 281
              yes                                              no


              Compactn >= 0.73                        ShapeFac >= 0.0068

                                    HOROZ          DERMASON         SEKER
                                .01 .01 .05 .00   .00 .00 .00 .88  .00 .00 .00 .06
        Perimete >= 896         .94 .00 .00       .00 .02 .09      .00 .87 .07
                                    14%               27%              15%


                              SIRA
                          .02 .00 .00 .06
     ShapeFac < 0.0041     .04 .06 .83
                              19%


       BOMBAY
   .00 1.00 .00 .00        Compactn >= 0.78
   .00 .00 .00
       4%

              BARBUNYA
          .86 .00 .13 .00
          .00 .01 .00         roundnes < 0.8
              8%


                    BARBUNYA          CALI
                .82 .00 .09 .00   .08 .00 .86 .00
                .09 .00 .01       .04 .00 .02
                    2%               12%
```

# 5. K-Nearest Neighbor Classification

Build K-Nearest Neighbor Model. Here we set the k to square root of n.

```
1. > nab = as.integer(sqrt(dim(data))[1])
2. > knn.pred <- knn(data.train[,-c(17)],data.valid[,-
              c(17)],cl=data.train$Class,k=nab)
3. > knn.pred
4.    [1] DERMASON DERMASON DERMASON DERMASON DERMASON DERMASON DERMASON DERMAS
              ON DERMASON
5.   [10] DERMASON DERMASON DERMASON DERMASON DERMASON DERMASON DERMASON DERMAS
              ON DERMASON
6.   [19] DERMASON DERMASON DERMASON DERMASON DERMASON DERMASON DERMASON DERMAS
              ON DERMASON
7.   ……
```

# 6. ROC and AUC Performance

Build ROC and AUC diagram.

```
1.  > data.valid$Class <- factor(data.valid$Class, ordered = T)
2.  > log.pred <- factor(predict(log.fit, data.valid), ordered = T)
3.  > log.roc <- multiclass.roc(data.valid$Class, log.pred,quiet=T)
4.  > log.rs <- log.roc[['rocs']]
5.  > svm.pred <- factor(predict(svm.fit, data.valid), ordered = T)
6.  > svm.roc <- multiclass.roc(data.valid$Class, svm.pred)
7.  > svm.rs <- svm.roc[['rocs']]
8.  > tree.pred <-
                factor(predict(tree.fit, data.valid,type = "class"), ordered
             = T)
9.  > tree.roc <- multiclass.roc(data.valid$Class, tree.pred)
10. There were 21 warnings (use warnings() to see them)
11. > tree.rs <- tree.roc[['rocs']]
12. > knn.pred <- factor(knn.pred, ordered = T)
13. > knn.roc <- multiclass.roc(data.valid$Class, knn.pred)
14. > knn.rs <- knn.roc[['rocs']]
15. > par(pty = "s")
16. > sec = 4
17. > plot.roc(log.rs[[sec]], col = "salmon", percent=TRUE,print.auc=T)
18. > plot.roc(svm.rs[[sec]], col = "goldenrod", print.auc=T, percent=TRUE,add
                = T, print.auc.y=.4)
19. > plot.roc(tree.rs[[sec]], col = "lightsteelblue", print.auc=T, percent=TR
                UE,add = T, print.auc.y=.3)
20. > plot.roc(knn.rs[[sec]], col = "darkolivegreen4", print.auc=T, percent=TR
                UE,add = T, print.auc.y=.2)
21. > legend("right",
22. +         legend=c("Logistic", "SVM", "Tree","KNN"),
23. +         col=c("salmon", "goldenrod", "lightsteelblue","darkolivegreen4"),
24. +         lwd=4, cex =0.4, xpd = TRUE, horiz = F)
```

Here we can see Logistic have highest AUC, SVM and decision tree were pretty close. But the AUC for KNN is not that good compared to the others.

Logistic and SVM both preformed pretty good in ROC followed by decision tree. KNN did better job at start but growth rate become slowing down after.

# 7. AUC Confusion Matrix

```
1.   > (log.cM <- table(log.pred, data.valid$Class))
2.
3.   log.pred    BARBUNYA BOMBAY CALI DERMASON HOROZ SEKER SIRA
4.      BARBUNYA      351      0    9        0     1     3    2
5.      BOMBAY          0    153    0        0     0     0    0
```

```
6.    CALI            21      0  456       0   13    0    0
7.    DERMASON         0      0    0    1019    3   11   64
8.    HOROZ            0      0    9       3  525    0   14
9.    SEKER            6      0    2      14    0  561    6
10.   SIRA           11      0    4      83   12   24  704
11. There were 50 or more warnings (use warnings() to see the first 50)
12. >
13. > (svm.cM <- table(svm.pred, data.valid$Class))
14.
15. svm.pred    BARBUNYA BOMBAY CALI DERMASON HOROZ SEKER SIRA
16.    BARBUNYA      346      0    6       0    1    1    3
17.    BOMBAY          0    153    0       0    0    0    0
18.    CALI           24      0  458       0    9    0    0
19.    DERMASON        0      0    0    1028    3   13   69
20.    HOROZ           0      0    9       3  530    0   12
21.    SEKER           6      0    2      14    0  558    5
22.    SIRA           13      0    5      74   11   27  701
23. >
24. > (tree.cM <- table(tree.pred, data.valid$Class))
25.
26. tree.pred   BARBUNYA BOMBAY CALI DERMASON HOROZ SEKER SIRA
27.    BARBUNYA      318      0   33       0    6    1    2
28.    BOMBAY          0    153    0       0    0    0    0
29.    CALI           45      0  400       0   27    0    3
30.    DERMASON        0      0    0    1032    2   36   98
31.    HOROZ           1      0   42       0  492    0    4
32.    SEKER           3      0    0      30    0  512   41
33.    SIRA           22      0    5      57   27   50  642
34. >
35. > (knn.cM <- table(knn.pred, data.valid$Class))
36.
37. knn.pred    BARBUNYA BOMBAY CALI DERMASON HOROZ SEKER SIRA
38.    BARBUNYA       48      0   41       0   33    0    0
39.    BOMBAY          0    153    0       0    0    0    0
40.    CALI          245      0  397       0   14    0    0
41.    DERMASON        0      0    0     952    8  188   44
42.    HOROZ          80      0   39       0  315   14   66
43.    SEKER           0      0    0     121   11  150   73
44.    SIRA           16      0    3      46  173  247  607
```

# 8. Accuracy

Same as before, Logistic get best accuracy. KNN did not so well.

```
1. > log.acc <- sum(diag(log.cM)) / sum(log.cM)
2. > svm.acc <- sum(diag(svm.cM)) / sum(svm.cM)
```

```
3. > tree.acc <- sum(diag(tree.cM)) / sum(tree.cM)
4. > knn.acc <- sum(diag(knn.cM)) / sum(knn.cM)
5. > acc <- data.frame(Models = c("Logistic", "SVM", "Tree","KNN"),
6. +                   Accuracy = c(log.acc, svm.acc, tree.acc, knn.acc))
7. > acc
8.    Models  Accuracy
9.         1 Logistic 0.9228697
10.        2      SVM 0.9240940
11.        3     Tree 0.8690010
12.        4      KNN 0.6420176
```

# 9. Specificity

```
1.  > log.conf <- confusionMatrix(log.pred, data.valid$Class)
2.  > svm.conf <- confusionMatrix(svm.pred, data.valid$Class)
3.  > tree.conf <- confusionMatrix(tree.pred, data.valid$Class)
4.  > knn.conf <- confusionMatrix(knn.pred, data.valid$Class)
5.  >
6.  > specificity <- data.frame(cbind(log.conf$byClass[,"Specificity"],
7.  +                                 svm.conf$byClass[,"Specificity"],
8.  +                                 tree.conf$byClass[,"Specificity"],
9.  +                                 svm.conf$byClass[,"Specificity"],
10. +                                 knn.conf$byClass[,"Specificity"]))
11. > names(specificity) <- c("Logistic", "SVM", "Tree","KNN")
12. > specificity
13.                  Logistic       SVM      Tree       KNN        NA
14. Class: BARBUNYA 0.9959405 0.9970230 0.9886333 0.9970230 0.9799729
15. Class: BOMBAY   1.0000000 1.0000000 1.0000000 1.0000000 1.0000000
16. Class: CALI     0.9905660 0.9908435 0.9791898 0.9908435 0.9281354
17. Class: DERMASON 0.9736931 0.9713322 0.9541315 0.9713322 0.9190556
18. Class: HOROZ    0.9926346 0.9932011 0.9866856 0.9932011 0.9436261
19. Class: SEKER    0.9919656 0.9922525 0.9787661 0.9922525 0.9411765
20. Class: SIRA     0.9593200 0.9605343 0.9511233 0.9605343 0.8527626
```

# 10. Precision

```
1.  > log.precision<- diag(log.cM) / rowSums(log.cM)
2.  > svm.precision<- diag(svm.cM) / rowSums(svm.cM)
```

```
3.  > tree.precision<- diag(tree.cM) / rowSums(tree.cM)
4.  > knn.precision<- diag(knn.cM) / rowSums(knn.cM)
5.  > precision <-
        data.frame(cbind(log.precision, svm.precision, tree.precision
        , knn.precision))
6.  > precision
7.          log.precision svm.precision tree.precision knn.precision
8.  BARBUNYA     0.9590164      0.9691877      0.8833333      0.3934426
9.  BOMBAY       1.0000000      1.0000000      1.0000000      1.0000000
10. CALI         0.9306122      0.9327902      0.8421053      0.6051829
11. DERMASON     0.9288970      0.9236298      0.8835616      0.7986577
12. HOROZ        0.9528131      0.9566787      0.9128015      0.6128405
13. SEKER        0.9524618      0.9538462      0.8737201      0.4225352
14. SIRA         0.8400955      0.8435620      0.7995019      0.5558608
```

# 11.  Sensitivity/Recall

```
1.  > log.recall <- diag(log.cM) / colSums(log.cM)
2.  > svm.recall <- diag(svm.cM) / colSums(svm.cM)
3.  > tree.recall <- diag(tree.cM) / colSums(tree.cM)
4.  > knn.recall <- diag(knn.cM) / colSums(knn.cM)
5.  >
6.  > recall <-
        data.frame(cbind(log.recall, svm.recall, tree.recall, knn.rec
        all))
7.  > recall
8.          log.recall svm.recall tree.recall knn.recall
9.  BARBUNYA  0.9023136  0.8894602   0.8174807  0.1233933
10. BOMBAY    1.0000000  1.0000000   1.0000000  1.0000000
11. CALI      0.9500000  0.9541667   0.8333333  0.8270833
12. DERMASON  0.9106345  0.9186774   0.9222520  0.8507596
13. HOROZ     0.9476534  0.9566787   0.8880866  0.5685921
14. SEKER     0.9365609  0.9315526   0.8547579  0.2504174
15. SIRA      0.8911392  0.8873418   0.8126582  0.7683544
```

## 12. Bias

```
1. > log.bias <- mean(predict(log.fit, data.train) != data.train$Class)
2. > svm.bias <- mean(predict(svm.fit, data.train) != data.train$Class)
3. > tree.bias <-
            mean(predict(tree.fit, data.train, type = "class") != data.tr
         ain$Class)
4. > bias <- data.frame(cbind(log.bias, svm.bias, tree.bias))
5. > bias
6.    log.bias   svm.bias tree.bias
7. 1 0.07232077 0.06623281 0.1241734
```

## 13. Variance

```
1. > log.var <- mean(log.pred != data.valid$Class)
2. > svm.var <- mean(svm.pred != data.valid$Class)
3. > tree.var <- mean(as.character(tree.pred) != data.valid$Class)
4. > knn.var <- mean(knn.pred != data.valid$Class)
5. > variance <- data.frame(cbind(log.var, svm.var, tree.var, knn.var))
6. > variance
7.    log.var    svm.var tree.var   knn.var
8. 1 0.07713026 0.07590597 0.130999 0.3579824
```

## 14. Compare all Model

In conclusion, logical regression and SVM preform similarly good in ROC and AUC test, decision tree is ok and KNN is not so well, so we optimal to select logical regression and SVM first. In Confusion Matrix, Accuracy, Specificity, Precision and Recall ( Sensitivity ), they both preform very similarly, but in Bias and Variance, we can see SVM is a little bit better, and SVM is better to prevent overfitting than logical regression, also we have small number of feature but large number of samples, in this case, I would recommend SVM over logical regression. Decision tree and KNN may be better if we use

ensemble techniques or other kernel to gives more optimize, but still hard to accuracy beyond SVM.