

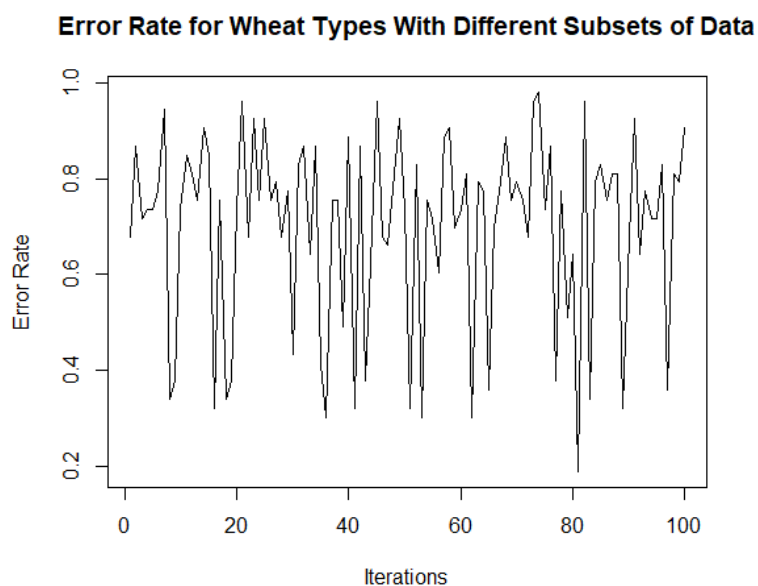
Please send your source codes and results (figures, replies etc.) in one text file (e.g., YourName.docx / YourName.pdf), then submit it to the Classroom.

Group homework is not allowed!

1. You will code a **decision tree** (DT) that will be similar with the one given in the lab examples. The input data is provided in the text file “wheat_types.txt”. The target feature is the “type” column. You will train and test the DT by using all features except “type”. Apply the following steps and reply the questions in your report.
 - You can use a different decision tree library apart from “tree”.
 - Use 75% of samples for training and 25% of them for test purposes.
 - Perform this sample division (cross-validation) process for “100” iterations.
 - Draw a plot to show how accuracy changes over 100 iterations.
 - Report average accuracy of DT after 100 iterations.
 - Show the DT (with parent and internal node’s decision criteria) which has the highest accuracy over 100 iterations.

ANSWERS FOR 1st –(The R Code is the bottom of the page)

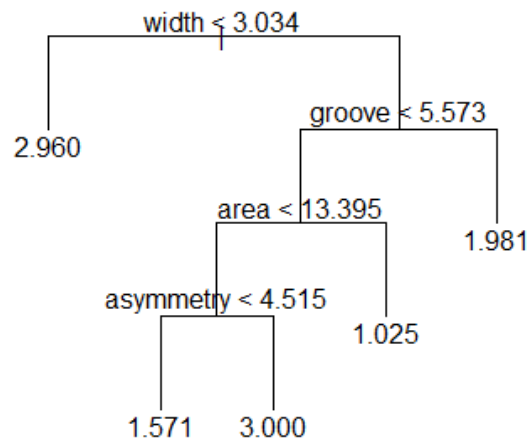
(Draw a plot to show how accuracy changes over 100 iterations.)



(Report average accuracy of DT after 100 iterations.)

```
> #Report average accuracy of DT after 100 iterations.  
> mean(dt_acc)  
[1] 0.2998113
```

(Show the DT which has the highest accuracy over 100 iterations.)



2. You will code a **k-nearest neighbor** (kNN) algorithm, which will be similar with the one given in the lab examples. You will use built-in data set “Smarket” that is provided in the “ISLR” library. The target feature is the “Direction” column. You will use various combinations of descriptive features and different *k* values. Apply the following steps and reply the questions in your report.

- You need to test several combinations of descriptive features to improve the prediction accuracy of the target feature. The combinations:
 - Lag1, Lag2, Lag3
 - Lag3, Lag4, Volume
 - Lag2, Lag5, Volume
- Use 60% of samples for training and 40% of them for test purposes (for each combination).
- Test different *k* values (1 to 30) for each combination.
- Make a plot to show how accuracy changes for different *k* values (for each combination).

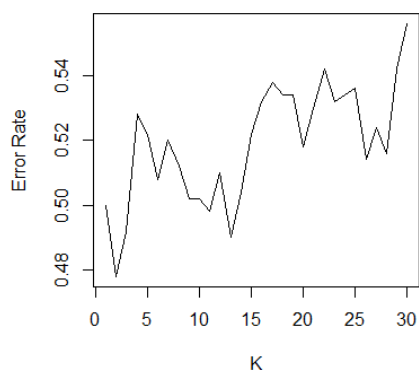
- Which k value did provide the highest accuracy for each combination?
Report both k values and the highest accuracy value.
- What is the highest accuracy to predict the “Direction” column when you set a specific k value and the descriptive feature combination?

ANSWERS FOR 2nd –(The R Code is the bottom of the page)

(plot to show how accuracy changes for different k values) (for each combination)

(Lag1, Lag2, Lag3)

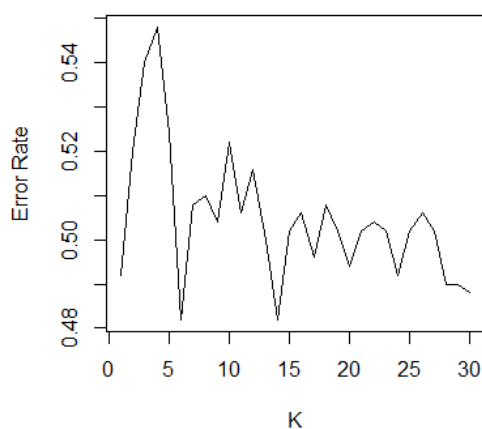
Error Rate for Smarket with varying K



```
> #highest accuracy
> highest_accuracy1
[1] 0.522
>
> #k value has highest accuracy
> highest_k_accuracy1
[1] 2
```

(Lag3, Lag4, Volume)

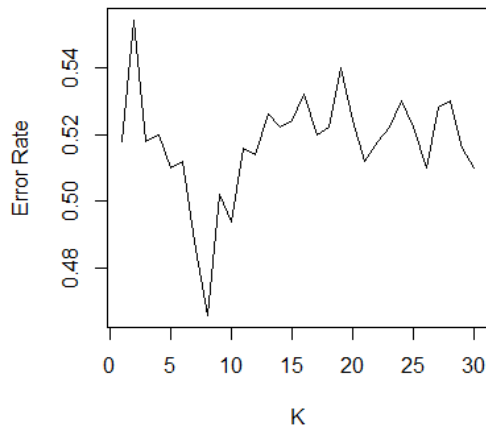
Error Rate for Smarket with varying K



```
> #highest accuracy
> highest_accuracy2
[1] 0.518
>
> #k value has highest accuracy
> highest_k_accuracy2
[1] 6
```

(Lag2, Lag5, Volume)

Error Rate for Smarket with varying K



```
> #highest accuracy
> highest_accuracy3
[1] 0.534
>
> #k value has highest accuracy
> highest_k_accuracy3
[1] 8
```

(The highest accuracy to predict the “Direction” column when set a specific k value and the descriptive feature combination)

```
+ j
The highest accuracy(0.534000) provided by the combinations that are Lag2, Lag5 and Volume
> |
```

R Codes For 1st Question

```
1. setwd('C:\\Users\\akifc\\Desktop')
2. options(max.print=999999)
3.
4. library(tree) # Contains the "tree" function
5.
6. wheat_types <- read.table("wheat_types.txt", header = TRUE, sep = ";", dec = ".")
7.
8.
9. #Set the seed for reproducibility
10. # Use 75% of samples for training, the rest for testing
11. # the indices (row ids) are saved in the "sub" vector
12. set.seed(579642)
13. sub <- sample(1:nrow(wheat_types), size=nrow(wheat_types)*0.75)
14.
15. wt.tr <- tree(type ~ ., data = wheat_types, subset = sub)
```

```

16. summary(wt.tr)
17.
18. #####
19. #Cross-validation version - Construct a new DT
20. #for different partitions of the samples - 100 times
21.
22. dt_acc <- numeric()
23. set.seed(1815850)
24.
25. highest_accuracy <- -99
26.
27. for(i in 1:100){
28.   sub <- sample(1:nrow(wheat_types), size=nrow(wheat_types)*0.75)
29.   fit2 <- tree(type ~ ., data = wheat_types, subset = sub)
30.   test_predict <- table(predict(fit2, wheat_types[-sub, ]), wheat_types[-
     sub, "type"])
31.   dt_acc <- c(dt_acc, sum(diag(test_predict)) / sum(test_predict))
32.
33.   if(highest_accuracy < dt_acc[i]){ #DT which has the highest accuracy over 100 iter
     ations.
34.     highest_accuracy <- dt_acc[i]
35.     highest_acc_tr <- fit2
36.   }
37. }
38.
39. #Report average accuracy of DT after 100 iterations.
40. mean(dt_acc)
41.
42.
43. #Draw a plot to show how accuracy changes over 100 iterations.
44. plot(1-
     dt_acc, type="l", ylab="Error Rate", xlab="Iterations", main="Error Rate for Wheat T
     ypes With Different Subsets of Data")
45.
46. #####
47.
48. #The value 'highest_acc_tr' filled above.
49. # Show the DT (with parent and internal node's decision criteria) which has the high
     est accuracy over 100 iterations.
50. plot(highest_acc_tr, type = "uniform")
51. text(highest_acc_tr)

```

R Codes For 2nd Question

```

1. install.packages("ISLR")
2. library(ISLR)
3.
4. install.packages('class')
5. library(class) # Contains the "knn" function
6.
7. #####
8. #####
9.
10. #Create partitions in the Smarket data set (60% for training, 40% for testing/evalua
     tion)
11. sm_sample <- sample(1:nrow(Smarket), size=nrow(Smarket)*0.6)
12. sm_train <- Smarket[sm_sample, ] #Select the 60% of rows
13. sm_test <- Smarket[-sm_sample, ] #Select the 40% of rows
14.

```

```

15. #####
16. #####
17.
18. #Here 2,3,4 in the vector represent Lag1, Lag2 and Lag3 features
19. train.X <- sm_train[,c(2,3,4)]
20. test.X <- sm_test[,c(2,3,4)]
21.
22. #Seed must set in order to get reproducible result
23. set.seed(4985912356) #Set the seed for reproducibility
24.
25. #First try to determine the right K-value
26. smarket_acc <- numeric() #holding variable
27.
28.
29. highest_accuracy1 <- -99
30. highest_k_accuracy1 <- -1
31.
32.
33. for(i in 1:30){
34.   predict <- knn(train=train.X, test=test.X, cl=sm_train$Direction, k=i)
35.   smarket_acc <- c(smarket_acc, mean(predict==sm_test$Direction))
36.   if(highest_accuracy1 < smarket_acc[i]){
37.     #To find the value which has the highest accuracy.
38.     highest_accuracy1 <- smarket_acc[i]
39.     highest_k_accuracy1 <- i
40.   }
41. }
42.
43. #Plot error rates for k=1 to 30
44. plot(1-
      smarket_acc, type="l", ylab="Error Rate", xlab="K", main="Error Rate for Smarket with varying K")
45.
46. #highest accuracy
47. highest_accuracy1
48.
49. #k value has highest accuracy
50. highest_k_accuracy1
51.
52. #####
53. #####
54.
55. #Here 4,5,7 in the vector represent Lag3, Lag4 and Volume features
56. train.X <- sm_train[,c(4,5,7)]
57. test.X <- sm_test[,c(4,5,7)]
58.
59. #Seed must set in order to get reproducible result
60.
61. #First try to determine the right K-value
62. smarket_acc <- numeric() #holding variable
63.
64.
65. highest_accuracy2 <- -99
66. highest_k_accuracy2 <- -1
67.
68.
69. for(i in 1:30){
70.   predict <- knn(train=train.X, test=test.X, cl=sm_train$Direction, k=i)
71.   smarket_acc <- c(smarket_acc, mean(predict==sm_test$Direction))
72.   if(highest_accuracy2 < smarket_acc[i]){
73.     #To find the value which has the highest accuracy.
74.     highest_accuracy2 <- smarket_acc[i]
75.     highest_k_accuracy2 <- i
76.   }
77. }
78.

```

```

79. #Plot error rates for k=1 to 30
80. plot(1-
      smarket_acc, type="l", ylab="Error Rate", xlab="K", main="Error Rate for Smarket wi
th varying K")
81.
82. #highest accuracy
83. highest_accuracy2
84.
85. #k value has highest accuracy
86. highest_k_accuracy2
87.
88. #####
89. #####
90.
91. #Here 3,6,7 in the vector represent Lag2, Lag5 and Volume features
92. train.X <- sm_train[,c(3,6,7)]
93. test.X <- sm_test[,c(3,6,7)]
94.
95. #Seed must set in order to get reproducible result
96.
97. #First try to determine the right K-value
98. smarket_acc <- numeric() #holding variable
99.
100.
101.     highest_accuracy3 <- -99
102.     highest_k_accuracy3 <- -1
103.
104.
105.     for(i in 1:30){
106.         predict <- knn(train=train.X, test=test.X, cl=sm_train$Direction, k=i)
107.         smarket_acc <- c(smarket_acc, mean(predict==sm_test$Direction))
108.         if(highest_accuracy3 < smarket_acc[i]){
109.             #To find the value which has the highest accuracy.
110.             highest_accuracy3 <- smarket_acc[i]
111.             highest_k_accuracy3 <- i
112.         }
113.     }
114.
115.     #Plot error rates for k=1 to 30
116.     plot(1-
          smarket_acc, type="l", ylab="Error Rate", xlab="K", main="Error Rate for Smarket wi
th varying K")
117.
118.     #highest accuracy
119.     highest_accuracy3
120.
121.     #k value has highest accuracy
122.     highest_k_accuracy3
123.
124.     #####
125.     max_value <- max(highest_k_accuracy1,highest_accuracy2,highest_k_accuracy3)
126.
127.     if(highest_k_accuracy1==max_value){
128.         cat(sprintf("The highest accuracy(%) provided by the combinations that are
Lag1, Lag2 and Lag3",highest_accuracy1))
129.     } else if(highest_k_accuracy2==max_value){
130.         cat(sprintf("The highest accuracy(%) provided by the combinations that are
Lag3, Lag4 and Volume",highest_accuracy2))
131.     } else {
132.         cat(sprintf("The highest accuracy(%) provided by the combinations that are
Lag2, Lag5 and Volume",highest_accuracy3))
133.     }

```
