Wesley Pryor

# Wine Classification

**Final Project**

**M5364–Data Mining   Fall 2016   Tarleton State Univ   Dr. Scott Cook   Assigned 2016-11-02   Due 2016-12-09**

## 1  Summary of Data and Models

The Wine Quality Data set is a set of data that contains observations of wines. This data comes form the UCI Machine Learning Repository. The data is divided into two sets. The first data set contains 4898 observations of white wines. The second data set has 1599 different red wines. The data contains 12 attributes. The data has 11 numeric attributes. The numeric attributes are fixed acidity, volatile acidity, citric acid, residual sugar, chlorides, free sulfur dioxide, total sulfur dioxide, density, pH, sulphates, and alcohol. The twelfth attribute is quality. This part of the data is usually what is trying to be calculated with a regression problem, or it is also determined via classification techniques. I was not sure how this attribute could be treated. I initially thought about treating it just as a normal numeric value. However, I wanted to see how to treat it as a factor.

I decided to tackle this set of data by combining the red and white rows of data. The combines red and white data had 6,497 different type of wines from Portugal. I wanted to classify the wines between red and white. For this project, I applied classification techniques that we have covered in class throughout the semester to determine the type of wine. The table below describes the techniques applied and the values of accuracy for the respective models. All original models were run with a simple 70-30 split and then cross-validation was applied to the models to determine the mean accuracy for each technique.

Table 1: Accuracy of Models Using Cross-Validation

| Model | Accuracy |
|---|---|
| R Part Decision Trees | 0.9804547 |
| C Tree Decision Trees | 0.9761422 |
| K Nearest Neighbors | 0.9938433 |
| Weighted K nearest Neighbors | 0.9943058 |
| Naive Bayes | 0.9752199 |
| Support Vector Machines | 0.9955366 |
| Artificial Neural Networks | 0.9723155 |

## 2  R Part Decision Trees

When trying to classify the model. I initially ran the decision tree with the `rpart` package. This decision tree started with a depth of 4 and and accuracy of 97.6%. Most of the data in the naive run of the tree was split from the counts of total sulfur dioxides and chlorides. The other factor involved was volatile acidity. Once this model was run, tuning was applied to the complexity parameter, minimum split, and the maximum depth of the trees.

For the complexity parameters, I decided that the best was to test small values from 0.0 to 0.2. After tuning was run, 0.001 was determined to be the most accurate parameter for the model. From our in class sources, I initially tried to test a minimum split of 5, 10, or 15. However, I was getting the same error rate for all of those values. After looking through the data more, I realized that trying to split with only 5, 10,, or 15 observations would be over-fitting the model and not very ideal. I then tried a splits of 50 125 by 5 and i saw that the errors began to change after 50 splits. 50 was determined to be the best minimum split. The maximum depth was decided from values between 1 and 10. The best tree had a depth of 6. Once complexity parameter, minimum split, and maximum depth were set to the optimal parameters, the accuracy of the rerun decision tree was 98%. The best decision tree included fixed acidity in splitting the data into the type of wines, along with the other attributes that were in the naive tree. The optimal decision tree was run using 10 fold cross-validation and the mean of those accuracies was 98%.

## 28 C Tree Decision Trees

29 In class this semester, our assignments led us to believe that the best decision tree was used with the `rpart` calculations instead
30 of the `ctree` calculations. I wanted to test the method and determine which method was better for this data. The first attempt
31 of the model split most frequently with density, alcohol, total sulfur dioxides, sulphates and chlorides. Other attributes used in
32 the model were pH, residual sugar, volatile acidity, fixed acidity, and chlorides. With these factors, the model had 69 nodes and
33 the model has a depth of 8. The initial accuracy was 98.1% The model was tuned for the minimum split and the maximum
34 depth. It was tested for values of minimum split at 50, 100, and 150. The optimal value was 50. At maximum depth, values
35 between 1 and 10 were tested. The best depth was calculated to be 10. The optimal tree was tested with these values and had
36 an accuracy of 98.3%. The tree looked similar to the first run of a decision tree. However, the overall depth was increased by
37 one. Once cross-validation was done with this model, the average accuracy was much lower at 97.6%. I was suprised to see this.
38 However, I believe that this was done from splits that had more data that was had similar values in certain categories. Overall
39 after cross-validation of both methods, the `rpart` calculations were proven to be better for this data.

## 40 K nearest neighbors

41 K nearest neighbors was standardized initially and was calculated with 3 nearest neighbors. This trial resulted in a 99.2%
42 accuracy. The trial was then tuned with nearest neighbors of 2 through 20. The method was also tuned with three different
43 types of sampling: cross-validation, fixed sampling, and bootstrapping. The best nearest neighbors for cross-validation was 5.
44 Fixed sampling had an optimal nearest neighbor count of 11. Bootstrapping was 12. I tested the model with all the values
45 of k that were determined in the sampling and then cross-validated. Through tuning and cross-validation, the overall average
46 accuracy was 99.3%

## 47 Weighted K Nearest Neighbors

48 I wanted to test to see how the kernels of K nearest neighbors affected the values of the model. The first trial was with a
49 triangular kernel and the accuracy was 99.2%. I tested all eight kernels available. I expected that the optimal kernel would
50 perform better than my first trial. After testing, the inverted kernel was determined to have minimal errors with 3 nearest
51 neighbors. The accuracy was run once for a 99.2% accuracy. cross-validation resulted in a mean 99.4% accuracy.

## 52 Naive Bayes

53 Naive Bayes has the assumption of normality. I had to first test the normality of all attributes of the data. The numerical
54 values of the data were not normal. At this point I tried to see if I could treat the quality as a factor instead of a numerical
55 value. All models had obvious curvature in all of the Q-Q Plots. Also, all Shapiro-Wilk tests returned the smallest possible
56 p-values for a normal distribution. The data had to be cut to assume normality with the `mycut` function that was developed in
57 this semester. The attributes that needed normalizing were cut into 10 factor levels. Due to time constraints, I was not able to
58 test multiple cuts. I wish that I had tested multiple cuts to find the best cut for the model. Due to the size of the data, I tried
59 10 and felt like this was a good estimate. After normalization, the first attempt at the model had a 97.8% accuracy. The model
60 then underwent 10 fold cross-validation and had a mean accuracy of 97.5%.

## 61 Support Vector Machines

62 Before we could start Support Vector Machines, I removed the the combined database that was adjusted for normality. Then I
63 recommitted the combined set of data and ran the support vector machines. The trial support vector machines had an accuracy
64 of 99.6% percent. The tuning then was tested for gamma and costs parameters. The first time that I tuned the function, the
65 upper extremities were the best parameters. I then attempted to tune again with the upper extremities as the lower boundaries

66  of the function. The second attempt at tuning the function resulted in the lower bounds, the original upper extremities, were

67  the optimal gamma and cost. Therefore these values were used and the model after tuning was 99.6%.

## 68 Artificial Neural Networks

69  The neural networks was run using the `nnet` package. The first time it was run with 2 levels and a maximum iteration at

70  100. The first accuracy was at 99.8% I changed the iterations to 500 and the size of the neural network to 10 hidden layers.

71  During the run of the neural network command. The first iterations stopped at 230 and the accuracy was 99.7%. The network

72  was tuned with sizes of hidden layers between 1 and 15. Also the tuning function was applied with controls of 5 repeats, cross

73  validation sampling, and 10 crosses. The tuning function ran cross validation and resulted in 97.2% accuracy. I was shocked at

74  this low value considering the first values that i got that were extremely high. I think that my first guess of 10 layers was the

75  best for the model and the sum from all other sizes of hidden layers caused the value to overall be lower.

# 76 Conclusion

77  Overall, Support Vector Machines were the optimal model for this data. Neural Networks and the C Tree Decision Trees were

78  the overall worst models for the data. As I was working with the data, I really did not consider the chemical make up of a red

79  wine and the chemical makeup of a white wine. I just wanted to run the data to see how it behaved on it's own. After looking

80  more into the data, there are certain key attributes that really determine the model. Alcohol, density, total sulfur dioxides, free

81  sulfur dioxides, chlorides, and pH were the best attributes that could determine the type of wine. Models where these were these

82  were used instead of the whole data could probably perform at similar accuracies to the models run with all attributes.
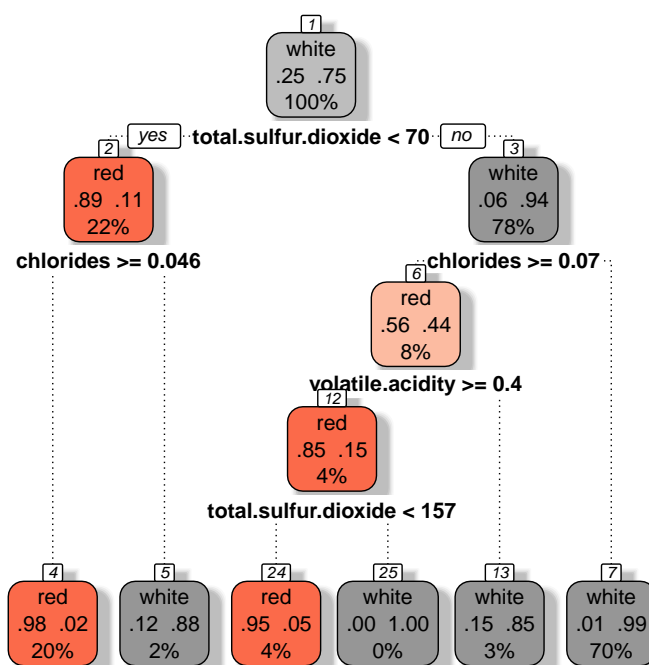
# 83 Code Appendix

```
84   set.seed(1506)
85   # Importing Data Sets
86   Red_Wine ← read.csv("~/Downloads/winequality−red.csv", sep = ";")
87   White_Wine ← read.csv("~/Downloads/winequality−white.csv", sep = ";")
88
89   # Adding the Type of Wine to the data set.
90   Red ← c(rep("red", nrow(Red_Wine)))
91   White ← c(rep("white", nrow(White_Wine)))
92   Red_Wine ← cbind(Red_Wine, Red)
93   White_Wine ← cbind(White_Wine, White)
94
95   # Changing the type of wine columns
96   colnames(Red_Wine)[13] ← "wine.type"
97   colnames(White_Wine)[13] ← "wine.type"
98
99   # Combining Both Red and White Wines into one data set
100  Total_Wine ← rbind(Red_Wine, White_Wine)
101  Total_Wine$quality ← as.factor(Total_Wine$quality)
102
103  # Item that will be predicted with the classification.
104  Wine.Type ← Total_Wine$wine.type
105
106  # Splitting Data for testing purposes
107  Wine.Split ← splitdata(Total_Wine, 0.7)
108  Wine.train ← Wine.Split$traindata
109  Wine.test ← Wine.Split$testdata
110
```

```
111   # Decision Trees to determine the type of wine, red or white.
112   wine.tree ← rpart(wine.type ∼ ., data = Wine.train)
113   plot.wine.tree ← fancyRpartPlot(wine.tree, palettes = c("Reds", "Greys"))
```



Rattle 2016−Dec−10 00:45:02 wdpryor1994

114

```
115   pred.tree ← predict(wine.tree, newdata = Wine.test, type = "class")
116
117   # Analysis of the decision tree
118   tree.conf ← confusion(Wine.test$wine.type, pred.tree, costs = NULL)
119   tree.conf
```

```
120
121   $counts
122          pred
123   true     red white
124     red    442    36
125     white   10  1461
126
127   $acc
128   [1] 0.9763982
129
130   $rates
131          pred
132   true            red        white
133     red    0.924686192 0.075313808
134     white  0.006798097 0.993201903
135
136   $sensitivities
137         red      white
138   0.9246862 0.9932019
139
140   $precisions
141         red      white
142   0.9778761 0.9759519
143
144   $F1s
```

```
145        red       white
146   0.9505376 0.9845013
147
148   $cost
149   [1] 46
150
```

```
151   # Tuning decision tree
152   tree.tune <- tune.rpart(wine.type ~ ., data = Total_Wine, cp = seq(0, 0.5, 0.001),
153       maxdepth = 1:10)
154   tree.tune.2 <- tune.rpart(wine.type ~ ., data = Total_Wine, minsplit = seq(50,
155       125, 5))
156   plot(tree.tune)
```

### Performance of 'rpart.wrap



```
157
158   plot(tree.tune.2)
```

## Performance of 'rpart.wrapper'



```
159

160   # The best parameters for the rpart decision tree
161   tree.tune$best.parameters

162
163           cp maxdepth
164   2507 0.001        6
165

166   tree.tune.2$best.parameters

167
168     minsplit
169   1       50
170

171   # Optimal rpart tree
172   tree.2 ← rpart(wine.type ~ ., data = Wine.train, control = rpart.control(minsplit =
173       tree.tune.2$best.parameters$minsplit,
174       cp = tree.tune$best.parameters$cp, maxdepth = tree.tune$best.parameters$maxdepth))
175   fancyRpartPlot(tree.2, palettes = c("Reds", "Greys"))
```

Rattle 2016–Dec–10 01:40:35 wdpryor1994

176

```
pred.tree.2 <- predict(tree.2, newdata = Wine.test, type = "class")

# Analysis of optimal decision tree with cross validation
tree.conf.2 <- confusion(Wine.test$wine.type, pred.tree.2, costs = NULL)
tree.conf.2
```

```
$counts
        pred
true      red white
   red    452    26
   white   12  1459

$acc
[1] 0.9805028

$rates
        pred
true             red         white
   red   0.945606695 0.054393305
   white 0.008157716 0.991842284

$sensitivities
       red      white
0.9456067 0.9918423

$precisions
       red      white
0.9741379 0.9824916

$F1s
       red      white
0.9596603 0.9871448

$cost
[1] 38
```

```
213  kflval ← function(k, data) {
214      folds = createfolds(nrow(data), k)
215      accvector = 1:k
216      for (k in 1:k) {
217          temptrain = data[folds != k, ]
218          temptest = data[folds == k, ]
219          temptree = rpart(wine.type ~ ., data = temptrain, control = rpart.control(minsplit =
220              tree.tune.2$best.parameters$minsplit,
221              cp = tree.tune$best.parameters$cp, maxdepth = tree.tune$best.parameters$maxdepth))
222          temppred = predict(temptree, newdata = temptest, type = "class")
223          analysis = confusion(temptest$wine.type, temppred)
224          accvector[k] = analysis$acc
225      }
226      return(mean(accvector))
227  }
228
229  kflval(10, Total_Wine)
```

```
[1] 0.9789119
```

```
233  # Ctree
234  wine.ctree ← ctree(wine.type ~ ., data = Wine.train)
235  plot(wine.ctree)
```



```
237  # Prediction and confusion of ctree decision tree
238  pred.tree.3 ← predict(wine.ctree, newdata = Wine.test)
239  wine.ctree.conf ← confusion(Wine.test$wine.type, pred.tree.3, costs = NULL)
240  wine.ctree.conf
```

```
$counts
        pred
true    red white
```

```
245    red      458      20
246    white    17    1454
247
248  $acc
249  [1] 0.9810159
250
251  $rates
252          pred
253  true           red        white
254    red    0.95815900 0.04184100
255    white 0.01155676 0.98844324
256
257  $sensitivities
258        red      white
259  0.9581590 0.9884432
260
261  $precisions
262        red      white
263  0.9642105 0.9864315
264
265  $F1s
266        red      white
267  0.9611752 0.9874363
268
269  $cost
270  [1]  37
271
```

```r
272  # Testing minsplit and maxdepth of ctree
273  maxdepth ← 1:10
274  # 50 minsplit
275  accvec ← NULL
276  for (i in 1:length(maxdepth)) {
277      temp.tree ← ctree(wine.type ~ ., data = Wine.train, controls = ctree_control(minsplit = 50,
278          maxdepth = maxdepth[i]))
279      temp.pred ← predict(temp.tree, newdata = Wine.test)
280      temp.conf ← confusion(Wine.test$wine.type, temp.pred, costs = NULL)
281      accvec[i] ← temp.conf$acc
282  }
283
284  # 100 misplit
285  accvec.2 ← NULL
286  for (i in 1:length(maxdepth)) {
287      temp.tree ← ctree(wine.type ~ ., data = Wine.train, controls = ctree_control(minsplit = 100,
288          maxdepth = maxdepth[i]))
289      temp.pred ← predict(temp.tree, newdata = Wine.test)
290      temp.conf ← confusion(Wine.test$wine.type, temp.pred, costs = NULL)
291      accvec.2[i] ← temp.conf$acc
292  }
293
294  # 150 minsplit
295  accvec.3 ← NULL
296  for (i in 1:length(maxdepth)) {
297      temp.tree ← ctree(wine.type ~ ., data = Wine.train, controls = ctree_control(minsplit = 150,
298          maxdepth = maxdepth[i]))
299      temp.pred ← predict(temp.tree, newdata = Wine.test)
300      temp.conf ← confusion(Wine.test$wine.type, temp.pred, costs = NULL)
301      accvec.3[i] ← temp.conf$acc
302  }
303
```

```
304  # Creating information matrix
305  max.depth ← as.vector(c(rep(list(1, 2, 3, 4, 5, 6, 7, 8, 9, 10), 3)))
306  min.split ← as.vector(c(rep(list(50, 100, 150), 10)))
307  accuracy ← as.vector(rbind(accvec, accvec.2, accvec.3))
308  acc.matrix ← cbind(max.depth, min.split, accuracy)
309  colnames(acc.matrix)[1] ← "maxdepth"
310  colnames(acc.matrix)[2] ← "misplit"
311  colnames(acc.matrix)[3] ← "accuracy"
312  which.max(accuracy)
```

```
313
314  [1] 19
315
```

```
316  acc.matrix[which.max(accuracy), ]
```

```
317
318  $maxdepth
319  [1] 9
320
321  $misplit
322  [1] 50
323
324  $accuracy
325  [1] 0.9712673
326
```

```
327  # Testing Optimal Ctree
328  opt.ctree ← ctree(wine.type ~ ., data = Wine.train, controls = ctree_control(minsplit = 50,
329      maxdepth = 10))
330  plot(opt.ctree)
```



```
331
```

```
332  opt.pred ← predict(opt.ctree, newdata = Wine.test)
333  opt.conf ← confusion(Wine.test$wine.type, opt.pred, costs = NULL)
334  opt.conf
```

```
335
336  $counts
337          pred
```

```
338  true     red  white
339    red     451     27
340    white    29   1442
341
342  $acc
343  [1] 0.9712673
344
345  $rates
346          pred
347  true              red        white
348    red    0.94351464 0.05648536
349    white 0.01971448 0.98028552
350
351  $sensitivities
352        red      white
353  0.9435146 0.9802855
354
355  $precisions
356        red      white
357  0.9395833 0.9816201
358
359  $F1s
360        red      white
361  0.9415449 0.9809524
362
363  $cost
364  [1] 56
365
```

```r
366  # Cross Validation with C Tree
367  ctree.kflval ← function(k, data) {
368      folds = createfolds(nrow(data), k)
369      accvector = 1:k
370      for (k in 1:k) {
371          temptrain = data[folds != k, ]
372          temptest = data[folds == k, ]
373          temptree = ctree(wine.type ~ ., data = temptrain, controls = ctree_control(minsplit = 50,
374              maxdepth = 10))
375          temppred = predict(temptree, newdata = temptest)
376          analysis = confusion(temptest$wine.type, temppred)
377          accvector[k] = analysis$acc
378      }
379      return(mean(accvector))
380  }
381
382  ctree.kflval(10, Total_Wine)
383
```

```
384  [1] 0.9738326
385
```

```r
386  # K nearest neighbors standardize the data
387  x = Total_Wine[, 1:11]
388  xbar = apply(x, 2, mean)
389  xbarMat = cbind(rep(1, nrow(Total_Wine))) %*% xbar
390  s = apply(x, 2, sd)
391  sMat = cbind(rep(1, nrow(Total_Wine))) %*% s
392  z = (x − xbarMat)/sMat
393
394  # K nearest neighbors sampling the data
395  z.split ← sample(nrow(z), round(nrow(z) * 0.7, 0))
396  z.train ← z[z.split, ]
```

```
397  z.test ← z[−z.split , ]
398
399  # K nearest neighbors Test
400  wine.knn ← knn(train = z.train , test = z.test , k = 3, cl = Wine.Type[z.split])
401  confusion (Wine.Type[−z.split], wine.knn , costs = NULL)
```
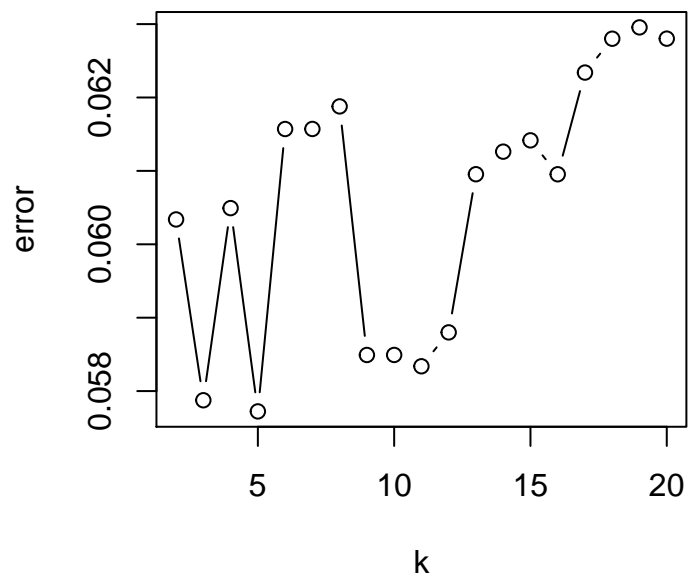
```
$counts
        pred
true      red white
  red     478       7
  white     6   1458

$acc
[1] 0.9933299

$rates
        pred
true              red         white
  red    0.985567010 0.014432990
  white 0.004098361 0.995901639

$sensitivities
      red     white
0.9855670 0.9959016

$precisions
      red     white
0.9876033 0.9952218

$F1s
      red     white
0.9865841 0.9955616

$cost
[1] 13
```

```
433  # Finding the optimal k for K nearest neighbors
434  x ← Total_Wine[, −13]
435  y ← Total_Wine[, 13]
436  knn.tune ← tune.knn(x, y, k = seq(2, 20, 1), tunecontrol = tune.control(sampling = "cross"))
437  knn.tune.2 ← tune.knn(x, y, k = seq(2, 20, 1), tunecontrol = tune.control(sampling = "fix"))
438  knn.tune.3 ← tune.knn(x, y, k = seq(2, 20, 1), tunecontrol = tune.control(sampling = "boot"))
439  plot(knn.tune)
```

**Performance of 'knn.wrapper'**



440

```
441   plot(knn.tune.2)
```

**Performance of 'knn.wrapper'**



442

```
443   plot(knn.tune.3)
```

## Performance of 'knn.wrapper'



```
444

445   knn.tune$best.parameters$k

446
447   [1] 5
448

449   knn.tune.2$best.parameters$k

450
451   [1] 4
452

453   knn.tune.3$best.parameters$k

454
455   [1] 14
456

457   wine.knn.2 ← knn(train = z.train, test = z.test, k = knn.tune$best.parameters$k,
458       cl = Wine.Type[z.split])
459   wine.knn.3 ← knn(train = z.train, test = z.test, k = knn.tune.2$best.parameters$k,
460       cl = Wine.Type[z.split])
461   wine.knn.4 ← knn(train = z.train, test = z.test, k = knn.tune.3$best.parameters$k,
462       cl = Wine.Type[z.split])
463   wine.knn.2.conf ← confusion(Wine.Type[−z.split], wine.knn.2, costs = NULL)
464   wine.knn.3.conf ← confusion(Wine.Type[−z.split], wine.knn.3, costs = NULL)
465   wine.knn.4.conf ← confusion(Wine.Type[−z.split], wine.knn.4, costs = NULL)
466
467   # Cross Validation for K nearest neighbors
468   cross.knn ← knn.cv(train = z, cl = y, k = knn.tune$best.parameters$k)
469   cross.knn.conf ← confusion(Wine.Type, cross.knn)
470
471   # Weighted K nearest neighbors
472   z ← cbind(z, Wine.Type)
473   z.train ← z[z.split, ]
474   z.test ← z[−z.split, ]
475   weight.knn.wine ← kknn(Wine.Type ~ ., train = z.train, test = z.test, k = knn.tune$best.parameters$k,
476       kernel = "triangular")
477   pred.weight ← predict(weight.knn.wine, newdata = z.test)
478   confusion(z.test$Wine.Type, pred.weight)
```

```
$counts
        pred
true      red white
   red     479     6
   white     4   1460

$acc
[1] 0.9948692

$rates
        pred
true              red        white
   red    0.98762887 0.01237113
   white 0.00273224 0.99726776

$sensitivities
       red      white
0.9876289 0.9972678

$precisions
       red      white
0.9917184 0.9959072

$F1s
       red      white
0.9896694 0.9965870

$cost
[1] 10
```

```r
# Finding the best kernels
tune.kknn <- function(K, kernels) {
    acc <- NULL
    for (i in 1:length(K)) {
        temp.model <- kknn(Wine.Type ~ ., train = z.train, test = z.test, k = i,
            kernel = kernels, distance = 2)
        pred.temp <- predict(temp.model, newdata = z.test)
        temp.conf <- confusion(z.test$Wine.Type, pred.temp)
        acc[i] <- temp.conf$acc
    }
    opt.k <- which.max(acc)
    opt.acc <- acc[opt.k]
    return(list(k = opt.k, accuracy = opt.acc))
}
K <- 2:30
rect <- tune.kknn(K, "rectangular")
triangle <- tune.kknn(K, "triangular")
ep <- tune.kknn(K, "epanechnikov")
biw <- tune.kknn(K, "biweight")
triw <- tune.kknn(K, "triweight")
cosine <- tune.kknn(K, "cos")
invert <- tune.kknn(K, "inv")
gauss <- tune.kknn(K, "gaussian")
rank.weight <- tune.kknn(K, "rank")
optimal <- tune.kknn(K, "optimal")
weighted.accuracy <- list(rect$accuracy, triangle$accuracy, ep$accuracy, biw$accuracy,
    triw$accuracy, cosine$accuracy, invert$accuracy, gauss$accuracy, rank.weight$accuracy,
    optimal$accuracy)
```

```
538  weighted.k ← list(rect$k, triangle$k, ep$k, biw$k, triw$k, cosine$k, invert$k,
539      gauss$k, rank.weight$k, optimal$k)
540
541  opt.weights ← cbind(weighted.k, weighted.accuracy)
542  opt.kknn ← kknn(Wine.Type ~ ., train = z.train, test = z.test, k = 3, kernel = "inv")
543  pred.kknn ← predict(opt.kknn, newdata = z.test)
544  confusion(z.test$Wine.Type, pred.kknn)
545
```

```
546  $counts
547          pred
548  true      red white
549    red     479      6
550    white     5   1459
551
552  $acc
553  [1] 0.9943561
554
555  $rates
556          pred
557  true              red        white
558    red    0.987628866 0.012371134
559    white 0.003415301 0.996584699
560
561  $sensitivities
562        red      white
563  0.9876289 0.9965847
564
565  $precisions
566        red      white
567  0.9896694 0.9959044
568
569  $F1s
570        red      white
571  0.9886481 0.9962445
572
573  $cost
574  [1]  11
575
```

```
576  # Cross Validation fro Weighted K nearest neighbors
577  kknn.kflval ← function(k, data) {
578      folds = createfolds(nrow(data), k)
579      accvector = 1:k
580      for (k in 1:k) {
581          temptrain = data[folds != k, ]
582          temptest = data[folds == k, ]
583          tempmodel = kknn(Wine.Type ~ ., train = temptrain, test = temptest,
584              k = 3, kernel = "inv")
585          temppred = predict(tempmodel, newdata = temptest, type = "raw")
586          analysis = confusion(temptest$Wine.Type, temppred)
587          accvector[k] = analysis$acc
588      }
589      return(mean(accvector))
590  }
591
592  kknn.kflval(10, z)
593
```

```
594  [1] 0.9941515
595
```

```
596  # Investigating qualities of data Naive Bayes
597  normality(Wine.train$fixed.acidity)
```

## Histogram of data



```
598
599
600
601          Shapiro-Wilk normality test
602
603   data:  data
604   W = 0.88403, p-value < 2.2e-16
605
```
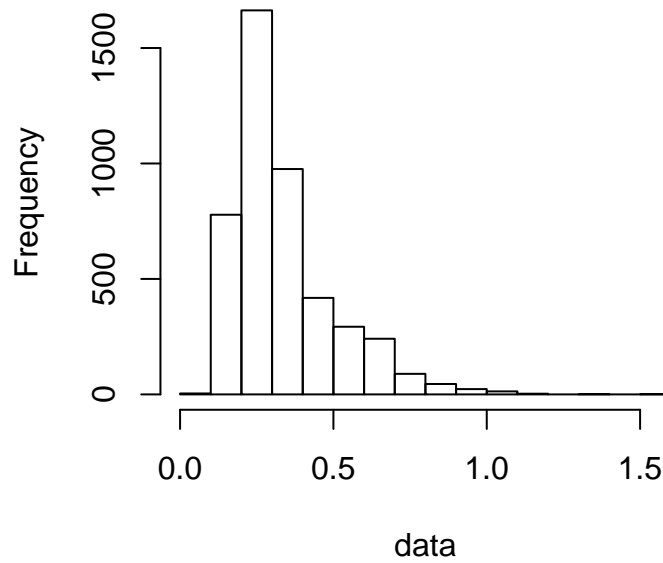
## Normal Q–Q Plot



```
606
607   normality(Wine.train$volatile.acidity)
```

## Histogram of data



```
608
609
610
611          Shapiro-Wilk normality test
612
613  data:  data
614  W = 0.87268, p-value < 2.2e-16
615
```
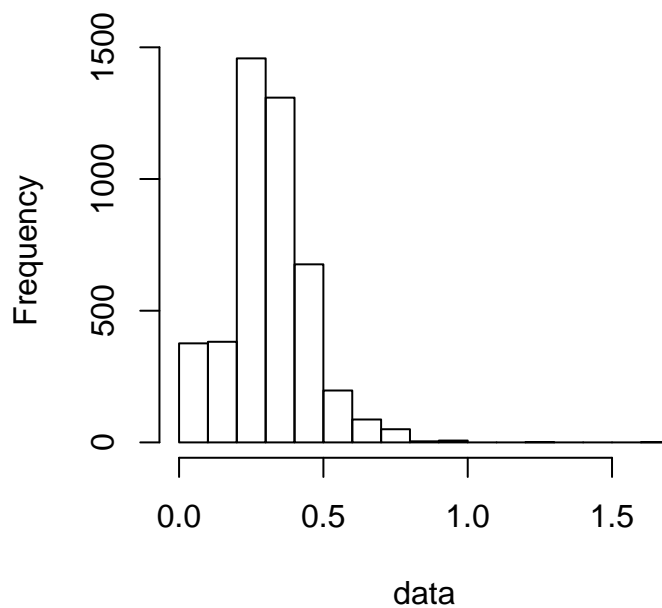
## Normal Q–Q Plot



```
616

617  normality(Wine.train$citric.acid)
```
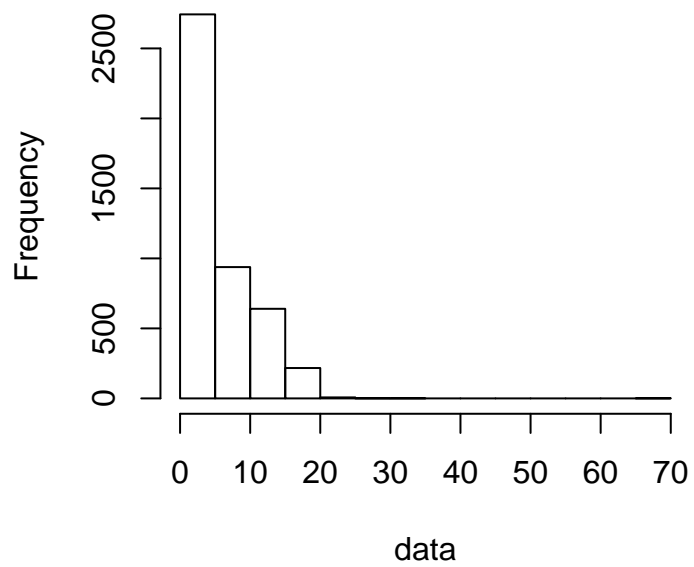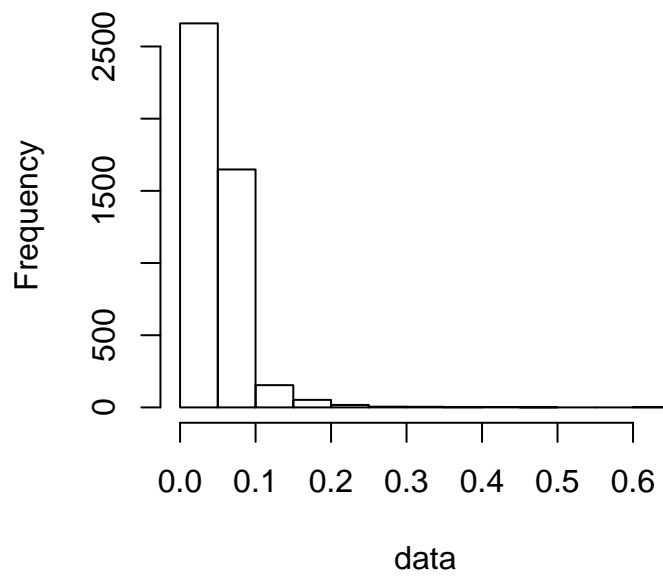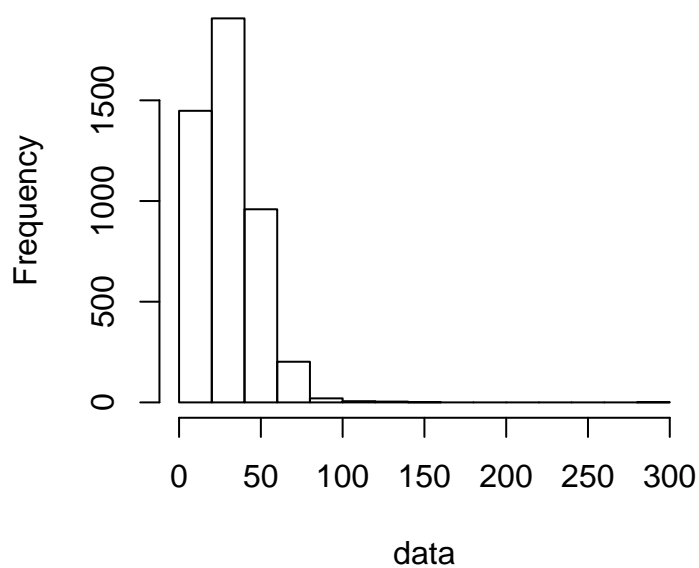
## Histogram of data

```
618
619
620
621          Shapiro - Wilk normality test
622
623   data:   data
624   W = 0.96462,  p-value < 2.2e-16
625
```

## Normal Q–Q Plot



```
626

627   normality ( Wine.train$residual.sugar )
```

## Histogram of data



```
628
629
630
631          Shapiro-Wilk normality test
632
633   data:  data
634   W = 0.82197, p-value < 2.2e-16
635
```

## Normal Q–Q Plot



```
636
637   normality(Wine.train$chlorides)
```

## Histogram of data

```
638
639
640
641         Shapiro-Wilk normality test
642
643    data:   data
644    W = 0.64501, p-value < 2.2e-16
645
```

## Normal Q–Q Plot



```
646

647    normality(Wine.train$free.sulfur.dioxide)
```

## Histogram of data



data

```
648
649
650
651         Shapiro-Wilk normality test
652
653  data:  data
654  W = 0.93333, p-value < 2.2e-16
655
```
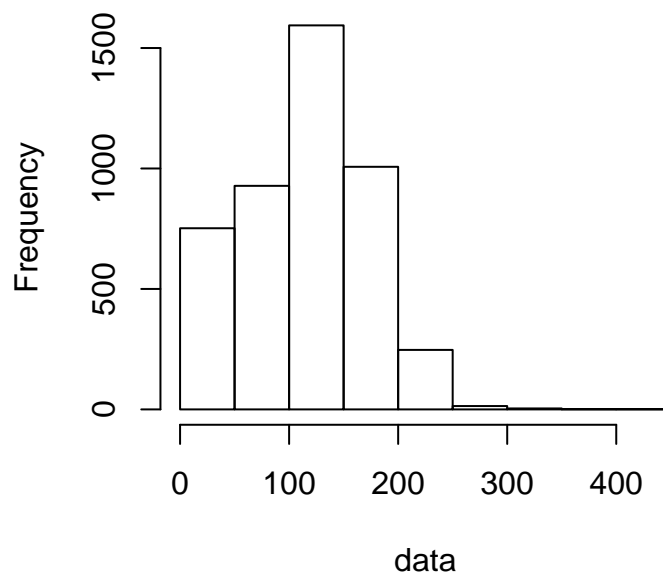
## Normal Q–Q Plot



Theoretical Quantiles

```
656

657  normality(Wine.train$total.sulfur.dioxide)
```

## Histogram of data



```
658
659
660
661            Shapiro-Wilk normality test
662
663   data:  data
664   W = 0.98268, p-value < 2.2e-16
665
```

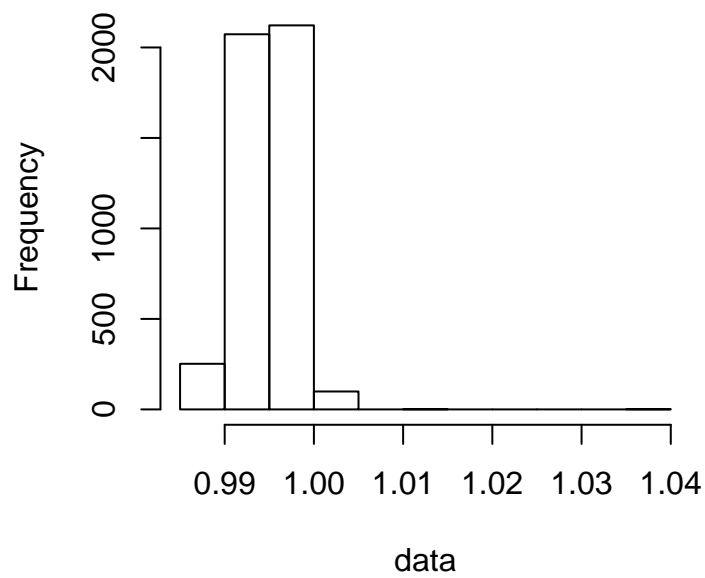## Normal Q–Q Plot



```
666
667   normality(Wine.train$density)
```

## Histogram of data



```
668
669
670
671        Shapiro-Wilk normality test
672
673  data:  data
674  W = 0.95975, p-value < 2.2e-16
675
```
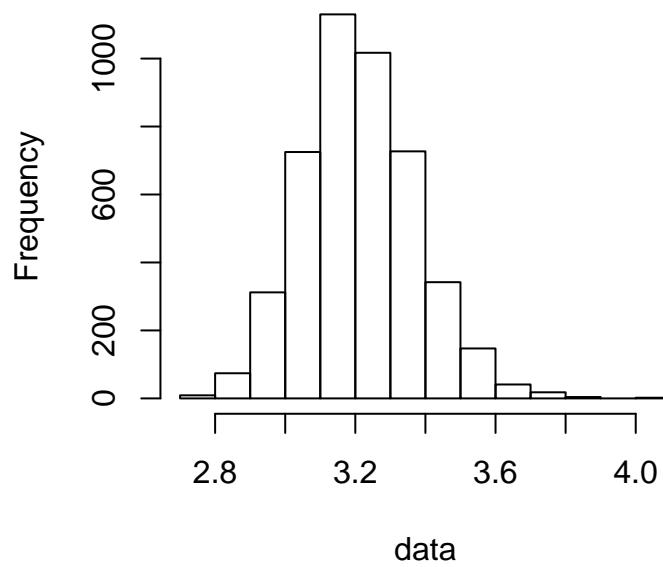
## Normal Q–Q Plot



```
676

677  normality(Wine.train$pH)
```
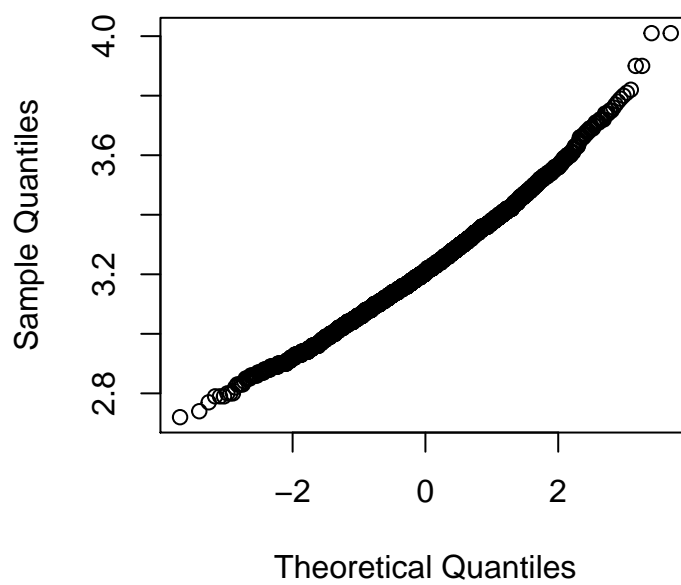
## **Histogram of data**



```
678
679
680
681         Shapiro-Wilk normality test
682
683   data:  data
684   W = 0.99181, p-value = 1.555e-15
685
```

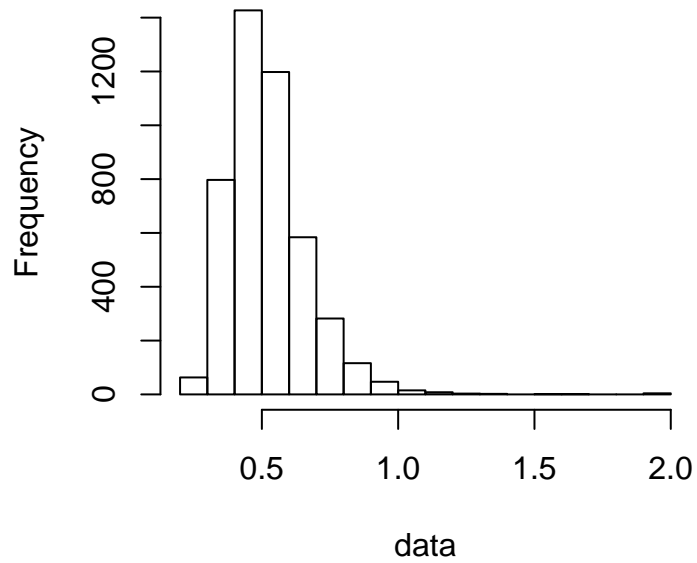## **Normal Q–Q Plot**



```
686
687   normality(Wine.train$sulphates)
```
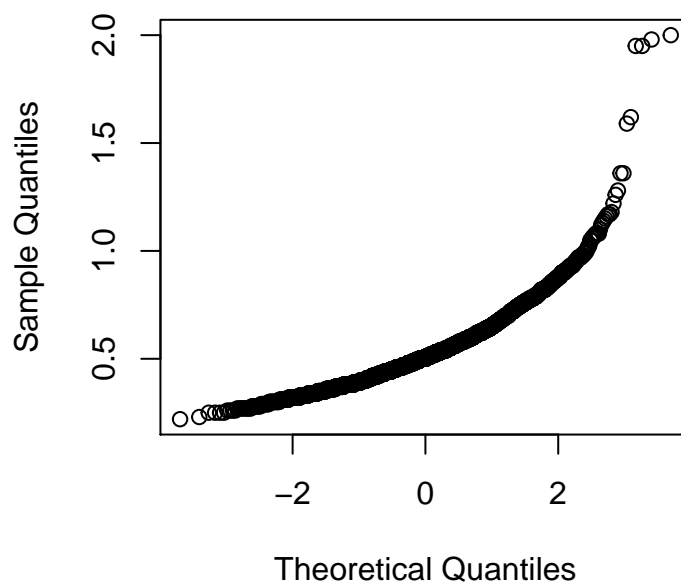
## Histogram of data



```
688
689
690
691         Shapiro - Wilk normality test
692
693 data:   data
694 W = 0.89657 , p-value < 2.2e-16
695
```

## Normal Q–Q Plot



```
696
697 normality ( Wine.train$alcohol )
```
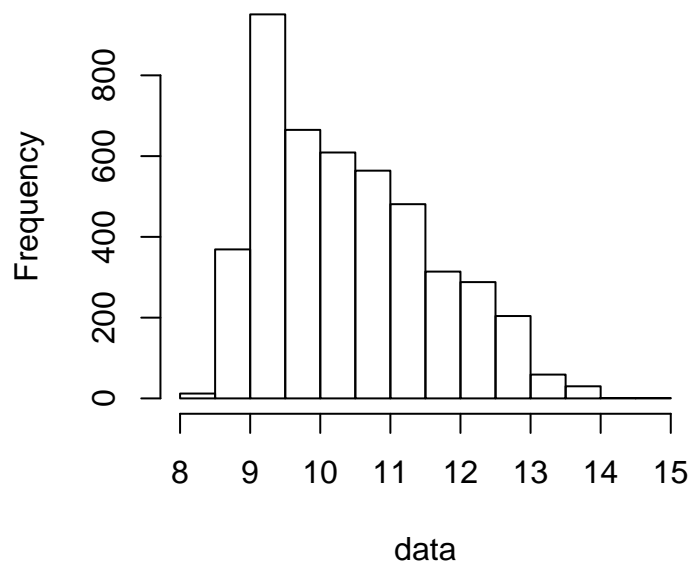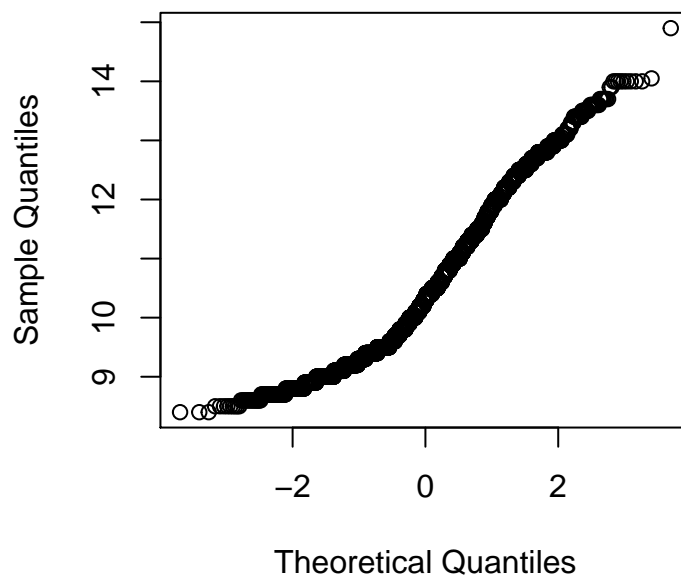
# Histogram of data



```
698
699
700
701        Shapiro-Wilk normality test
702
703   data:  data
704   W = 0.95221, p-value < 2.2e-16
705
```

# Normal Q–Q Plot



```
706

707   # Discretizing data
708   Total_Wine$ fixed.acidity   ← mycut( Total_Wine$ fixed.acidity , 10)
709   Total_Wine$ volatile.acidity ← mycut( Total_Wine$ volatile.acidity , 10)
```

```
710  Total_Wine$citric.acid ← mycut(Total_Wine$citric.acid, 10)
711  Total_Wine$residual.sugar ← mycut(Total_Wine$residual.sugar, 10)
712  Total_Wine$chlorides ← mycut(Total_Wine$chlorides, 10)
713  Total_Wine$free.sulfur.dioxide ← mycut(Total_Wine$free.sulfur.dioxide, 10)
714  Total_Wine$total.sulfur.dioxide ← mycut(Total_Wine$total.sulfur.dioxide, 10)
715  Total_Wine$density ← mycut(Total_Wine$density, 10)
716  Total_Wine$pH ← mycut(Total_Wine$pH, 10)
717  Total_Wine$sulphates ← mycut(Total_Wine$sulphates, 10)
718  Total_Wine$alcohol ← mycut(Total_Wine$alcohol, 10)
719  Total_Wine$quality ← as.factor(Total_Wine$quality)
720
721  Wine.Split ← splitdata(Total_Wine, 0.7)
722  Wine.train ← Wine.Split$traindata
723  Wine.test ← Wine.Split$testdata
724
725  # Applying Naive Bayes
726  nb.wine ← naiveBayes(wine.type ~ ., data = Wine.train)
727  pred.nb.wine ← predict(nb.wine, newdata = Wine.test, type = "class")
728  nb.conf ← confusion(Wine.test$wine.type, pred.nb.wine)
729
730  # Cross Validation
731  nb.kflval ← function(k, data) {
732      folds = createfolds(nrow(data), k)
733      accvector = 1:k
734      for (k in 1:k) {
735          temptrain = data[folds != k, ]
736          temptest = data[folds == k, ]
737          tempmodel = naiveBayes(wine.type ~ ., data = temptrain)
738          temppred = predict(tempmodel, newdata = temptest, type = "class")
739          analysis = confusion(temptest$wine.type, temppred)
740          accvector[k] = analysis$acc
741      }
742      return(mean(accvector))
743  }
744
745  nb.kflval(10, Total_Wine)
```

```
[1] 0.9898424
```

```
749  # Support Vector Machines
750  rm(Total_Wine)
751  Total_Wine ← rbind(Red_Wine, White_Wine)
752  wine.svm ← svm(wine.type ~ ., data = Wine.train, kernel = "linear", cost = 1,
753      scale = T)
754  pred.svm ← predict(wine.svm, newdata = Wine.test)
755  svm.conf ← confusion(Wine.test$wine.type, pred.svm)
756  plot(wine.type ~ ., wine.svm, data = Wine.train)
```

wine.type

fixed.acidity = Wine.train, kernel = "linear", cost = 1, scale =

[3.8,6]   (6.5,6.8]   (7.2,7.5]   (8.8,15.9]

0.0   0.2   0.4   0.6   0.8   1.0

fixed.acidity

757

wine.type

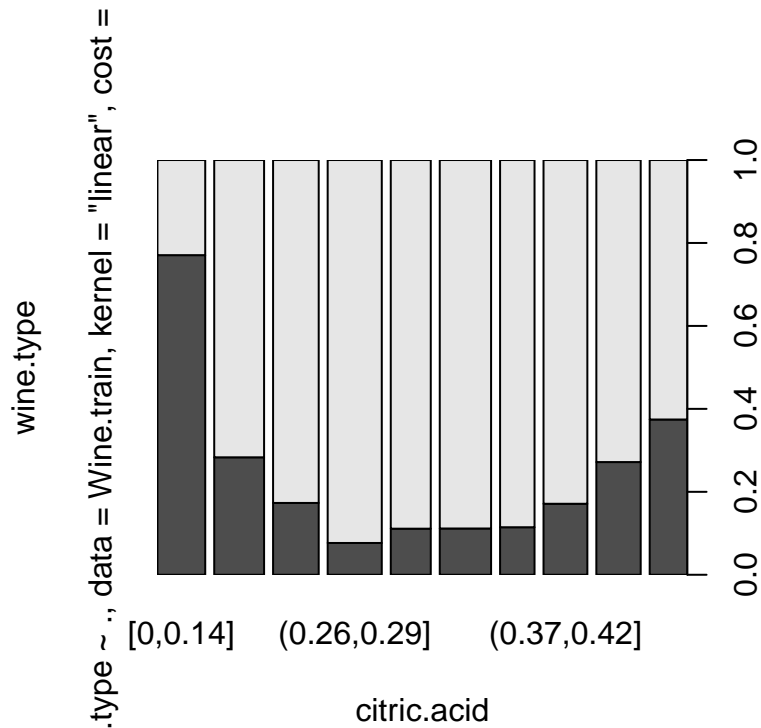type ~ ., data = Wine.train, kernel = "linear", cost = 1, s  .,  data = Wine.train, kernel = "linear", cost = 1, s  .,  data = Wine.train, kernel = "linear", cost = 1, scale =

[0.08,0.18]   (0.27,0.29]   (0.45,0.59]

0.0   0.2   0.4   0.6   0.8   1.0

volatile.acidity

758

wine.type

= cost ,"linear" = kernel Wine.train, = data ,~ine.type s 1, = cost ,"linear" = kernel Wine.train, = data .,~ type

citric.acid

[0,0.14]   (0.26,0.29)   (0.37,0.42]

759

wine.type

ype ~ ., data = Wine.train, kernel = "linear", cost = 1, s

residual.sugar

[0.6,1.3]   (2,2.3]   (5,7.1]   (13,65.8]

760

wine.type

data = Wine.train, kernel = "linear", cost = 1, s

type ~ ., data = Wine.train, kernel = "linear", cost =

[0.009,0.031]   (0.044,0.047]   (0.086,0.611]

chlorides

761

wine.type

ne.type ~ ., data = Wine.train, kernel = "linear", cost =

[1,9]   (15,19)   (29,33)   (45,54)

free.sulfur.dioxide

762

wine.type

total.sulfur.dioxide

[6,30]     (89,105]     (132,148]

763



wine.type

density

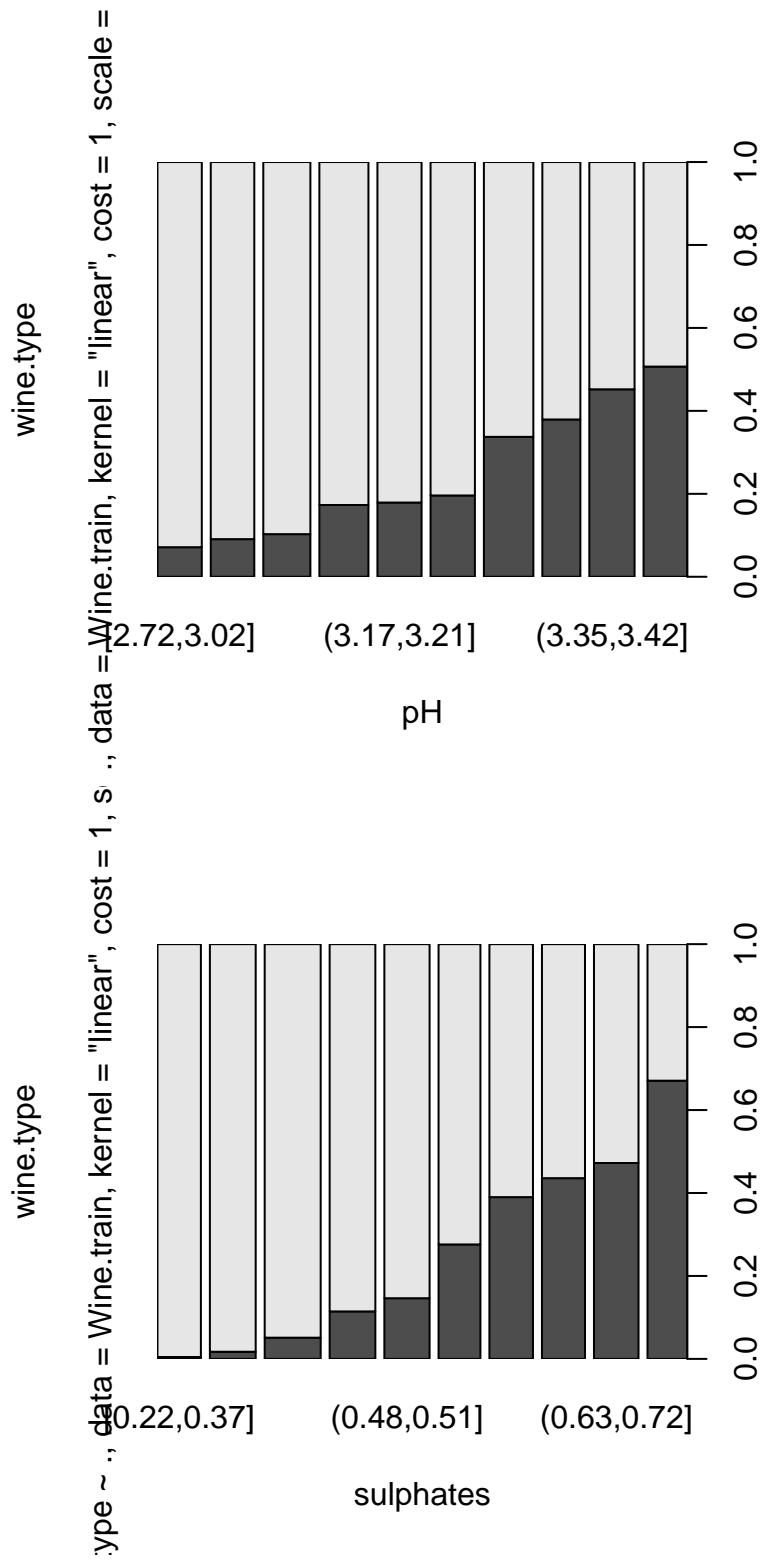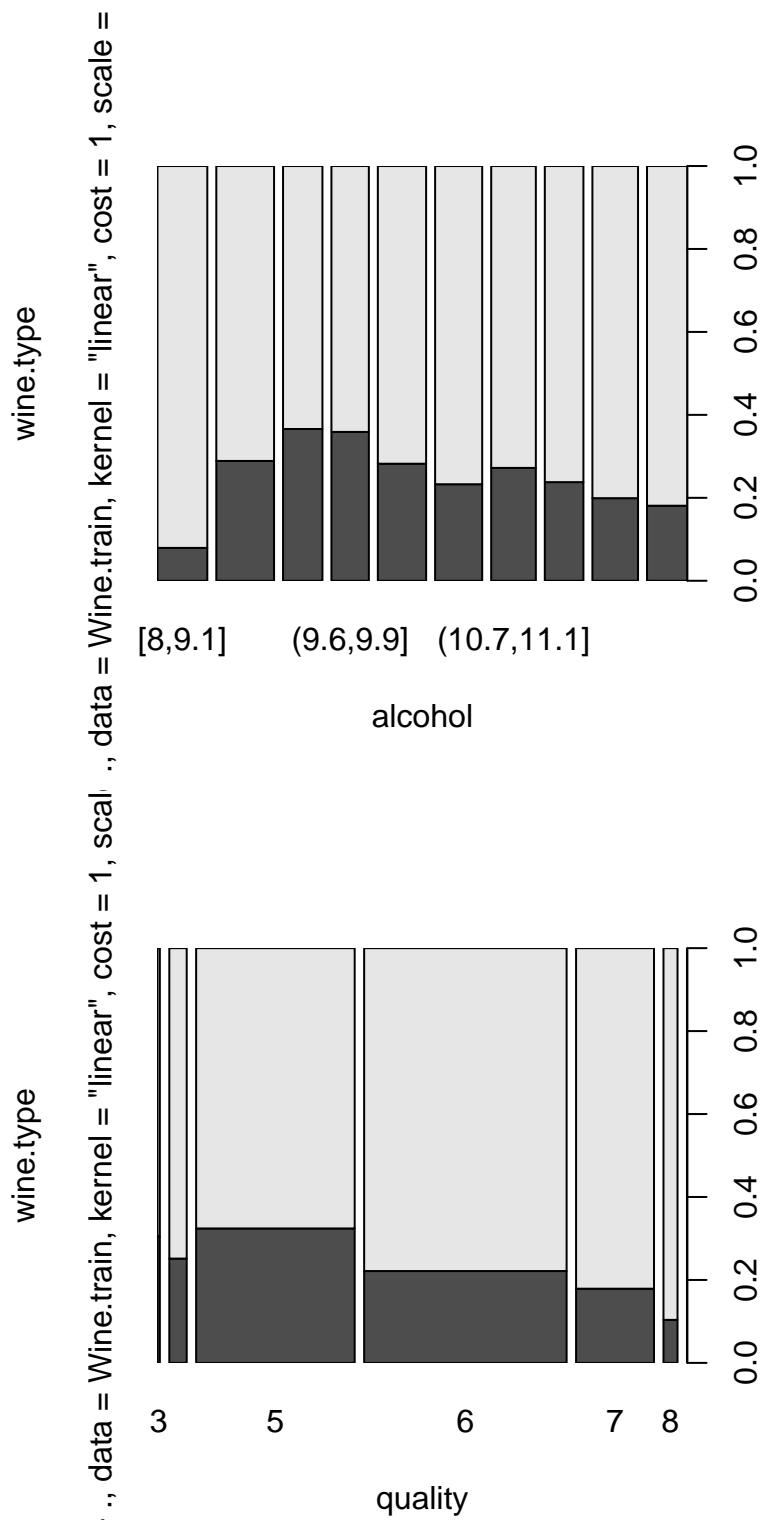[0.9871,0.9907]     (0.9949,0.9957]

764

765



766

767

768

```r
# Tuning gamma and cost on SVM
svm.tuning <- tune.svm(wine.type ~ ., data = Total_Wine, gamma = 10^(-6:-3),
    cost = 10^(1:2))
svm.tuning$best.parameters
```

```
  gamma cost
8 0.001  100
```

```r
svm.tuning.2 <- tune.svm(wine.type ~ ., data = Total_Wine, gamma = 10^(-3:2),
```

```r
778        cost = 10^(2:4))
779
780 # Optmal SVM
781 opt.svm <- svm(wine.type ~ ., data = Wine.train, kernel = "linear", gamma =
782     svm.tuning.2$best.parameters$gamma,
783     cost = svm.tuning.2$best.parameters$cost)
784
785 # Cross Validation of SVM
786 svm.kflval <- function(k, data) {
787     folds = createfolds(nrow(data), k)
788     accvector = 1:k
789     for (k in 1:k) {
790         temptrain = data[folds != k, ]
791         temptest = data[folds == k, ]
792         tempmodel = svm(wine.type ~ ., data = Wine.train, kernel = "linear",
793             gamma = svm.tuning.2$best.parameters$gamma, cost = svm.tuning.2$best.parameters$cost)
794         temppred = predict(tempmodel, newdata = temptest)
795         analysis = confusion(temptest$wine.type, temppred)
796         accvector[k] = analysis$acc
797     }
798     return(mean(accvector))
799 }
800
801 svm.kflval(10, Total_Wine)
```

```
802
803 Error in predict.svm(tempmodel, newdata = temptest): test data does not match model !
804
```

```r
805 # Neural Networks
806 wine.nnet <- nnet(wine.type ~ ., data = Wine.test, size = 10, linout = FALSE,
807     maxit = 500)
```

```
808
809 Error in nnet.default(x, y, w, entropy = TRUE, ...): too many (1071) weights
810
```

```r
811 pred.nnet <- predict(wine.nnet, newdata = Wine.test, type = "class")
```

```
812
813 Error in predict(wine.nnet, newdata = Wine.test, type = "class"): object 'wine.nnet' not found
814
```

```r
815 confusion(Wine.test$wine.type, pred.nnet)
```

```
816
817 Error in table(true, pred): object 'pred.nnet' not found
818
```

```r
819 nnet.tuning <- tune.nnet(wine.type ~ ., data = Total_Wine, size = 1:15, trace = FALSE,
820     tunecontrol = tune.control(nrepeat = 5, sampling = "cross", cross = 10,
821         ))
822 1 - mean(nnet.tuning$performances$error)
```

```
823
824 [1] 0.9754159
825
```

```r
826 # Plotting Neaural Networks import the function from Github
827 library(devtools)
828 source_url("https://gist.githubusercontent.com/fawda123/7471137/raw/466c1474d0a505ff044412703516c34f1a4684a5/nnet_
```

```
829
830 SHA-1 hash of file is 74c80bd5ddbc17ab3ae5ece9c0ed9beb612e87ef
831
```

```r
832 plot.nnet(wine.nnet)
```

```
833
834 Loading required package: scales
835
```

836
837
838
839

```
Attaching package: 'scales'
```

840
841

```
The following object is masked from 'package:kernlab':

    alpha
```

842
843
844

845
846
847

```
Error in match(x, table, nomatch = 0L): object 'wine.nnet' not found
```