**Douglas Rowe**

# Final Project

**M5364–Data Mining   Fall 2016   Tarleton State Univ   Dr. Scott Cook   Assigned 2016-09-30   Due 2016-12-10**

---

1  This data set is called the abalone data set. Abalone's are a form of sea snails whose age can be determined by boring into
2  the shell and, through a "boring and time consuming" process, by counting the rings much like a tree. In order to save our
3  conchologist friends time, we will attempt to predict the number of rings in the shell by measuring:

4  1) Sex

5  2) Length

6  3) Diameter

7  4) Height

8  5) Whole Weight

9  6) Shucked Weight

10  7) Viscera weight

11  8) Shell weight

12  Preprocessing:

```
data = read.csv("C:/Users/Rowe/Documents/My Repo/Final Project DM/abalone.data.txt",
    header = FALSE)
# We have enough data from the UCI website to know that none of the data
# points are missing values or Nan's.
table(data[, 1])
```

```
   F    I    M
1307 1342 1528
```

```
# Because there is almost even number of 'M', 'F' and 'I' values here, there
# is minimal chance of a class imbalance problem coming from this variable.
my_shapiro(data[, 2:8])
```

```
[1] 7.442090e-29 1.648335e-28 1.181265e-47 1.013778e-27 9.340986e-32
[6] 1.777103e-29 1.565014e-28
```

```
# From this function, we can confidently say that none of the continuous
# variables are normal.
```

32  Beginning:

```
# This function calculates a 1x2 matrix whose values are the Root Mean
# Square Error (RMSE) and the Mean Absolute Error (MAE) respectively.
RMSE_MAE = function(true, predicted) {
    error_matrix = matrix(nrow = 1, ncol = 2)
    error_matrix[1, 1] = sqrt((1/length(true)) * sum((predicted - true)^2))
    error_matrix[1, 2] = (1/length(true)) * sum(abs(true - predicted))
    return(error_matrix)
}
# I am creating master train and test set so we can compare the acccuracies
# of the models directly
set.seed(5364)
split_data = splitdata(data, 0.7)
train = split_data$train
test = split_data$test
train_rows = split_data$train_rows
```

48   Artificial Neural Networks:

```
49   Error_matrix = matrix(ncol = 2, nrow = 15)
50   for (i in 1:15) {
51       crude_ann = nnet(V9 ~ ., data = train, size = as.numeric(i), maxit = 1000,
52           linout = TRUE, trace = FALSE)
53       predicted_ann = predict(crude_ann, newdata = test)
54       Error_matrix[i, ] = RMSE_MAE(test[, 9], predicted_ann)
55   }
56   Error_matrix[argmin(Error_matrix[, 1]), ]
```

```
57
58   [1] 2.087510 1.464726
59
```

```
60   Error_matrix[argmin(Error_matrix[, 2]), ]
```

```
61
62   [1] 2.087510 1.464726
63
```

```
64   RMSE_size = argmin(Error_matrix[, 1])
65   MAE_size = argmin(Error_matrix[, 2])
66   ann_cv = function(size, d, m, data_frame) {
67       error_matrix = matrix(nrow = m + 1, ncol = 2)
68       for (i in 1:m) {
69           d_records = sample(nrow(data_frame), d)
70           temp_train = data_frame[d_records, ]
71           temp_test = data_frame[-d_records, ]
72           model = nnet(V9 ~ ., data = temp_train, size = size, maxit = 1000, linout = TRUE,
73               trace = FALSE)
74           predicted = predict(model, temp_test)
75           error_matrix[i, ] = RMSE_MAE(predicted, temp_test[, 9])
76       }
77       error_matrix[m + 1, 1] = mean(error_matrix[1:m, 1])
78       error_matrix[m + 1, 2] = mean(error_matrix[1:m, 2])
79       return(error_matrix)
80   }
81   ann_cv(RMSE_size, round(0.7 * nrow(data)), 25, data)
```

```
82
83              [,1]      [,2]
84    [1,] 2.095105 1.485442
85    [2,] 2.207351 1.541343
86    [3,] 2.039443 1.449004
87    [4,] 2.065539 1.448992
88    [5,] 2.064905 1.466190
89    [6,] 2.002281 1.447325
90    [7,] 2.123396 1.480538
91    [8,] 2.070027 1.465895
92    [9,] 2.086959 1.464638
93   [10,] 2.070136 1.474192
94   [11,] 2.005938 1.410567
95   [12,] 2.150864 1.485463
96   [13,] 2.177452 1.523892
97   [14,] 2.172037 1.501955
98   [15,] 1.966703 1.418677
99   [16,] 2.111945 1.465941
100  [17,] 1.962255 1.428519
101  [18,] 2.083109 1.427756
102  [19,] 1.980216 1.427262
103  [20,] 2.082509 1.491201
104  [21,] 2.059137 1.480417
105  [22,] 2.068889 1.454277
106  [23,] 2.017523 1.434391
```

```
107  [24,]  2.146891  1.493057
108  [25,]  2.134369  1.493402
109  [26,]  2.077799  1.466413
110
```

```
111  ann_cv(MAE_size, round(0.7 * nrow(data)), 25, data)
```

```
112
113            [,1]        [,2]
114  [1,]    2.829343  1.497518
115  [2,]  165.295623  9.595944
116  [3,]    2.155146  1.506562
117  [4,]    2.038264  1.465907
118  [5,]    2.036226  1.422680
119  [6,]    2.074686  1.468382
120  [7,]    2.136545  1.486556
121  [8,]    2.132773  1.493238
122  [9,]    2.097541  1.504436
123  [10,]   2.083242  1.471270
124  [11,]   2.155716  1.520520
125  [12,]   2.111108  1.493397
126  [13,]   2.071499  1.478480
127  [14,]   2.046546  1.438408
128  [15,]   2.038515  1.423895
129  [16,]   1.974339  1.396308
130  [17,]   2.105372  1.502165
131  [18,]   2.082981  1.478904
132  [19,]   2.127992  1.514478
133  [20,]   2.071086  1.455224
134  [21,]   2.243506  1.582867
135  [22,]   2.017569  1.465865
136  [23,]   2.208430  1.496113
137  [24,]   2.159128  1.542066
138  [25,]   2.192894  1.563071
139  [26,]   8.659443  1.810570
140
```

```
141  ann_n_fold_cross = function(n, data_frame, size) {
142      chopping_point = floor(nrow(data_frame)/n)
143      stopping_points = c(1/chopping_point, 1:n) * chopping_point
144      error_matrix = matrix(ncol = 2, nrow = n + 1)
145      for (i in 1:n) {
146          lower_bound = stopping_points[i]
147          upper_bound = stopping_points[i + 1]
148          temp_test = data_frame[lower_bound:upper_bound, ]
149          temp_train = data_frame[-(lower_bound:upper_bound), ]
150          model = nnet(V9 ~ ., data = temp_train, size = size, maxit = 1000, linout = TRUE,
151              trace = FALSE)
152          predicted = predict(model, temp_test)
153          error_matrix[i, ] = RMSE_MAE(predicted, temp_test[, 9])
154      }
155      error_matrix[n + 1, 1] = mean(error_matrix[1:n, 1])
156      error_matrix[n + 1, 2] = mean(error_matrix[1:n, 2])
157      return(error_matrix)
158  }
159  ann_n_fold_cross(10, data, RMSE_size)
```

```
160
161            [,1]       [,2]
162  [1,]  2.648594  1.955683
163  [2,]  3.260283  2.327803
164  [3,]  1.612628  1.217336
165  [4,]  2.050177  1.595570
```

```
166    [5,]  1.587534 1.187856
167    [6,]  2.845384 2.046782
168    [7,]  1.385622 1.054785
169    [8,]  2.188472 1.620859
170    [9,]  1.639786 1.196326
171   [10,]  1.957278 1.408840
172   [11,]  2.117576 1.561184
173
```

```
174   ann_n_fold_cross(10, data, MAE_size)
```

```
175
176            [,1]       [,2]
177    [1,]  2.587194 1.911639
178    [2,]  3.179929 2.274739
179    [3,]  1.722674 1.349745
180    [4,]  2.050171 1.595562
181    [5,]  1.583929 1.190128
182    [6,]  2.856398 2.044390
183    [7,]  1.349745 1.027931
184    [8,]  3.225724 2.373436
185    [9,]  1.656324 1.219490
186   [10,]  1.959487 1.417368
187   [11,]  2.217157 1.640443
188
```

189   Decision Trees:

```
190   # First, we're going to create a decision tree to find a good guess as how
191   # accurate we can get it.
192   first_tree = rpart(V9 ~ ., data = train)
193   first_predicted = predict(first_tree, newdata = test)
194   RMSE_MAE(first_predicted, test[, 9])
```

```
195
196            [,1]      [,2]
197   [1,]  2.452406 1.754041
198
```

```
199   # Now that we have an estimate for accuracy, we're going to use randomForest
200   # to get a more accurate model using the same test and train data as the
201   # decision tree did. Note: like all tune functions, it does include a
202   # 10-fold cross validation test of accuracy.
203   forestmodel = tune.randomForest(V9 ~ ., data = train)
204   forestmodel
```

```
205
206
207   Error estimation of 'randomForest' using 10-fold cross validation: 4.709902
208
```

```
209   randomforest_cv = function(d, m, data_frame) {
210       error_matrix = matrix(nrow = m + 1, ncol = 2)
211       for (i in 1:m) {
212           d_records = sample(nrow(data_frame), d)
213           temp_train = data_frame[d_records, ]
214           temp_test = data_frame[-d_records, ]
215           model = randomForest(V9 ~ ., data = temp_train)
216           predicted = predict(model, temp_test)
217           error_matrix[i, ] = RMSE_MAE(predicted, temp_test[, 9])
218       }
219       error_matrix[m + 1, 1] = mean(error_matrix[1:m, 1])
220       error_matrix[m + 1, 2] = mean(error_matrix[1:m, 2])
221       return(error_matrix)
222   }
223   randomforest_cv(round(0.7 * nrow(data)), 25, data)
```

```
224
225          [,1]      [,2]
226   [1,]  2.189761  1.559031
227   [2,]  2.032656  1.451119
228   [3,]  2.147769  1.502647
229   [4,]  2.160276  1.517783
230   [5,]  2.155134  1.504519
231   [6,]  2.237676  1.545699
232   [7,]  2.156094  1.521302
233   [8,]  2.065127  1.496686
234   [9,]  2.080927  1.502607
235  [10,]  2.097192  1.522242
236  [11,]  2.146814  1.506259
237  [12,]  2.029483  1.470464
238  [13,]  2.089718  1.517656
239  [14,]  2.149859  1.547441
240  [15,]  2.225827  1.560536
241  [16,]  2.151263  1.543244
242  [17,]  2.139978  1.500482
243  [18,]  2.183016  1.509353
244  [19,]  2.049982  1.459659
245  [20,]  2.161172  1.491655
246  [21,]  2.316180  1.586989
247  [22,]  2.089543  1.470437
248  [23,]  1.925467  1.424065
249  [24,]  2.140897  1.521362
250  [25,]  2.250510  1.603325
251  [26,]  2.134893  1.513462
252
```

Support Vector Machines:

```r
# Here I'm making a for loop that runs the tune function five times to find
# a good guess at the best parameters.
best_parameters = matrix(ncol = 2, nrow = 5)
obj <- tune.svm(V9 ~ ., data = train, gamma = 2^(0:4), cost = 2^(0:4))
best_parameters[1, ] = as.matrix(obj$best.parameters)
for (i in 1:4) {
    gamma = as.numeric(obj$best.parameters)[1]
    cost = as.numeric(obj$best.parameters)[2]
    obj <- tune.svm(V9 ~ ., data = train, gamma = (0.25 * gamma):(1.75 * gamma),
        cost = (0.5 * cost):(1.5 * cost))
    best_parameters[i + 1, ] = as.matrix(obj$best.parameters)
}
best_model = obj$best.model
svm_cv = function(d, m, data_frame) {
    error_matrix = matrix(nrow = m + 1, ncol = 2)
    for (i in 1:m) {
        d_records = sample(nrow(data_frame), d)
        temp_train = data_frame[d_records, ]
        temp_test = data_frame[-d_records, ]
        model = svm(V9 ~ ., data = temp_train, gamma = best_parameters[5, 1],
            cost = best_parameters[5, 2])
        predicted = as.numeric(predict(model, temp_test))
        error_matrix[i, ] = RMSE_MAE(predicted, temp_test[, 9])
    }
    error_matrix[m + 1, 1] = mean(error_matrix[1:m, 1])
    error_matrix[m + 1, 2] = mean(error_matrix[1:m, 2])
    return(error_matrix)
}
svm_cv(round(0.7 * nrow(data)), 25, data)
```

```
           [,1]      [,2]
 [1,]  2.290255  1.554808
 [2,]  2.192552  1.505933
 [3,]  2.186716  1.510708
 [4,]  2.179284  1.517026
 [5,]  2.182027  1.498752
 [6,]  2.219791  1.533720
 [7,]  2.357776  1.592112
 [8,]  2.323458  1.560145
 [9,]  2.254139  1.545554
[10,]  2.300648  1.561159
[11,]  2.265284  1.539249
[12,]  2.299996  1.529469
[13,]  2.180669  1.470896
[14,]  2.240785  1.510838
[15,]  2.263901  1.561225
[16,]  2.268788  1.529638
[17,]  2.271419  1.541745
[18,]  2.218457  1.482707
[19,]  2.240139  1.545508
[20,]  2.333338  1.609474
[21,]  2.234419  1.524602
[22,]  2.338490  1.618267
[23,]  2.205766  1.516662
[24,]  2.337352  1.556476
[25,]  2.135988  1.479137
[26,]  2.252857  1.535832
```

```
obj
```

```
Parameter tuning of 'svm':

- sampling method: 10-fold cross validation

- best parameters:
      gamma     cost
 0.00390625  1.9375

- best performance: 5.145983
```

Naive Bayes:

```
nrings = data[, 9]
rings = as.factor(nrings)
sex = data[, 1]
nb_data = data[, -1]
nb_train = train[, -1]
nb_test = test[, -1]
Error_matrix = matrix(nrow = 14, ncol = 2)
# Since the Shapiro-Wilks test suggests that all of the numeric variables
# are nonnormal, I have to find out the best way to turn them into discrete
# factors using a for loop. Thus, I am keeping the train and test set
# constant and only changing number of levels for the newly-factorized
# variables.
colnames(Error_matrix) = c("RMSE", "MAE")
for (i in 2:15) {
    discretized_train = discretizer(nb_train[, -8], i)
    discretized_train = cbind(sex[train_rows], discretized_train)
```

```
342        discretized_test = discretizer(nb_test[, -8], i)
343        discretized_test = cbind(sex[-train_rows], discretized_test)
344        model = naiveBayes(rings[train_rows] ~ ., data = discretized_train)
345        predicted_naive_Bayes = predict(model, test)
346        predicted_naive_Bayes = as.numeric(predicted_naive_Bayes)
347        Error_matrix[i - 1, ] = RMSE_MAE(predicted_naive_Bayes, test[, 9])
348 }
349 Error_matrix[argmin(Error_matrix[, 1]), ]
```

```
350
351      RMSE        MAE
352  3.547534 2.709497
353
```

```
354 Error_matrix[argmin(Error_matrix[, 2]), ]
```

```
355
356      RMSE        MAE
357  3.662323 2.604948
358
```

```
359 # We can see that the factor levels that minimizes the Root Mean Square
360 # error and the Mean Absolute error are very close but different and both of
361 # them don't change all that much when switching between the two. So, we
362 # could confidently pick either one and would recieve similar results. I am
363 # now going to use 10-fold cross-validation to double check my numbers.
364 data_RMSE = discretizer(nb_data[, -8], argmin(Error_matrix[, 1]))
365 data_RMSE = cbind(sex, data_RMSE)
366 data_MAE = discretizer(nb_data[, -8], argmin(Error_matrix[, 2]))
367 data_MAE = cbind(sex, data_MAE)
368 data = cbind(sex, data)
369 naiveBayes_cv = function(d, m, data_frame) {
370     error_matrix = matrix(nrow = m + 1, ncol = 2)
371     for (i in 1:m) {
372         d_records = sample(nrow(data_frame), d)
373         temp_test = data_frame[d_records, ]
374         temp_train = data_frame[-d_records, ]
375         model = naiveBayes(rings[-d_records] ~ ., data = temp_train)
376         predicted = as.numeric(predict(model, temp_test))
377         error_matrix[i, ] = RMSE_MAE(predicted, nrings[d_records])
378     }
379     error_matrix[m + 1, 1] = mean(error_matrix[1:m, 1])
380     error_matrix[m + 1, 2] = mean(error_matrix[1:m, 2])
381     return(error_matrix)
382 }
383 naiveBayes_cv(round(0.7 * nrow(data_RMSE)), 25, data_RMSE)
```

```
384
385              [,1]       [,2]
386   [1,] 3.060773 2.007182
387   [2,] 3.401939 2.136115
388   [3,] 3.450897 2.155609
389   [4,] 3.254482 2.045828
390   [5,] 3.449410 2.161081
391   [6,] 3.033838 1.975718
392   [7,] 3.504297 2.223324
393   [8,] 3.248224 2.117305
394   [9,] 3.370275 2.136457
395  [10,] 3.421436 2.158345
396  [11,] 2.901837 1.914501
397  [12,] 3.318583 2.152531
398  [13,] 3.087362 2.022914
399  [14,] 2.921160 1.927839
400  [15,] 3.483693 2.117647
```

```
401  [16,] 3.427777 2.121067
402  [17,] 3.236569 2.064295
403  [18,] 2.952253 1.929207
404  [19,] 3.490411 2.189124
405  [20,] 3.168868 2.058824
406  [21,] 3.285336 2.110807
407  [22,] 3.155403 2.049590
408  [23,] 2.829696 1.879959
409  [24,] 3.365401 2.092681
410  [25,] 3.086531 2.023940
411  [26,] 3.236258 2.070876
412
```

```
413  naiveBayes_cv(round(0.7 * nrow(data_MAE)), 25, data_MAE)
```

```
414
415            [,1]      [,2]
416   [1,] 3.011322 1.985978
417   [2,] 3.144491 2.037620
418   [3,] 3.029495 2.009576
419   [4,] 3.050924 1.981190
420   [5,] 3.057866 2.020862
421   [6,] 3.150847 2.102257
422   [7,] 3.083594 2.044118
423   [8,] 3.015464 1.987004
424   [9,] 3.121732 2.037278
425  [10,] 3.077377 2.019494
426  [11,] 3.058985 1.993502
427  [12,] 3.053669 2.015048
428  [13,] 3.073818 2.011286
429  [14,] 3.190701 2.067031
430  [15,] 3.012514 1.978796
431  [16,] 3.085866 2.027360
432  [17,] 3.190433 2.049590
433  [18,] 3.079376 2.021546
434  [19,] 3.086254 2.037278
435  [20,] 3.029100 2.003762
436  [21,] 3.043461 2.017100
437  [22,] 2.996521 1.964090
438  [23,] 3.256793 2.127223
439  [24,] 3.266492 2.122777
440  [25,] 3.044416 2.005130
441  [26,] 3.088461 2.026676
442
```

```
443  nb_n_fold_cross = function(n, data_frame) {
444      chopping_point = floor(nrow(data_frame)/n)
445      stopping_points = c(1/chopping_point, 1:n) * chopping_point
446      error_matrix = matrix(ncol = 2, nrow = n + 1)
447      for (i in 1:n) {
448          lower_bound = stopping_points[i]
449          upper_bound = stopping_points[i + 1]
450          temp_test = data_frame[lower_bound:upper_bound, ]
451          temp_train = data_frame[-(lower_bound:upper_bound), ]
452          model = naiveBayes(rings[-(lower_bound:upper_bound)] ~ ., data = temp_train)
453          predicted = as.numeric(predict(model, temp_test))
454          error_matrix[i, ] = RMSE_MAE(predicted, nrings[lower_bound:upper_bound])
455      }
456      error_matrix[n + 1, 1] = mean(error_matrix[1:n, 1])
457      error_matrix[n + 1, 2] = mean(error_matrix[1:n, 2])
458      return(error_matrix)
459  }
```

```
460  nb_n_fold_cross(10, data_RMSE)
```

```
461
462            [,1]      [,2]
463   [1,]  4.236120  3.074519
464   [2,]  5.044064  3.930622
465   [3,]  3.055311  1.602871
466   [4,]  2.486807  1.270335
467   [5,]  2.845293  1.507177
468   [6,]  4.653033  3.516746
469   [7,]  2.730739  1.337321
470   [8,]  3.635945  2.497608
471   [9,]  2.833920  1.653110
472  [10,]  3.049041  1.827751
473  [11,]  3.457027  2.221806
474
```

```
475  nb_n_fold_cross(10, data_MAE)
```

```
476
477            [,1]      [,2]
478   [1,]  3.964688  2.843750
479   [2,]  5.021961  3.861244
480   [3,]  1.599192  1.093301
481   [4,]  1.537412  1.143541
482   [5,]  1.547494  1.078947
483   [6,]  4.469460  3.401914
484   [7,]  1.409978  1.069378
485   [8,]  3.138357  2.174641
486   [9,]  2.336009  1.476077
487  [10,]  2.664672  1.775120
488  [11,]  2.768922  1.991791
489
```

490  Knn function:

```
491  # Because knn can't calculate distance with factor variables, we have to
492  # create two dummy variables. Also, because knn relies on distance, I scaled
493  # the data so that variables who have smaller ranges aren't more important.
494  male = as.numeric(data[, 1] == "M")
495  female = as.numeric(data[, 1] == "F")
496  knn_data = cbind(male, female, data[, -1])
497  knn_train = cbind(male[train_rows], female[train_rows], train[, -1])
498  knn_test = cbind(male[-train_rows], female[-train_rows], test[, -1])
499  Error_matrix = matrix(ncol = 2, nrow = 50)
500  for (i in 1:50) {
501      predicted_knn = as.numeric(knn(knn_train[, -10], knn_test[, -10], knn_train[,
502          10], k = i))
503      Error_matrix[i, ] = cbind(RMSE_MAE(knn_test[, 10], predicted_knn))
504  }
505  Error_matrix[argmin(Error_matrix[, 1]), ]
```

```
506
507  [1] 3.519414 2.777334
508
```

```
509  Error_matrix[argmin(Error_matrix[, 2]), ]
```

```
510
511  [1] 3.585686 2.774142
512
```

```
513  k_RMSE = argmin(Error_matrix[, 1])
514  k_MAE = argmin(Error_matrix[, 2])
515  # Delete-d cross validation
516  knn_cv = function(d, m, data_frame, k) {
517      error_matrix = matrix(nrow = m + 1, ncol = 2)
```

```
518        for (i in 1:m) {
519            d_records = sample(nrow(data_frame), d)
520            temp_test = data_frame[-d_records, ]
521            temp_train = data_frame[d_records, ]
522            predicted_knn = as.numeric(knn(temp_train[, -10], temp_test[, -10],
523                cl = temp_train[, 10], k = k))
524            error_matrix[i, ] = RMSE_MAE(temp_test[, 10], predicted_knn)
525        }
526        error_matrix[m + 1, 1] = mean(error_matrix[1:m, 1])
527        error_matrix[m + 1, 2] = mean(error_matrix[1:m, 2])
528        return(error_matrix)
529 }
530 knn_cv(round(0.7 * nrow(knn_data)), 25, knn_data, k = k_RMSE)
```

```
531
532 Warning in knn(temp_train[, -10], temp_test[, -10], cl = temp_train[, 10], : NAs introduced by coercion
533
```

```
534
535 Warning in knn(temp_train[, -10], temp_test[, -10], cl = temp_train[, 10], : NAs introduced by coercion
536
```

```
537
538 Error in knn(temp_train[, -10], temp_test[, -10], cl = temp_train[, 10], : NA/NaN/Inf in foreign function
539    call (arg 6)
540
```

```
541 knn_cv(round(0.7 * nrow(knn_data)), 25, knn_data, k = k_MAE)
```

```
542
543 Warning in knn(temp_train[, -10], temp_test[, -10], cl = temp_train[, 10], : NAs introduced by coercion
544
```

```
545
546 Warning in knn(temp_train[, -10], temp_test[, -10], cl = temp_train[, 10], : NAs introduced by coercion
547
```

```
548
549 Error in knn(temp_train[, -10], temp_test[, -10], cl = temp_train[, 10], : NA/NaN/Inf in foreign function
550    call (arg 6)
551
```

```
552 # 10-fold cross validation
553 knn_n_fold_cross = function(n, data_frame, k) {
554        chopping_point = floor(nrow(data_frame)/n)
555        stopping_points = c(1/chopping_point, 1:n) * chopping_point
556        error_matrix = matrix(ncol = 2, nrow = n + 1)
557        for (i in 1:n) {
558            lower_bound = stopping_points[i]
559            upper_bound = stopping_points[i + 1]
560            temp_test = data_frame[lower_bound:upper_bound, ]
561            temp_train = data_frame[-(lower_bound:upper_bound), ]
562            predicted = as.numeric(knn(temp_train[, -10], temp_test[, -10], cl = temp_train[,
563                10], k = k))
564            error_matrix[i, ] = RMSE_MAE(predicted, temp_test[, 10])
565        }
566        error_matrix[n + 1, 1] = mean(error_matrix[1:n, 1])
567        error_matrix[n + 1, 2] = mean(error_matrix[1:n, 2])
568        return(error_matrix)
569 }
570 knn_n_fold_cross(10, knn_data, k_RMSE)
```

```
571
572 Warning in knn(temp_train[, -10], temp_test[, -10], cl = temp_train[, 10], : NAs introduced by coercion
573
```

```
574
575 Warning in knn(temp_train[, -10], temp_test[, -10], cl = temp_train[, 10], : NAs introduced by coercion
576
```

```
577
578 Error in knn(temp_train[, -10], temp_test[, -10], cl = temp_train[, 10], : NA/NaN/Inf in foreign function
579    call (arg 6)
580
```

```
581 knn_n_fold_cross(10, knn_data, k_MAE)
```

```
Warning in knn(temp_train[, -10], temp_test[, -10], cl = temp_train[, 10], : NAs introduced by coercion
```

```
Warning in knn(temp_train[, -10], temp_test[, -10], cl = temp_train[, 10], : NAs introduced by coercion
```

```
Error in knn(temp_train[, -10], temp_test[, -10], cl = temp_train[, 10], : NA/NaN/Inf in foreign function
    call (arg 6)
```

Weighted Knn-function:

```r
k_vector = c(1:50)
# Since the optimal kernel does not always provide the best result, all of
# them are being tested with k values between 1 and 50. Once the optimal
# kernal and k value are calculated, they will cross validated.
optimal_test = train.kknn(V9 ~ ., data = knn_data, ks = k_vector, kernel = "optimal",
    ykernel = -1.2:1.2)
triangular_test = train.kknn(V9 ~ ., data = knn_data, ks = k_vector, kernel = "triangular",
    ykernel = -1.2:1.2)
gaussian_test = train.kknn(V9 ~ ., data = knn_data, ks = k_vector, kernel = "gaussian",
    ykernel = -1.2:1.2)
epanechnikov_test = train.kknn(V9 ~ ., data = knn_data, ks = k_vector, kernel = "epanechnikov",
    ykernel = -1.2:1.2)
biweight_test = train.kknn(V9 ~ ., data = knn_data, ks = k_vector, kernel = "biweight",
    ykernel = -1.2:1.2)
triweight_test = train.kknn(V9 ~ ., data = knn_data, ks = k_vector, kernel = "triweight",
    ykernel = -1.2:1.2)
cosine_test = train.kknn(V9 ~ ., data = knn_data, ks = k_vector, kernel = "cos",
    ykernel = -1.2:1.2)
inverted_test = train.kknn(V9 ~ ., data = knn_data, ks = k_vector, kernel = "inv",
    ykernel = -1.2:1.2)
rectangular_test = train.kknn(V9 ~ ., data = knn_data, ks = k_vector, kernel = "rectangular",
    ykernel = -1.2:1.2)
optimal_data_frame = rbind(c("optimal", min(optimal_test$MEAN.ABS), min(optimal_test$MEAN.SQU),
    optimal_test$best.parameters$k), c("triangular", min(triangular_test$MEAN.ABS),
    min(triangular_test$MEAN.SQU), triangular_test$best.parameters$k), c("epanechnikov",
    min(epanechnikov_test$MEAN.ABS), min(epanechnikov_test$MEAN.SQU), epanechnikov_test$best.parameters$k),
    c("biweight", min(biweight_test$MEAN.ABS), min(biweight_test$MEAN.SQU),
        biweight_test$best.parameters$k), c("triweight", min(triweight_test$MEAN.ABS),
        min(triweight_test$MEAN.SQU), triweight_test$best.parameters$k), c("cos",
        min(cosine_test$MEAN.ABS), min(cosine_test$MEAN.SQU), cosine_test$best.parameters$k),
    c("inv", min(inverted_test$MEAN.ABS), min(inverted_test$MEAN.SQU), inverted_test$best.parameters$k),
    c("gaussian", min(gaussian_test$MEAN.ABS), min(gaussian_test$MEAN.SQU),
        gaussian_test$best.parameters$k), c("rectangular", min(rectangular_test$MEAN.ABS),
        min(rectangular_test$MEAN.SQU), rectangular_test$best.parameters$k))
# Since train.kknn includes a leave-one-out cross validation technique, I
# will use that method.
optimal_data_frame[argmin(optimal_data_frame[, 2]), ]
```

```
[1] "gaussian"          "1.51924144935422" "4.8439507372699"
[4] "18"
```

```r
optimal_data_frame[argmin(optimal_data_frame[, 3]), ]
```

```
[1] "biweight"          "1.52023104377889" "4.82416270611321"
[4] "43"
```

```r
optimal_k_abs = optimal_data_frame[argmin(optimal_data_frame[, 2]), 4]
optimal_k_abs = as.numeric(optimal_k_abs)
optimal_kernel_abs = optimal_data_frame[argmin(optimal_data_frame[, 2]), 1]
```

```
642  optimal_kernel_abs = as.character(optimal_kernel_abs)
643  optimal_k_rmse = optimal_data_frame[argmin(optimal_data_frame[, 3]), 4]
644  optimal_k_rmse = as.numeric(optimal_k_rmse)
645  optimal_kernel_rmse = optimal_data_frame[argmin(optimal_data_frame[, 3]), 1]
646  optimal_kernel_rmse = as.character(optimal_kernel_rmse)
647  # 10 fold cross validation using convenience command
648  fold_RMSE = cv.kknn(V9 ~ ., data = knn_data, kcv = 10, k = optimal_k_rmse, kernel = optimal_kernel_rmse)
649  RMSE_MAE(fold_RMSE[[1]][1:4177], fold_RMSE[[1]][4178:(2 * 4177)])
```

```
650
651          [,1]      [,2]
652  [1,] 2.206253 1.526847
653
```

```
654  fold_MAE = cv.kknn(V9 ~ ., data = knn_data, kcv = 10, k = optimal_k_abs, kernel = optimal_kernel_abs)
655  RMSE_MAE(fold_MAE[[1]][1:4177], fold_MAE[[1]][4178:(2 * 4177)])
```

```
656
657          [,1]      [,2]
658  [1,] 2.205438 1.528683
659
```