# CS 199
# HW2 - Predicting Crime Rates

April 8, 2014

**Sam Laane** <laane2@illinois.edu>
**José Vicente Ruiz** <ruizcep2@illinois.edu>

# Contents

# 1   Removing variables with missing values

## 1.1   Implementation

```
 1  # Read the data.
 2  data <- read.csv('communitiesH.data', h=T) # Headers = True.
 3
 4  # Remove non-predictive attributes (including state number).
 5  clean_data <- data[, -c(1:5)]
 6
 7  # Remove the attributes with question marks (unknowns) of them.
 8  unkw <- clean_data == '?'
 9  unkw_per_attr <- apply(unkw, 2, sum) # Sum the columns (2).
10  attr_with_unkw <- which(unkw_per_attr > 0) # Equivalent to find in
       Matlab.
11  no_unkw_data <- clean_data[, -attr_with_unkw]
```

# 2   Basics - Linear regression

## 2.1   Implementation

```
 1  # Compute the regression with the data free of unknowns.
 2  linear_regr <- lm(ViolentCrimesPerPop ~ ., data=no_unkw_data)
 3
 4  # Evaluate the regression looking at the mean-squared error on the
       training data.
 5  residuals <- resid(linear_regr)  # Function to get the residuals.
 6  plot(predict(linear_regr, no_unkw_data), residuals)
 7
 8  mse_residuals <- sum((residuals - mean(residuals)) ^ 2) / length(
       residuals)
 9  printf(" - Mean-squared error on the whole data: %.2e", mse_residuals)
10
11  # Split off some test data and compute a new regression on the rest of
       the data
12  # (train data).
13  folds <- cvFolds(nrow(no_unkw_data), K=5)
14  train_data <- no_unkw_data[folds$subsets[folds$which != 1], ]
15  test_data  <- no_unkw_data[folds$subsets[folds$which == 1], ]
16
17  linear_regr_test <- lm(ViolentCrimesPerPop ~ ., data=train_data)
18
19  # Evaluate the regression looking at the mean-squared error on that
       test data.
20  residuals_test <- test_data$ViolentCrimesPerPop - predict(linear_regr_
       test, test_data)
21  plot(test_data$ViolentCrimesPerPop, residuals_test)
22
23  mse_residuals_test <- sum((residuals_test - mean(residuals_test)) ^ 2)
       / length(residuals_test)
24  printf(" - Mean-squared error on the test data (20%%): %.2e", mse_
       residuals_test)
25
```

```r
26  # Prepare the data for the Box-Cox tranformation substituting zero
        values by
27  # a very small number.
28  boxcox_data <- no_unkw_data
29  boxcox_data$ViolentCrimesPerPop[ which(boxcox_data$ViolentCrimesPerPop
        == 0) ] <- 1e-100
30
31  # Apply Box-Cox and get a list of the lamdbas and their log likelihood
        and obtain
32  # the value with the highest likelihood.
33  linear_regr_mod <- lm(ViolentCrimesPerPop ~ ., data=boxcox_data)
34  lambdas <- boxcox(linear_regr_mod, plotit=F)
35  max_lambda <- lambdas$x[ which(lambdas$y == max(lambdas$y)) ]
36
37  # Change the very low values to zero again.
38  boxcox_data$ViolentCrimesPerPop[which(boxcox_data$ViolentCrimesPerPop
        == 1e-100)] <- 0
39
40  # Transform the data with the given value of lambda.
41  if (max_lambda == 0) {
42    boxcox_data$ViolentCrimesPerPop <- log(boxcox_data$
        ViolentCrimesPerPop)
43  } else {
44    boxcox_data$ViolentCrimesPerPop <- (boxcox_data$ViolentCrimesPerPop^
        max_lambda - 1)/max_lambda
45  }
46
47  # Compute the linear regression again and plot it.
48  boxcox_regr <- lm(ViolentCrimesPerPop ~ ., data=boxcox_data)
49  plot(boxcox_regr)
```

## 2.2  Results

## 2.3  Conclusions

# 3  Basics - Nearest Neighbour regression

## 3.1  Implementation

```r
 1  # Compute the regression and plot it.
 2  nn_regr <- knn.reg(train_data[!(names(train_data)) %in% c("
        ViolentCrimesPerPop")],
 3                     test = NULL,
 4                     train_data$ViolentCrimesPerPop, k = 1,
 5                     algorithm = c("kd_tree"))
 6  plot(nn_regr$pred, nn_regr$residuals)
 7
 8  # Evaluate the regression on the above training data with mean-squared
        error.
 9  mse_residuals_knn <- sum((nn_regr$residuals - mean(nn_regr$residuals))
        ^ 2) /
10      length(nn_regr$residuals)
11  printf(" - Mean-squared error on the training data (80%%): %.2e", mse_
        residuals_knn)
```

```
12
13  # Compute the regression on the above test data.
14  klabels <- knn(train_data[!(names(train_data)) %in% c("
        ViolentCrimesPerPop")],
15                 test_data [!(names(test_data))  %in% c("
                     ViolentCrimesPerPop")],
16                 train_data$ViolentCrimesPerPop, k = 1,
17                 prob = FALSE, algorithm = c("kd_tree"))
18  k_test_residuals <- test_data$ViolentCrimesPerPop - as.numeric(levels(
        klabels))[klabels]
19  plot(as.numeric(levels(klabels))[klabels], k_test_residuals)
20
21  # Evaluate the previous regression with mean-squared error and print
        the result.
22  mse_residuals_knn_test <- sum((k_test_residuals - mean(k_test_residuals
        )) ^ 2) /
23      length(k_test_residuals)
24  printf(" - Mean-squared error on the test data (20%%): %.2e", mse_
        residuals_knn_test)
```

## 3.2  Results

## 3.3  Conclusions

# 4  Dealing with missing values

## 4.1  Implementation

```
 1  printf(" - Mean-squared error on the test data (20%%): %.2e", mse_
        residuals_knn_test)
 2
 3  # Impute unknown values by covering all the columns with unknowns.
 4  interpolated_data <- clean_data
 5  for(i in attr_with_unkw) {
 6      # Divide data by rows in two sets: one with the samples that
            contain a question
 7      # mark in that column and one with the rest.
 8      unk_indices <- which(clean_data[,i] == '?')
 9      i_unk_data <- no_unkw_data[unk_indices, ]
10      i_no_unk_data <- no_unkw_data[-unk_indices, ]
11
12      # Compute the nearest neighbours in the set of elements with
            question mark of
13      # the elements with a question mark in that column.
14      i_klabels <- knn(i_no_unk_data,
15                      i_unk_data,
16                      i_no_unk_data$ViolentCrimesPerPop, k = 1,
17                      prob = FALSE, algorithm=c("kd_tree"))
18      indices <- attr(i_klabels, "nn.index")
19      nn_subs <- clean_data[-unk_indices, ][indices,]
20
21      # Substitute the question mark values by the value of the nearest
            neighbour.
```

```
22        interpolated_data[unk_indices, i] <- nn_subs[,i]
23  }
24
25  # Split off new test data.
26  folds <- cvFolds(nrow(interpolated_data), K=5)
27  train_data <- interpolated_data[folds$subsets[folds$which != 1], ]
28  test_data  <- interpolated_data[folds$subsets[folds$which == 1], ]
29
30  # Compute a linear regression again and evaluate it with mean-squares.
31  printf("Linear Regression (imputed missing values):")
32  linear_regr_test <- lm(ViolentCrimesPerPop ~ ., data=train_data)
33
34  residuals_test <- test_data$ViolentCrimesPerPop - predict(linear_regr_
        test, test_data)
35  plot(test_data$ViolentCrimesPerPop, residuals_test)
36
37  mse_residuals_test <- sum((residuals_test - mean(residuals_test)) ^ 2)
        / length(residuals_test)
38  printf(" - Mean-squared error on the test data (20%%): %.2e", mse_
        residuals_test)
39
40  # Also, compute a nearest neighbour regression and evaluate it.
41  printf("Nearest Neighbours (imputed missing values):")
42  klabels <- knn(train_data[!(names(train_data)) %in% c("
        ViolentCrimesPerPop")],
43                 test_data [!(names(test_data))  %in% c("
                    ViolentCrimesPerPop")],
44                 train_data$ViolentCrimesPerPop, k = 1,
45                 prob = FALSE, algorithm = c("kd_tree"))
46  k_test_residuals <- test_data$ViolentCrimesPerPop - as.numeric(levels(
        klabels))[klabels]
47  plot(as.numeric(levels(klabels))[klabels], k_test_residuals)
48
49  # Compute the mean-squared error and print the result.
50  mse_residuals_knn_test <- sum((k_test_residuals - mean(k_test_residuals
        )) ^ 2) /
51      length(k_test_residuals)
52  printf(" - Mean-squared error on the test data (20%%): %.2e", mse_
        residuals_knn_test)
```

## 4.2  Results

## 4.3  Conclusions

# 5   Modified Nearest Neighbours

## 5.1   Implementation

## 5.2   Results

## 5.3   Conclusions