

DATA SCIENCE

GROUP 6

# UNIVERSITY OF SOUTHERN DENMARK



## Report

Niels Alexander Hvid  
nhvid13@student.sdu.dk

Matic Novak  
maticnova@gmail.com

Lucas Østergaard Jensen  
lujen14@student.sdu.dk

Simon Hjortshøj Larsen  
sila114@student.sdu.dk

Mads Stagelund Mogensen  
mmoge14@student.sdu.dk

# 1 Alcohol related car accidents

At the beginning, we prepared the received data to use it in the language *R*. We exported the .xlsx file to a .csv file. All categorical values were replaced with numbers and used as a factor:

- Gender: Male = 1; Female = 0
- Accident: Yes = 1; No = 0
- Socioeconomic\_status: Lower = 0; Middle = 1; Upper = 2

We also removed the column *Subject* because it does not bring any value to the data - it is not a meaningful attribute.

## Question 1

The three other variables given apart from BAC are age, gender and socioeconomic status. The likelihood that a certain group has in getting into an accident could be seen bellow.

$$\begin{aligned} \text{Age } <= 40 : \frac{22}{125} = 0.18 \\ \text{Age } > 40 : \frac{35}{71} = 0.49 \\ \text{Females} : \frac{11}{85} = 0.13 \\ \text{Males} : \frac{46}{111} = 0.41 \end{aligned} \tag{1}$$

$$\text{Lower socioeconomic status} : \frac{19}{70} = 0.27$$

$$\text{Middle socioeconomic status} : \frac{14}{49} = 0.29$$

$$\text{Upper socioeconomic status} : \frac{24}{77} = 0.31$$

This shows that males are more than 3 times as likely to be in an accident compared to females, and people above 40 are more than twice as likely to be in an accident compared to people below 40.

This suggests that for gender, male is a risk factor while female is a protective factor. With age, people at 41 or older is a risk factor while 40 or below is a protective factor.

## Question 2 and 4

After loading the data in R, we used the GLM function to get the logistics models. The crude model was given the BAC data, while the adjusted model was given all the data. We then extracted the betas and standard errors to calculate OR, higher, and lower bounds. The code for this can be seen in appendix 1.

The plots for crude and adjusted models on a linear scale can be seen on figure 1 and 2 respectively along with the adjusted model on a logarithmic scale on figure 3.

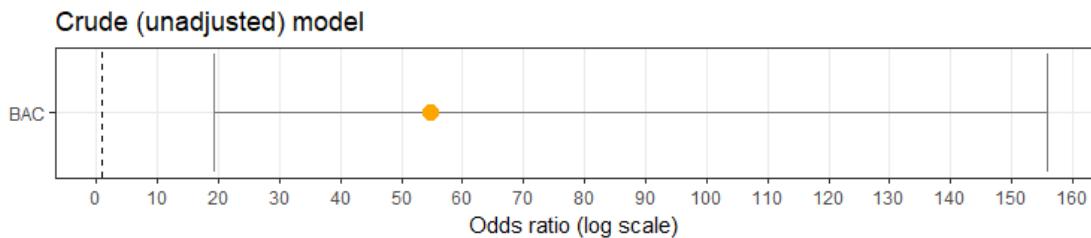


Figure 1: Crude model (linear scale)

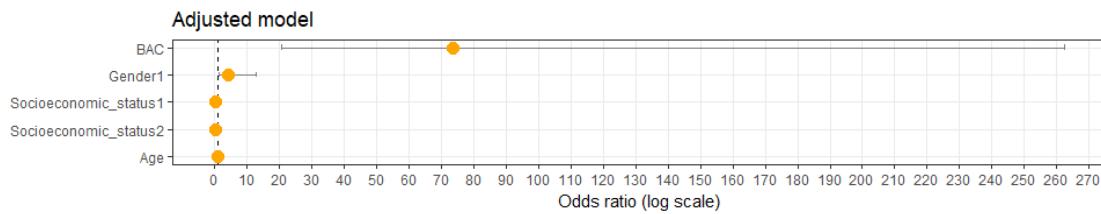


Figure 2: Adjusted model (linear scale)

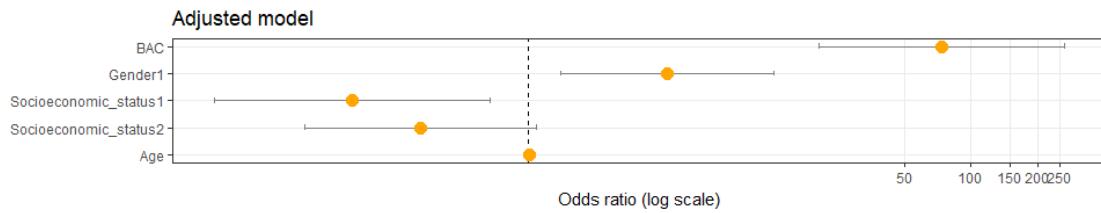


Figure 3: Adjusted model (logarithmic scale)

The adjusted model indicates that BAC has more significance than the crude model indicates. This is interesting as it is easy to make the assumption that other factors such as age and gender would decrease the significance of BAC. Males are also more likely than females to be in an accident. The socioeconomic status does not have a significant association with accidents.

The OR of the adjusted model for BAC is 73 while the unadjusted is 54.

In question 1 we got the impression that age could be a risk factor, but the model predicts that age has little impact on the probability of crashing. We made a scatterplot of Age versus BAC to see if this was really true.

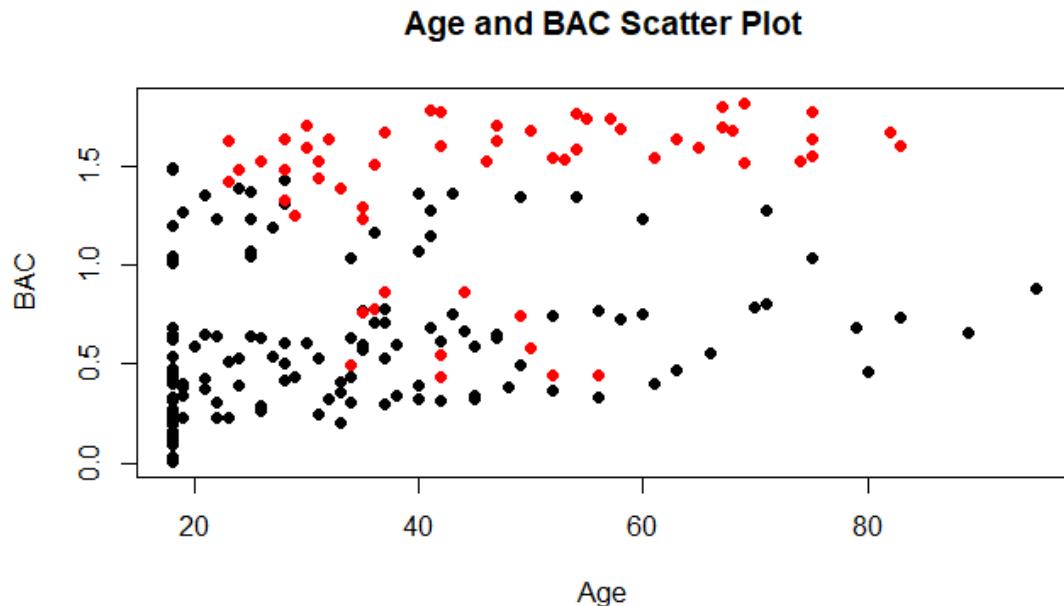


Figure 4: A plot of age versus BAC. Red dots are accidents.

It can be seen that higher age does not necessarily correspond to higher probability of crashing. We can also see that there are many samples of 18 year olds with low BAC that did not crash, which could explain our initial thoughts in question 1.

### Question 3

It's a recurring problem in Denmark that a high number of people choose to drive with a high BAC during Christmas. Especially when going home from Christmas dinners which are notoriously known for being celebrated with a high alcohol consumption. On top of the higher BAC, the Christmas dinners are held during the winter where the roads are more likely to be slippery hence more accidents could happen, not because of high BAC, but because of poor road conditions. Other holidays or festivity seasons could likewise influence the amount of accidents. One way to circumvent this is by including the season in which the drive took place.

Another possible confounding variable is drugs. Some drugs have a much more significant impact on cognitive abilities than alcohol and people who take drugs likely consumes alcohol at the same time. A false correlation between alcohol and traffic accidents could be found. This could be accounted for by including what types of drugs (if any) were taken when the drive took place.

The drugs could be one of the confounders because they fulfill these three conditions:

- People with higher BAC are more likely to take drugs than nondrinkers.
- Drugs are a risk factor for the accidents, even for the people with low BAC (or who do not drink).
- Drugs are not an intermediate factor (i.e. BAC or drinking does not cause taking drugs).

## Question 5

The model gave the following predictions for a man with a BAC of 1.0%:

- 40 years old: 59.5%
- 50 years old: 63.5%
- 60 years old: 67.3%
- 70 years old: 71%
- 80 years old: 74.3%

For each 10 year that passes the probability increases by roughly 3.8 percentage points. Although the change is almost linear, it decreases slowly. This makes sense as the probability would otherwise eventually exceed 100% which of course is not possible.

## Question 6

- Accuracy: 76.47%
- Sensitivity: 88.89%
- Specificity: 62.50%
- Precision: 72.73%

	Actual NO	Actual YES
Predicted NO	8	3
Predicted YES	1	5

Table 1: Confusion Matrix

The performance of the model is decent but not super good. This is both because of the small training set *and* the small validation set. Maybe it's also because of the aforementioned confounding variables.

## 2 LDA, k-NN, GLM

The *Boston* dataset is used for this part of the assignment. The median crime rate is calculated and applied as a class to each row. Rows with crime rate above the median are classified with class 1 and those below the median with class 0. The crime rate column is then removed from the dataset and all data except the class is scaled. A seed is chosen so the results are reproducible. The dataset is divided into training (75%) and testing (25%) sets. For training the rows with either class are equally represented.

### Question 1a

The Boston dataset comprises 14 variables. Each one of them has a short description as listed below:

crim	per capita crime rate by town.
zn	proportion of residential land zoned for lots over 25,000 sq.ft.
indus	proportion of non-retail business acres per town.
chas	Charles River dummy variable (= 1 if tract bounds river; 0 otherwise).
nox	nitrogen oxides concentration (parts per 10 million).
rm	average number of rooms per dwelling.
age	proportion of owner-occupied units built prior to 1940.
dis	weighted mean of distances to five Boston employment centres.
rad	index of accessibility to radial highways.
tax	full-value property-tax rate per \$10,000.
ptratio	pupil-teacherratio by town.
black	$1000(Bk - 0.63)^2$ where Bk is the proportion of blacks by town.
lstat	lower status of the population (percent).
medv	median value of owner-occupied homes in \$1000s.

Table 2: Boston dataset variable description as found 12/12/2017 on <https://www.cs.toronto.edu/~delve/data/boston/bostonDetail.html>

### Question 1b

Looking at the correlation matrix we find that certain variables have high correlation.

The correlation between *nox* and *indus* was 0.778. This makes sense, as many industries release more nitrogen oxide gasses than the average household.

The correlation between *indus* and *tax* was 0.774. This could indicate that the government in the area enforces higher property taxes on industrial properties.

As a result there is also a pretty high correlation between *nox* and *tax* (of 0.707). To make a subset, it would make sense to remove *nox* and *tax*, as *indus* is so highly

correlated with both of them.

## Question 1c

The following table shows the accuracy of the different models. Both the test set and validation set is shown as well as the accuracy for the test set of the subset. The test error rates can simply be calculated as  $1 - \text{accuracy}$ .

	Logistic Regression	LDA	k-NN
Test set	87.67%	85.95%	87.80%
Test subset	84.76%	83.12%	85.38%
Validation set	92.86%	86.51%	90.48%

Table 3: Accuracy for the different models

The Logistic regression performed the best closely followed by k-NN.

	Logistic Regression	LDA	k-NN
Sensitivity	92.96%	95.24%	85.71%
Specificity	93.65%	77.78%	90.48%
Precision	92.06%	93.55%	90.00%

Table 4: performance of the different models

The performances shown in table 4 were found using the validation set.

By performing calculations of odds ratios and boundaries of the attributes in logistic regression, we get the impression that especially *dis* and *rad* have high influence on the crime rate.

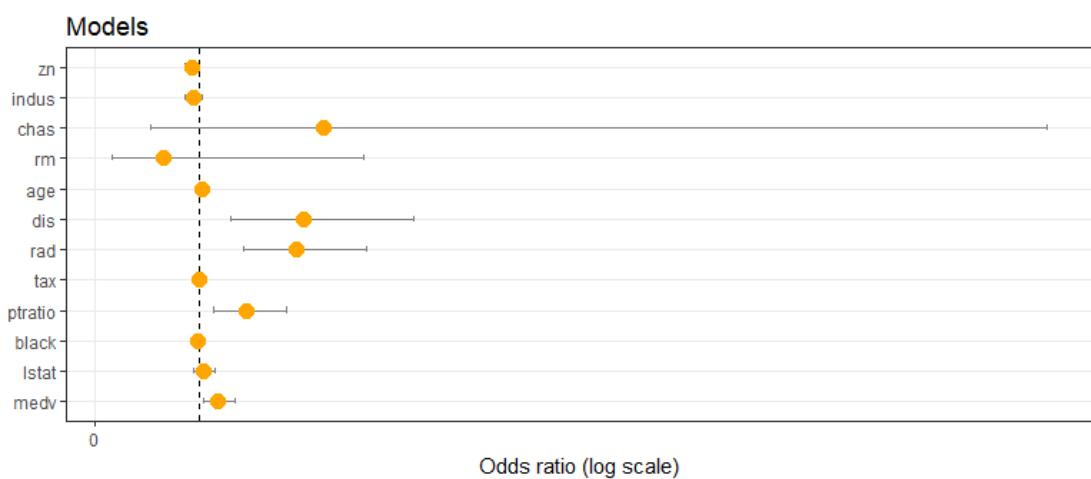


Figure 5: Odds ratios with boundaries

## Question 2a

$$\beta_0 = -6, \beta = \begin{bmatrix} 0.05 \\ 1 \end{bmatrix}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{array}{l} \text{time} \\ \text{grade point average} \end{array}$$

$x_1 = 40$  hours ;  $x_2 = 3.5$  grade => scoring A

$$\ln\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$\ln\left(\frac{p}{1-p}\right) = 6 + 0.05 \cdot 40 + 1 \cdot 3.5 = -0.5 \quad (2)$$

$$\frac{p}{1-p} = e^{-0.5}$$

$$p = (1-p) \cdot 0.6$$

$$p = 0.6 - 0.6p$$

$$p + 0.6p = 0.6$$

$$p = \frac{0.6}{1.6} = 0.377 = \underline{\underline{\underline{37.8\%}}}$$

Figure 6: The probability that a student working 40 hours with a C<sup>+</sup> (3.5) average grade scores an A (5) in the finals.

## Question 2b

$$\begin{aligned}
 p &= 80\% \\
 \ln\left(\frac{p}{1-p}\right) &= -6 + 0.05 \cdot H + 1 \cdot 3.5 \\
 \ln\left(\frac{0.8}{0.2}\right) &= -2.5 + 0.05H \quad (3) \\
 1.38 &= -2.5 + 0.05H \\
 3.88 &= 0.05H \\
 H &= \underline{\underline{77.73}} \text{ Hours per week}
 \end{aligned}$$

Figure 7: The number of hours the student must spend per week to have an 80% probability of scoring an A.

$$\begin{aligned}
 \text{Test:} \\
 \ln\left(\frac{p}{1-p}\right) &= -6 + 0.05 \cdot 77.73 + 3.5 \\
 \ln\left(\frac{p}{1-p}\right) &= 1.38 \\
 \frac{p}{1-p} &= e^{1.38} \\
 p &= 4(1-p) \\
 p &= 4 - 4p \\
 5p &= 4 \\
 \underline{\underline{p = 0.8}}
 \end{aligned} \tag{4}$$

Figure 8: Test to confirm the 80% probability of scoring an A.

### 3 Artist Identification

Paintings from the different artist from the *impressionism* period were downloaded from [www.wikiart.org](http://www.wikiart.org). The different artists drew different amount of paintings:

- Edouard Manet: 121
- Edgar Degas: 605
- Pierre-Auguste Renoir: 1377
- Claude Monet: 1324

To train a network to answer "Is this painting the work of  $x$ ?" four different networks has to be trained, one for each artist. The networks are all made using Matlab with these steps:

1. Load paintings of the artist to be detected.
2. Load same amount of paintings spread across the 3 other artists.
3. Load the pre-trained CNN *AlexNet*.
4. Resize all paintings to 227 by 227 to fit *AlexNet*.
5. Divide the dataset into a training set and a test set (75% and 25% respectively).
6. Extract features of training and test datasets from the layer before the classifier in *AlexNet*.
7. Train an SVM classifier using the extracted features.
8. Predict the labels of the test set to evaluate the performance of the classifier.

The amount of paintings vary from artist to artist. Manet only has 121 paintings. So the network designed to detect Manets paintings is trained on a dataset consisting of just 242 paintings (50% from Manet and 50% from others). The other networks can use datasets of size 726.

After the 4 networks finished training, we run the 4 paintings through all the networks resulting in the following response. We used Google reverse image search tool to find the real author of the paintings.

	Manet network	Degas network	Renior network	Monet network
Painting 1	<b>No</b>	<b>No</b>	<b>Yes</b>	<b>No</b>
Painting 2	<b>No</b>	<b>No</b>	<b>No</b>	<i>Yes</i>
Painting 3	<b>Yes</b>	<b>No</b>	<b>No</b>	<b>No</b>
Painting 4	<i>Yes</i>	<b>Yes</b>	<b>No</b>	<b>No</b>

Table 5: Network response, bold text is correct response

The accuracy of the networks could likely be increased by increasing the size of the datasets used to train the networks. Furthermore another pretrained network instead of *AlexNet* could be used which could result in a better performance (GoogleNet, ResNet50...).

	Manet network	Degas network	Renior network	Monet network
Accuracy	85%	86%	86%	90%

Table 6: Accuracy of the networks on the test set



Figure 9: (P1) Pierre Auguste Renoir



Figure 10: (P2) Jackson Pollock



Figure 11: (P3) Manet



Figure 12: (P4) Edgar Degas

## Appendix 1 - Car accidents

```
1 library(car)
2 library("caret")
3 library("ggplot2")
4
5 # Read and prepare the data.
6 filePath = "Data_Car_accidents_196.csv"
7 data <- read.csv(file=filePath, header=TRUE, sep=";", stringsAsFactors = TRUE)
8 data <- data[,2:6] # remove IDs
9 data$Gender <- as.factor(data$Gender)
10 data$Accident <- as.factor(data$Accident)
11 data$Socioeconomic_status <- as.factor(data$Socioeconomic_status)
12 summary(data)
13
14 plot(data$Age, data$BAC, main="Age and BAC Scatter Plot",
15       xlab="Age", ylab="BAC", col=data$Accident, pch=19)
16
17
18 ### METHODS
19 getBetas <- function(logistic_model) {
20   betas <- c()
21   for (i in 1:length(logistic_model$coefficients)) {
22     betas <- c(betas, logistic_model$coefficients[i])
23   }
24   return (betas)
25 }
26
27 getOddsRatios <- function(betas) {
28   oddsRatios <- c()
29   for (i in 1:length(betas)) {
30     oddsRatios <- c(oddsRatios, exp(betas[i]))
31   }
32   return (oddsRatios)
33 }
34
35 getStandardErrors <- function(logistic_model) {
36   standardErrors <- c()
37   for (i in 1:length(summary(logistic_model)$coefficients[, 2])) {
```

```

38     standardErrors <- c(standardErrors, summary(logistic_model)
39     $coefficients[, 2][i])
40   }
41   return (standardErrors)
42 }
43
44 getBounds <- function(oddsRatios, betas, standardErrors, confidence_number = 1.96)
45 {
46   bounds <- data.frame()
47   for (i in 1:length(betas)) {
48     low_bound <- exp(betas[i] - confidence_number * standardErrors[i])
49     high_bound <- exp(betas[i] + confidence_number * standardErrors[i])
50     entry <- data.frame(OR = oddsRatios[i], LB = low_bound, HB = high_bound, beta = betas[i])
51     bounds <- rbind(bounds, entry)
52   }
53   return (bounds)
54 }
55
56 plotBounds <- function(bounds, numf_of_columns_from_bounds = 3) {
57   zeros <- rep(0, numf_of_columns_from_bounds)
58   num_of_plots <- nrow(bounds)
59   attach(mtcars)
60   par(mfrow=c(num_of_plots, 1)) # number of rows and columns on the same graph
61
62   for (i in 1:nrow(bounds)) {
63     dots <- as.numeric(bounds[i,1:numf_of_columns_from_bounds])
64     plot(dots, zeros, main=rownames(bounds)[i])
65     #boxplot(dots)
66   }
67 }
68
69 forestPlot <- function(df, boxLabels, yAxis, modelTitle) {
70   #dev.off()
71   p <- ggplot(df, aes(x = boxOdds, y = yAxis))
72   p + geom_vline(aes(xintercept = 1), size = .25, linetype = "dashed") +
73     geom_errorbarh(aes(xmax = boxCIHigh, xmin = boxCILow), size = .5, height = .2,
74     color = "gray50") +
75     geom_point(size = 3.5, color = "orange") +
76     theme_bw() +

```

```

77   theme(panel.grid.minor = element_blank()) +
78   scale_y_continuous(breaks = yAxis, labels = boxLabels) +
79   scale_x_continuous(breaks = seq(0,270,10) ) +
80   #coord_trans(x = "log10") +
81   ylab("") +
82   xlab("Odds ratio (log scale)") +
83   #annotate(geom = "text", y =1.1, x = 3.5, label ="Model p < 0.001\nPseudo
84   R^2 = 0.10", size = 3.5, hjust = 0) +
85   ggtitle(modelTitle)
86 }
87 ####
88
89
90 #### 1) Identify demographic characteristics of the drivers that are risk
91 ####      (or protective) factors of car accidents.
92 onlyAccidents <- data[data$Acciden == 1, ]
93
94 gender <- data[data$Gender == 1, ]
95 genderAcc <- onlyAccidents[onlyAccidents$Gender == 1, ]
96 nrow(genderAcc)/nrow(gender)
97
98 genderF <- data[data$Gender == 0, ]
99 genderFAcc <- onlyAccidents[onlyAccidents$Gender == 0, ]
100 nrow(genderFAcc)/nrow(genderF)
101
102 ageLessThan41 <- data[data$Age < 41, ]
103 ageLessThan41Acc <- onlyAccidents[onlyAccidents$Age < 41, ]
104 nrow(ageLessThan41Acc)/nrow(ageLessThan41)
105
106 ageMoreThan40 <- data[data$Age > 40, ]
107 ageMoreThan40Acc <- onlyAccidents[onlyAccidents$Age > 40, ]
108 nrow(ageMoreThan40Acc)/nrow(ageMoreThan40)
109
110 status0 <- data[data$Socioeconomic_status == 0, ]
111 status0Acc <- status0[status0$Accident == 1, ]
112 nrow(status0Acc)/nrow(status0)
113
114 status1 <- data[data$Socioeconomic_status == 1, ]
115 status1Acc <- status1[status1$Accident == 1, ]

```

```

116 nrow(status1Acc)/nrow(status1)

117
118 status2 <- data[data$Socioeconomic_status == 2, ]
119 status2Acc <- status2[status2$Accident == 1, ]
120 nrow(status2Acc)/nrow(status2)

121
122
123 ### 2) Obtain the model relating BAC and car accidents.
124 ### a) Crude (not adjusted) OR's.
125 lmBac <- glm(Accident ~ BAC, family=binomial(logit), data=data)
126 summary(lmBac)
127 betasBac <- getBetas(lmBac)
128 oddsRatiosBac <- getOddsRatios(betasBac)
129 standardErrorsBac <- getStandardErrors(lmBac)
130 boundsBac <- getBounds(oddsRatiosBac, betasBac, standardErrorsBac)
131 plotBounds(boundsBac)

132
133 dev.off()
134 boxLabels = c("BAC")
135 yAxis = length(boxLabels):1
136 df <- data.frame(
137   yAxis = length(boxLabels):1,
138   boxOdds = c(boundsBac[2,1]),
139   boxCILow = c(boundsBac[2,2]),
140   boxCIHigh = c(boundsBac[2,3]))
141 )
142 forestPlot(df, boxLabels, yAxis, "Crude (unadjusted) model")

143
144
145 ### b) Adjusted OR's.
146 lmAdj <- glm(Accident ~ BAC + Gender + Socioeconomic_status + Age, family=binomial)
147 summary(lmAdj)
148 betasAdj <- getBetas(lmAdj)
149 oddsRatiosAdj <- getOddsRatios(betasAdj)
150 standardErrorsAdj <- getStandardErrors(lmAdj)
151 boundsAdj <- getBounds(oddsRatiosAdj, betasAdj, standardErrorsAdj)
152 plotBounds(boundsAdj)

153
154 dev.off()

```

```

155 boxLabels = c("BAC", "Gender1", "Socioeconomic_status1", "Socioeconomic_status2",
156 yAxis = length(boxLabels):1
157 df <- data.frame(
158   yAxis = length(boxLabels):1,
159   boxOdds = c(boundsAdj[2,1],boundsAdj[3,1],boundsAdj[4,1],boundsAdj[5,1],
160   boundsAdj[6,1]),
161   boxCILow = c(boundsAdj[2,2],boundsAdj[3,2],boundsAdj[4,2],boundsAdj[5,2],
162   boundsAdj[6,2]),
163   boxCIHigh = c(boundsAdj[2,3],boundsAdj[3,3],boundsAdj[4,3],boundsAdj[5,3],
164   boundsAdj[6,3]))
165 )
166 forestPlot(df, boxLabels, yAxis, "Adjusted model")
167
168
169 ### 5) What is the probability that a 40 yr male whose BAC is >1%, causes a car ac-
170 ###      What will be the probability, 10, 20, 30 and 40 years later? Is this change
171 linear?
172 probabilityOfX <- function(BAC, Gender, Age, betas) {
173   topFraction <- exp(betas[1] + betas[2] * BAC + betas[3] * Gender + betas[6] *
174     Age)
175   return (topFraction / (1 + topFraction))
176 }
177 p <- probabilityOfX(1, 1, 40, betasAdj)
178 # p = 0.595
179 # BetaAge = 1.017061034 --> exp(BetaAge) = 2.765056
180 # Every year increase in age, makes response being true more likely with
181 factor 2.765056.
182
183 p10 <- probabilityOfX(1, 1, 50, betasAdj)
184 p20 <- probabilityOfX(1, 1, 60, betasAdj)
185 p30 <- probabilityOfX(1, 1, 70, betasAdj)
186 p40 <- probabilityOfX(1, 1, 80, betasAdj)
187
188
189 ### 6) Evaluate the predictive performance of the model.
190 testFilePath = "Data_Car_accidents_17.csv"
191 testData <- read.csv(file=testFilePath, header=TRUE, sep=";", ,
192 stringsAsFactors = TRUE)
193 testLabels <- testData[,3]

```

```

194 testData <- testData[,c("Gender", "Age", "Socioeconomic_status", "BAC")]
195 testData$Gender <- as.factor(testData$Gender)
196 testData$Socioeconomic_status <- as.factor(testData$Socioeconomic_status)
197 summary(testData)
198 predicted <- predict(lmAdj, testData, type="response")
199
200 classFromProbabilities <- function(x) {
201   if (x > 0.5) return (1)
202   else return (0)
203 }
204
205 predictedClass <- sapply(predicted, classFromProbabilities)
206 performanceTable <- table(predictedClass, testLabels)
207 confussionMatrix <- confusionMatrix(performanceTable)
208
209 confussionMatrix
210 # precision = TP / (TP + FP)

```

## Appendix 2 - Crime Rate

```
1 library(ISLR)
2 library(caret)
3 library(car)
4 library(MASS)
5 set.seed(7)
6 data <- Boston
7
8 # Calculate median crime rate and set class to each row regarding
9 # if it is above or below median crime rate. Scale the orher data.
10 median_crim_rate <- median(data$crim)
11 classFromCrimRate <- function(x) {
12   if (x > median_crim_rate) return (1)
13   else return (0)
14 }
15 calculatedClass <- sapply(data$crim, classFromCrimRate)
16 # data <- data.frame(scale(data)) #Scaling is done before learning by method train
17 data$Class <- calculatedClass
18 data$Class <- as.factor(data$Class)
19 data <- data[,-1]
20 nrow(data[data$Class == 1,])
21
22
23 # Prepare training and test set.
24 TrainingDataIndex <- createDataPartition(data$Class, p=0.75, list=FALSE)
25 trainingData <- data[TrainingDataIndex,]
26 testData <- data[-TrainingDataIndex,]
27 trainingLabels <- trainingData$Class
28 testLabels <- testData$Class
29 testDataWithoutClass <- testData[,-14]
30 prop.table(table(trainingLabels))
31
32
33 ## Find the correlations between variables.
34 nearZeroVariables <- nearZeroVar(trainingData, saveMetrics = TRUE)
35 correlations <- cor(trainingData[,-14])
36 highCorrelations <- findCorrelation(correlations, cutoff = 0.75)
37
```

```

38 trainingDataCorr <- trainingData[,-drop(c(4,9))]
39 testDataCorr <- testData[,-drop(c(4,9))]
40 testDataCorrWithoutClass <- testDataCorr[,-12]
41
42
43 # Logistic Regression
44 lr_model_all <- train(Class~., data=trainingData,
45                         method='glm', family=binomial(link='logit'),
46                         preProcess=c('scale', 'center'))
47 lr_pred_all <- predict(lr_model_all, testData[,-14])
48 confusionMatrix(lr_pred_all, testData$Class)
49
50 lr_model_corr <- train(Class ~ rad + nox,
51                         data=trainingData,
52                         method='glm', family=binomial(link='logit'),
53                         preProcess=c('scale', 'center'))
54 lr_pred_corr <- predict(lr_model_corr, testData[,-14])
55 confusionMatrix(lr_pred_corr, testData$Class)
56 vif(lr_model_corr$finalModel)
57
58
59 # LDA
60 LDA_model_all <- train(Class~., data=trainingData,
61                         method='lda',
62                         preProcess=c('scale', 'center'))
63 LDA_pred_all <- predict(LDA_model_all, testData[,-14])
64 confusionMatrix(LDA_pred_all, testData$Class)
65
66 LDA_model_corr <- train(Class~., data=trainingDataCorr,
67                         method='lda',
68                         preProcess=c('scale', 'center'))
69 LDA_pred_corr <- predict(LDA_model_corr, testDataCorr[,-12])
70 confusionMatrix(LDA_pred_corr, testDataCorr$Class)
71
72
73 # k-NN
74 knnGrid <- expand.grid(.k=c(2))
75 knn_model_all <- train(x=trainingData[,-14], method='knn',
76                         y=trainingData$Class,

```

```

77         preProcess=c('center', 'scale'),
78         tuneGrid = knnGrid)
79 knn_model_corr <- train(x=trainingDataCorr[,-12], method='knn',
80                         y=trainingDataCorr$Class,
81                         preProcess=c('center', 'scale'),
82                         tuneGrid = knnGrid)
83 knn_pred <- predict(knn_model_all, testData[, -14])
84 confusionMatrix(knn_pred, testLabels)

85
86
87 # Plot the odds ratios and boundaries for logistic regression.
88 betasBac <- getBetas(summary(lr_model_all))
89 oddsRatiosBac <- getOddsRatios(betasBac)
90 standardErrorsBac <- getStandardErrors(lr_model_all)
91 boundsBac <- getBounds(oddsRatiosBac, betasBac, standardErrorsBac)
92 boundsBac<-boundsBac[!(boundsBac$OR > 10000),]

93
94 dev.off()
95 boxLabels = c("zn", "indus", "chas", "rm", "age", "dis", "rad", "tax", "ptratio",
96 yAxis = length(boxLabels):1
97 df <- data.frame(
98     yAxis = length(boxLabels):1,
99     boxOdds = c(boundsBac[2,1], boundsBac[3,1], boundsBac[4,1], boundsBac[5,1],
100    boundsBac[6,1], boundsBac[7,1], boundsBac[8,1], boundsBac[9,1],
101    boundsBac[10,1], boundsBac[11,1], boundsBac[12,1], boundsBac[13,1]),
102    boxCILow = c(boundsBac[2,2], boundsBac[3,2], boundsBac[4,2], boundsBac[5,2],
103    boundsBac[6,2], boundsBac[7,2], boundsBac[8,2], boundsBac[9,2],
104    boundsBac[10,2], boundsBac[11,2], boundsBac[12,2], boundsBac[13,2]),
105    boxCIHigh = c(boundsBac[2,3], boundsBac[3,3], boundsBac[4,3], boundsBac[5,3],
106    boundsBac[6,3], boundsBac[7,3], boundsBac[8,3], boundsBac[9,3],
107    boundsBac[10,3], boundsBac[11,3], boundsBac[12,3], boundsBac[13,3]))
108 )
109 forestPlot(df, boxLabels, yAxis, "Models")

```

## Appendix 3 - Artist Identification

```
1  %% Load Images
2  rootFolder = fullfile(pwd, 'images');
3  categories = {'degas', 'notDegas'};
4  % ImageDatastore operates on image file locations and help you manage the data.
5  images = imageDatastore(fullfile(rootFolder, categories),
6  'LabelSource', 'foldernames');
7  % Summarize the number of images per category.
8  numOfImagesPerCategory = countEachLabel(images)
9
10
11 %% Balance the number of images in the training set
12 % Determine the smallest amount of images in a category.
13 minNumOfImagesPerCategory = min(numOfImagesPerCategory(:,2));
14 % Use splitEachLabel method to trim the set.
15 images = splitEachLabel(images, minNumOfImagesPerCategory, 'randomize');
16 % Notice that each set now has exactly the same number of images.
17 countEachLabel(images)
18
19
20 %% Load Pre-trained AlexNet Network
21 net = alexnet()
22
23 %% Pre-process Images For CNN
24 % net can only process RGB images that are 227-by-227.
25 % imds.ReadFcn pre-process images on-the-fly.
26 % The imds.ReadFcn is called every time an image is read from the ImageDatastore.
27 images.ReadFcn = @(filename)readAndPreprocessImage(filename);
28
29
30 %% Prepare Training and Test Image Sets
31 % Pick 75% of images from each set for the training data and the remainder,
32 % 25%, for the validation data. Randomize the split to avoid biasing the results.
33 [train, test] = splitEachLabel(images, 0.75, 'randomize');
34
35
36 %% Extract Training Features Using CNN
37 % Extract features from one of the deeper layers using the activations method.
```

```

38 % Layer right before the classification layer.
39 % In net, this layer is named 'fc7'.
40 featureLayer = 'fc7';
41 trainingFeatures = activations(net, train, featureLayer, ...
42     'MiniBatchSize', 32, 'OutputAs', 'columns');
43
44
45 %% Train A Multiclass SVM Classifier Using CNN Features
46 % Use the CNN image features to train a multiclass SVM classifier.
47
48 % Get training labels from the trainingSet.
49 trainLabels = train.Labels;
50
51 % Train multiclass SVM classifier using a fast linear solver, and set
52 % 'ObservationsIn' to 'columns' to match the arrangement used for training
53 % features.
54 classifier = fitcecoc(trainingFeatures, trainLabels, ...
55     'Learners', 'Linear', 'Coding', 'onevsall', 'ObservationsIn', 'columns');
56
57
58 %% Evaluate Classifier
59 % Extract image features from test set using the CNN.
60 testFeatures = activations(net, test, featureLayer, 'MiniBatchSize', 32);
61
62 % Pass CNN image features to trained classifier.
63 [predictedLabels] = predict(classifier, testFeatures);
64
65 % Get the known labels.
66 testLabels = test.Labels;
67
68 % Tabulate the results using a confusion matrix.
69 confusionMatrix = confusionmat(testLabels, predictedLabels);
70
71 % Convert confusion matrix into percentage form.
72 confusionMatrix = bsxfun(@rdivide, confusionMatrix, sum(confusionMatrix, 2));
73
74 % Display the mean accuracy.
75 mean(diag(confusionMatrix))
76

```

```

77
78 %% Try the Newly Trained Classifier on "Validation" Images
79
80 % Load test images
81 testImagesFolder = fullfile(rootFolder, 'testImages');
82 d = dir([testImagesFolder, '\*.bmp']);
83 fileNames = char(d.name);
84
85 % Classify all files.
86 for idx = 1:length(d)
87     newImage = fullfile(testImagesFolder, fileNames(idx, :))
88
89     % Pre-process the images as required for the CNN.
90     img = readAndPreprocessImage(newImage);
91
92     % Extract image features using the CNN.
93     imageFeatures = activations(net, img, featureLayer);
94
95     % Make a prediction using the classifier.
96     label = predict(classifier, imageFeatures)
97 end
98
99
100 %% Helper Functions
101 function preparedImage = readAndPreprocessImage(file)
102     img = imread(file);
103
104     % Some images may be grayscale.
105     % Replicate the image 3 times to create an RGB image.
106     if ismatrix(img)
107         img = cat(3, img, img, img);
108     end
109
110     % Resize the image as required for the CNN.
111     preparedImage = imresize(img, [227 227]);
112 end

```