

Homework_7

Claire Kraft

2024-10-08

Question 10.1

Using the same crime data set `uscrime.txt` as in Questions 8.2 and 9.1, find the best model you can using (a) a regression tree model, and (b) a random forest model. In R, you can use the `tree` package or the `rpart` package, and the `randomForest` package. For each model, describe one or two qualitative takeaways you get from analyzing the results (i.e., don't just stop when you have a good model, but interpret it too).

In the regression tree it seems there are 4 "generations". The children are Po1, Pop, LF, NW. In the model validation section the R^2 result is 72.445%. R^2 is the metric used to measure the fitness of the model. The closer that the R^2 gets to 100 the better fit the model is. 72% is closer to 100% as the model fits nearly perfectly.

In the random forest i chose 5 predictors. The mean of squared residuals is 84022.61 and percent of variances explained is 42.61%. The model doesn't seem that strong as the variance explainability is sub 50%. But the R^2 is 90 which suggests the model fits well. Given the contradiction, poor variances explanation score (42.61) and good R^2 makes me think there could be an overfitting? Perhaps redoing the model may lead to a better result.

```
# Set up
rm(list = ls())

# Helper
#install.packages("caret")
library(tree)
library(caret)
```

```
## Loading required package: ggplot2
```

```
## Loading required package: lattice
```

```
library(randomForest)
```

```
## randomForest 4.7-1.2
```

```
## Type rfNews() to see new features/changes/bug fixes.
```

```
##
## Attaching package: 'randomForest'
```

```
## The following object is masked from 'package:ggplot2':  
##  
##   margin
```

```
#install.packages("devtools") # Look at reference 1 and 2  
library(devtools)
```

```
## Loading required package: usethis
```

```
#devtools::install_github('araastat/reprtree') # Look at reference 2  
library(reprtree)
```

```
## Loading required package: plotrix
```

```
## Registered S3 method overwritten by 'reprtree':  
##   method      from  
##   text.tree tree
```

```
#install.packages("dplyr")  
library(dplyr)
```

```
##  
## Attaching package: 'dplyr'
```

```
## The following object is masked from 'package:randomForest':  
##  
##   combine
```

```
## The following objects are masked from 'package:stats':  
##  
##   filter, lag
```

```
## The following objects are masked from 'package:base':  
##  
##   intersect, setdiff, setequal, union
```

```
#install.packages("magrittr")  
library(magrittr)  
#install.packages("pROC")  
library(pROC)
```

```
## Type 'citation("pROC")' for a citation.
```

```
##
## Attaching package: 'pROC'
```

```
## The following objects are masked from 'package:stats':
##
##      cov, smooth, var
```

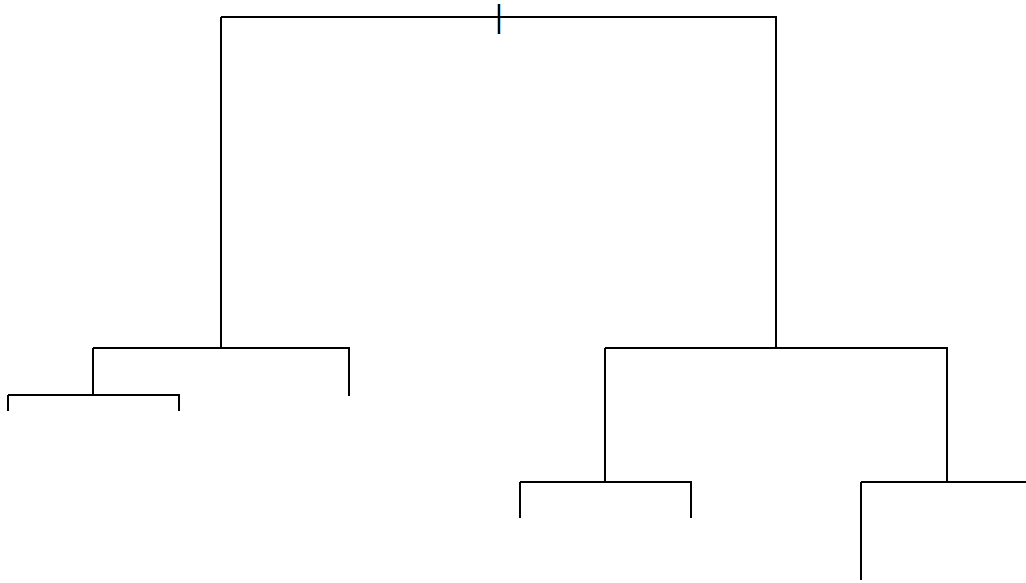
```
# Read in data
uscrime <- read.table("~/GitHub/omsa/ISYE 6501/Homework 07/uscrime.txt", stringsAsFactors = FALSE, header = TRUE)
head(uscrime)
```

```
##      M So  Ed Po1 Po2  LF  M.F Pop  NW  U1 U2 Wealth Ineq  Prob
## 1 15.1  1  9.1  5.8  5.6 0.510 95.0  33 30.1 0.108 4.1  3940 26.1 0.084602
## 2 14.3  0 11.3 10.3  9.5 0.583 101.2 13 10.2 0.096 3.6  5570 19.4 0.029599
## 3 14.2  1  8.9  4.5  4.4 0.533  96.9 18 21.9 0.094 3.3  3180 25.0 0.083401
## 4 13.6  0 12.1 14.9 14.1 0.577  99.4 157  8.0 0.102 3.9  6730 16.7 0.015801
## 5 14.1  0 12.1 10.9 10.1 0.591  98.5  18  3.0 0.091 2.0  5780 17.4 0.041399
## 6 12.1  0 11.0 11.8 11.5 0.547  96.4  25  4.4 0.084 2.9  6890 12.6 0.034201
##      Time Crime
## 1 26.2011    791
## 2 25.2999   1635
## 3 24.3006    578
## 4 29.9012   1969
## 5 21.2998   1234
## 6 20.9995    682
```

```
# Fit a reg tree
crime_tree <- tree(Crime~., data = uscrime)
summary(crime_tree)
```

```
##
## Regression tree:
## tree(formula = Crime ~ ., data = uscrime)
## Variables actually used in tree construction:
## [1] "Po1" "Pop" "LF" "NW"
## Number of terminal nodes:  7
## Residual mean deviance:  47390 = 1896000 / 40
## Distribution of residuals:
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -573.900 -98.300  -1.545   0.000  110.600  490.100
```

```
# Plot tree
plot(crime_tree)
text(crime_tree, pretty = 0)
```



crime_tree

```

## node), split, n, deviance, yval
##      * denotes terminal node
##
## 1) root 47 6881000  905.1
##    2) Po1 < 7.65 23  779200  669.6
##      4) Pop < 22.5 12  243800  550.5
##        8) LF < 0.5675 7  48520  466.9 *
##        9) LF > 0.5675 5  77760  667.6 *
##      5) Pop > 22.5 11  179500  799.5 *
##    3) Po1 > 7.65 24 3604000 1131.0
##      6) NW < 7.65 10  557600  886.9
##        12) Pop < 21.5 5  146400 1049.0 *
##        13) Pop > 21.5 5  147800  724.6 *
##      7) NW > 7.65 14 2027000 1305.0
##        14) Po1 < 9.65 6  170800 1041.0 *
##        15) Po1 > 9.65 8 1125000 1503.0 *

```

```

# Read in data
uscrime <- read.table("~/GitHub/omsa/ISYE 6501/Homework 07/uscrime.txt", stringsAsFactors = FALSE, header = TRUE)
head(uscrime)

```

```
##      M So   Ed Po1  Po2    LF   M.F Pop   NW   U1  U2 Wealth Ineq    Prob
## 1 15.1   1   9.1  5.8   5.6 0.510 95.0  33 30.1 0.108 4.1   3940 26.1 0.084602
## 2 14.3   0  11.3 10.3   9.5 0.583 101.2 13 10.2 0.096 3.6   5570 19.4 0.029599
## 3 14.2   1   8.9  4.5   4.4 0.533 96.9  18 21.9 0.094 3.3   3180 25.0 0.083401
## 4 13.6   0  12.1 14.9  14.1 0.577 99.4 157  8.0 0.102 3.9   6730 16.7 0.015801
## 5 14.1   0  12.1 10.9  10.1 0.591 98.5  18  3.0 0.091 2.0   5780 17.4 0.041399
## 6 12.1   0  11.0 11.8  11.5 0.547 96.4  25  4.4 0.084 2.9   6890 12.6 0.034201
##      Time Crime
## 1 26.2011    791
## 2 25.2999   1635
## 3 24.3006    578
## 4 29.9012   1969
## 5 21.2998   1234
## 6 20.9995    682
```

```
# Fit a reg tree
crime_tree <- tree(Crime~., data = uscrime)

# Make predictions on the training data
predicted_values <- predict(crime_tree, newdata = uscrime)
summary(predicted_values)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##    466.9   696.1   799.5   905.1 1049.2  1502.9
```

```
# Calculate Mean Squared Error (MSE)
mse <- mean((uscrime$Crime - predicted_values)^2)
mse
```

```
## [1] 40334.5
```

```
# Calculate R-squared
# total sum of squares
sst <- sum((uscrime$Crime - mean(uscrime$Crime))^2)
sst
```

```
## [1] 6880928
```

```
# sum of square errors
sse <- sum((uscrime$Crime - predicted_values)^2)
sse
```

```
## [1] 1895722
```

```
r2 <- 1 - (sse / sst) #1- 1895722/40334.5
r2
```

```
## [1] 0.7244962
```

```
# Set up
rm(list = ls())

# Read in data
uscrime <- read.table("~/GitHub/omsa/ISYE 6501/Homework 07/uscrime.txt", stringsAsFactors = FALSE, header = TRUE)
head(uscrime)
```

```
##      M So   Ed Po1  Po2   LF   M.F Pop   NW   U1  U2 Wealth Ineq   Prob
## 1 15.1   1   9.1  5.8   5.6 0.510  95.0  33 30.1 0.108 4.1   3940 26.1 0.084602
## 2 14.3   0  11.3 10.3   9.5 0.583 101.2  13 10.2 0.096 3.6   5570 19.4 0.029599
## 3 14.2   1   8.9  4.5   4.4 0.533  96.9  18 21.9 0.094 3.3   3180 25.0 0.083401
## 4 13.6   0  12.1 14.9  14.1 0.577  99.4 157   8.0 0.102 3.9   6730 16.7 0.015801
## 5 14.1   0  12.1 10.9  10.1 0.591  98.5  18   3.0 0.091 2.0   5780 17.4 0.041399
## 6 12.1   0  11.0 11.8  11.5 0.547  96.4  25   4.4 0.084 2.9   6890 12.6 0.034201
##      Time Crime
## 1 26.2011    791
## 2 25.2999   1635
## 3 24.3006    578
## 4 29.9012   1969
## 5 21.2998   1234
## 6 20.9995    682
```

```
# Set seed for reproducibility
set.seed(123)
num_pred <- 5

# Fit a random forest model
crime_forest <- randomForest(Crime~., data = uscrime, mtry = num_pred, importance = TRUE, ntree = 500)
crime_forest
```

```
##
## Call:
## randomForest(formula = Crime ~ ., data = uscrime, mtry = num_pred, importance = TRUE, ntree = 500)
##
##           Type of random forest: regression
##           Number of trees: 500
## No. of variables tried at each split: 5
##
##           Mean of squared residuals: 84022.61
##           % Var explained: 42.61
```

```
# Make predictions on the training data
crime_forest_pred <- predict(crime_forest, newdata = uscrime)
```

```
# Calculate R-squared
# Calculate Total Sum of Squares (SST)
sst <- sum((uscrime$Crime - mean(uscrime$Crime))^2)
sst
```

```
## [1] 6880928
```

```
# Calculate Sum of Squared Errors (SSE)
sse <- sum((crime_forest_pred - uscrime$Crime)^2)
sse
```

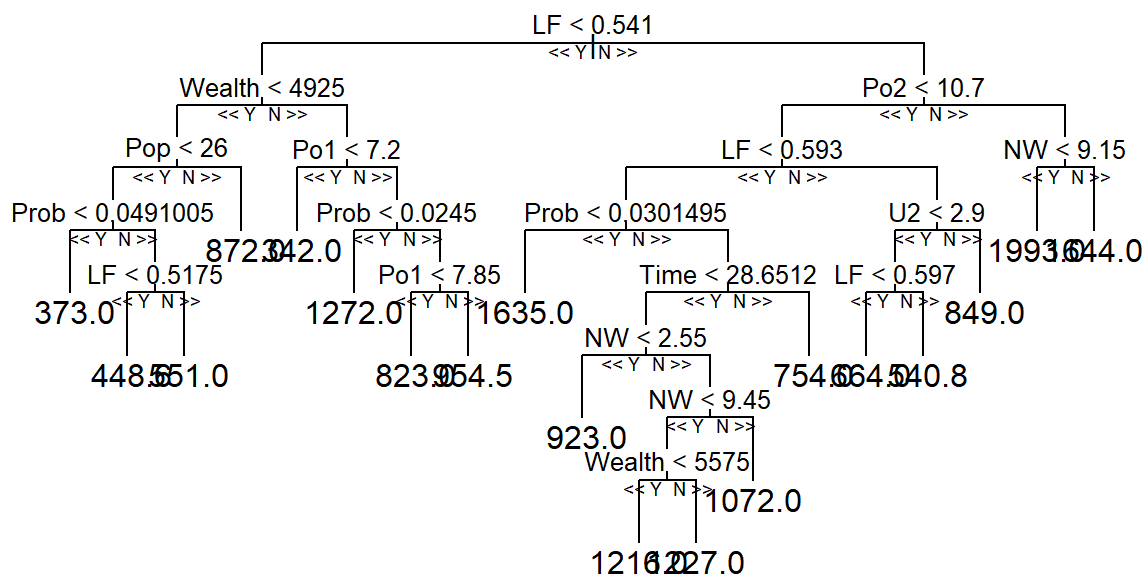
```
## [1] 684962
```

```
# Calculate R-squared
r2 <- 1 - (sse / sst)
r2
```

```
## [1] 0.900455
```

```
# Plot the random forest model
reptree::plot.getTree(crime_forest, k = 1, labelVar = TRUE) # Look at reference 1
```

```
## Warning in text.default(xy$x[leaves], xy$y[leaves] - 0.5 * charht, labels =
## stat, : "labelVar" is not a graphical parameter
```



Question 10.2

Describe a situation or problem from your job, everyday life, current events, etc., for which a logistic regression model would be appropriate. List some (up to 5) predictors that you might use.

Logistic regression is like linear regression except it is more useful for discrete variables. Discrete variables are like integers where the values are finite. Linear regression is good for continuous variables which are like floats. Some possible predictors could be - Work at home (1) or go in the office (0) - Go climbing (1) or stay at home (0) - Travel abroad (1) or stay-cation (0) - Cook at home (1) or take out (0)

Question 10.3.1

Using the GermanCredit data set `germancredit.txt` from <http://archive.ics.uci.edu/ml/machine-learning-databases/statlog/german> (<http://archive.ics.uci.edu/ml/machine-learning-databases/statlog/german>) / (description at <http://archive.ics.uci.edu/ml/datasets/Statlog+%28German+Credit+Data%29>) (<http://archive.ics.uci.edu/ml/datasets/Statlog+%28German+Credit+Data%29>), use logistic regression to find a good predictive model for whether credit applicants are good credit risks or not. Show your model (factors used and their coefficients), the software output, and the quality of fit. You can use the `glm` function in R. To get a logistic regression (logit) model on data where the response is either zero or one, use `family=binomial(link="logit")` in your `glm` function call.

Wow most variables are irrelevant. There is only one with the twinkles meaning they carry some level of statistical significance. So I will just extract variable `V14A143`. Looking at the ROC .8 looks to be the appropriate threshold.


```
# Set up
rm(list = ls())

# Read in data
ger_credit <- read.table("~/GitHub/omsa/ISYE 6501/Homework 07/german_credit.txt")

# Convert variables to binary
ger_credit$V21[ger_credit$V21==1]<-0
ger_credit$V21[ger_credit$V21==2]<-1
head(ger_credit)
```

```
##      V1 V2  V3  V4   V5  V6  V7 V8  V9  V10 V11  V12 V13  V14  V15 V16  V17 V18
## 1 A11  6 A34 A43 1169 A65 A75  4 A93 A101  4 A121  67 A143 A152  2 A173  1
## 2 A12 48 A32 A43 5951 A61 A73  2 A92 A101  2 A121  22 A143 A152  1 A173  1
## 3 A14 12 A34 A46 2096 A61 A74  2 A93 A101  3 A121  49 A143 A152  1 A172  2
## 4 A11 42 A32 A42 7882 A61 A74  2 A93 A103  4 A122  45 A143 A153  1 A173  2
## 5 A11 24 A33 A40 4870 A61 A73  3 A93 A101  4 A124  53 A143 A153  2 A173  2
## 6 A14 36 A32 A46 9055 A65 A73  2 A93 A101  4 A124  35 A143 A153  1 A172  2
##      V19  V20 V21
## 1 A192 A201  0
## 2 A191 A201  1
## 3 A191 A201  0
## 4 A191 A201  0
## 5 A191 A201  1
## 6 A192 A201  0
```

```
# Split data
ger_credit_train <- ger_credit[1:700,]
ger_credit_test <- ger_credit[701:1000,]
table(ger_credit_train$V21)
```

```
##
##      0      1
## 493 207
```

```
table(ger_credit_test$V21)
```

```
##
##      0      1
## 207  93
```

```
# Create logistical regression model
ger_credit_logreg = glm(V21~., data = ger_credit_train, family = binomial(link = "logit"))
summary(ger_credit_logreg)
```

```
##
## Call:
## glm(formula = V21 ~ ., family = binomial(link = "logit"), data = ger_credit_train)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  8.094e-01  1.318e+00   0.614 0.539040
## V1A12        -1.634e-01  2.627e-01  -0.622 0.533900
## V1A13        -1.039e+00  4.368e-01  -2.379 0.017345 *
## V1A14        -1.742e+00  2.908e-01  -5.990 2.09e-09 ***
## V2           2.888e-02  1.109e-02   2.605 0.009177 **
## V3A31         6.540e-01  6.927e-01   0.944 0.345051
## V3A32        -6.821e-01  5.092e-01  -1.339 0.180424
## V3A33        -9.090e-01  5.558e-01  -1.635 0.101968
## V3A34        -1.528e+00  5.305e-01  -2.880 0.003981 **
## V4A41        -1.873e+00  4.937e-01  -3.793 0.000149 ***
## V4A410       -1.556e+00  8.474e-01  -1.836 0.066314 .
## V4A42        -8.275e-01  3.200e-01  -2.586 0.009715 **
## V4A43        -8.720e-01  3.037e-01  -2.872 0.004085 **
## V4A44        -1.269e-01  9.259e-01  -0.137 0.890996
## V4A45        -5.639e-01  6.640e-01  -0.849 0.395753
## V4A46         4.279e-02  4.545e-01   0.094 0.924983
## V4A48        -2.364e+00  1.318e+00  -1.794 0.072847 .
## V4A49        -7.918e-01  4.184e-01  -1.892 0.058447 .
## V5           1.153e-04  5.572e-05   2.069 0.038552 *
## V6A62        -2.399e-01  3.460e-01  -0.693 0.488086
## V6A63        -4.461e-01  5.256e-01  -0.849 0.396045
## V6A64        -1.641e+00  6.364e-01  -2.579 0.009922 **
## V6A65        -8.197e-01  3.163e-01  -2.592 0.009548 **
## V7A72        -2.960e-01  5.418e-01  -0.546 0.584866
## V7A73        -4.289e-01  5.131e-01  -0.836 0.403237
## V7A74        -1.198e+00  5.565e-01  -2.153 0.031289 *
## V7A75        -5.071e-01  5.145e-01  -0.986 0.324340
## V8           3.565e-01  1.079e-01   3.304 0.000954 ***
## V9A92        -4.955e-01  4.594e-01  -1.079 0.280776
## V9A93        -1.299e+00  4.492e-01  -2.891 0.003843 **
## V9A94        -3.874e-01  5.379e-01  -0.720 0.471446
## V10A102       7.803e-01  4.842e-01   1.612 0.107066
## V10A103      -1.099e+00  5.261e-01  -2.090 0.036646 *
## V11          1.389e-02  1.048e-01   0.133 0.894503
## V12A122       4.053e-01  3.135e-01   1.293 0.196068
## V12A123       1.573e-01  2.837e-01   0.555 0.579128
## V12A124       8.694e-01  5.126e-01   1.696 0.089905 .
## V13          -1.376e-02  1.140e-02  -1.208 0.227097
## V14A142       6.347e-02  5.019e-01   0.126 0.899384
## V14A143      -6.359e-01  2.917e-01  -2.180 0.029250 *
## V15A152      -2.121e-01  2.939e-01  -0.722 0.470525
## V15A153      -9.052e-01  5.847e-01  -1.548 0.121597
## V16          3.427e-01  2.253e-01   1.521 0.128271
## V17A172      -7.079e-02  8.640e-01  -0.082 0.934695
## V17A173       1.246e-02  8.326e-01   0.015 0.988064
## V17A174       1.442e-01  8.212e-01   0.176 0.860619
```

```
## V18          4.924e-01  3.096e-01  1.591 0.111703
## V19A192      -3.059e-01  2.500e-01 -1.223 0.221152
## V20A202      -1.418e+00  8.186e-01 -1.732 0.083288 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 850.06  on 699  degrees of freedom
## Residual deviance: 612.64  on 651  degrees of freedom
## AIC: 710.64
##
## Number of Fisher Scoring iterations: 5
```

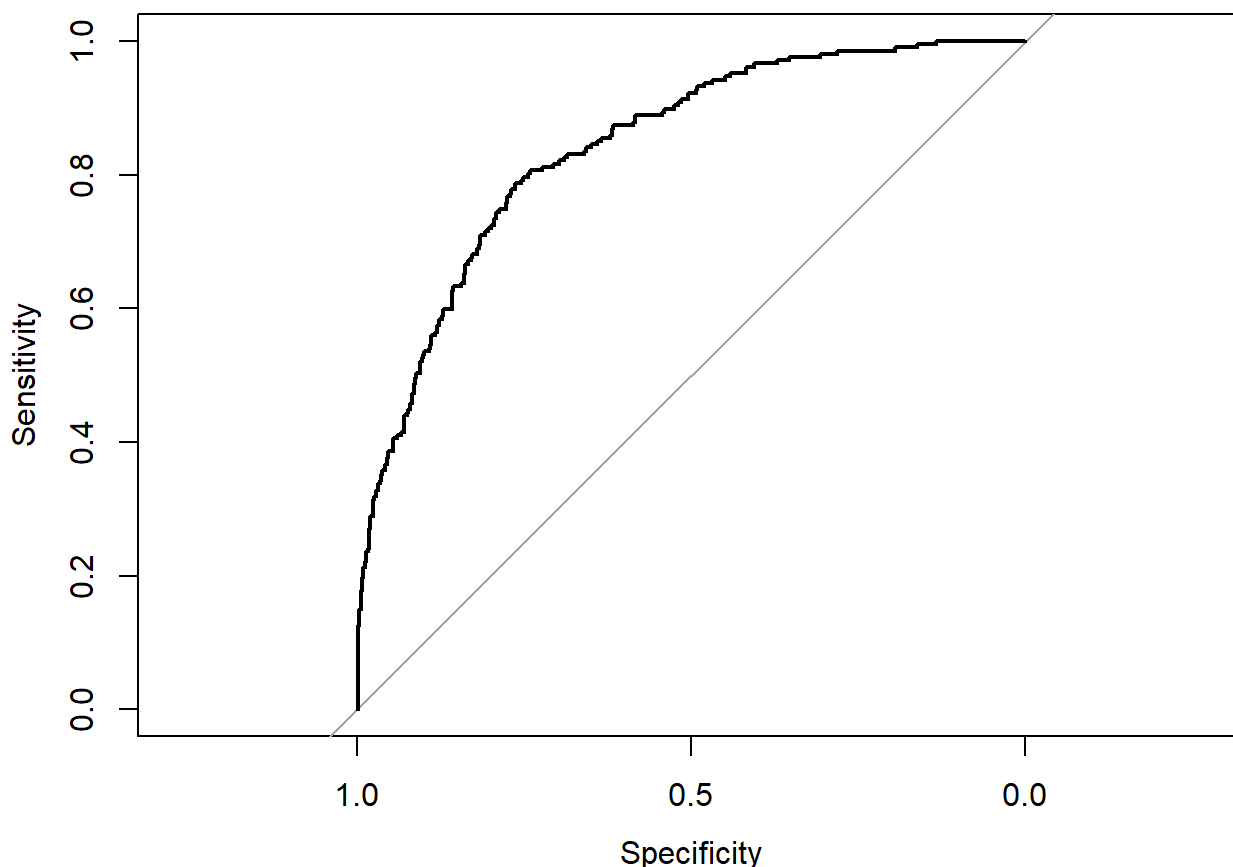
```
# get the probability for each data point
fitted = predict(ger_credit_logreg, type = 'response')%>%as.data.frame()
#fitted

# plot the ROC
# Look at reference 4
datas = data.frame(ornigial = ger_credit_train$V21, prob=fitted$.)
roc = roc(ornigial~prob, data = datas)
```

```
## Setting levels: control = 0, case = 1
```

```
## Setting direction: controls < cases
```

```
plot(roc)
```



Question 10.3.2

Because the model gives a result between 0 and 1, it requires setting a threshold probability to separate between “good” and “bad” answers. In this data set, they estimate that incorrectly identifying a bad customer as good, is 5 times worse than incorrectly classifying a good customer as bad. Determine a good threshold probability based on your model.

Looking at the ROC 8 looks to be the appropriate threshold.

References:

[1] alejandro_hagan. (2024). How to visualize random forest model output in R? Stack Overflow.

<https://stackoverflow.com/questions/73898275/how-to-visualize-random-forest-model-output-in-r>

(<https://stackoverflow.com/questions/73898275/how-to-visualize-random-forest-model-output-in-r>)

[2] araastat. (2018, July 17). Error in library(“reptree”) : there is no package called “reptree” · Issue #12 ·

araastat/reptree. GitHub. <https://github.com/araastat/reptree/issues/12>

(<https://github.com/araastat/reptree/issues/12>)

[3] Hofmann, H. (1994). Statlog (German Credit Data) [Dataset]. UCI Machine Learning Repository.

<https://doi.org/10.24432/C5NC77> (<https://doi.org/10.24432/C5NC77>).

[4] roc function - RDocumentation. (n.d.). [Www.rdocumentation.org](http://www.rdocumentation.org).

<https://www.rdocumentation.org/packages/pROC/versions/1.18.5/topics/roc>

(<https://www.rdocumentation.org/packages/pROC/versions/1.18.5/topics/roc>)

[5] APA citation