# Logistic Regression

Load the train data set on which ROSE was applied

```r
train_rose <- read.csv("data/train_rose.csv")
test <- read.csv("data/test.csv")
```

Remove the column MOSTYPE.41 created during one-hot encoding because it is redundant.

```r
train_rose  <- subset(train_rose, select = -c(MOSTYPE.41))
test <- subset(test, select = -c(MOSTYPE.41))
```

## Complete model

Fit a logistic regression model using all the variables

```r
complete_model <- glm(
    CARAVAN ~ .,
    data = train_rose,
    family = "binomial"
    )
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```r
summary(complete_model)
```

```
##
## Call:
## glm(formula = CARAVAN ~ ., family = "binomial", data = train_rose)
##
## Coefficients:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept) -5.758e+00  8.404e-01  -6.851 7.32e-12 ***
## MOSTYPE.1    7.198e+00  5.382e+00   1.337 0.181094
## MOSTYPE.2    4.951e+00  3.423e+00   1.446 0.148074
## MOSTYPE.3    6.259e+00  3.020e+00   2.072 0.038238 *
## MOSTYPE.4   -6.074e+00  3.896e+00  -1.559 0.119013
## MOSTYPE.5    5.204e+00  3.095e+00   1.681 0.092676 .
## MOSTYPE.6    6.211e+00  5.376e+00   1.155 0.247947
## MOSTYPE.7    3.189e+00  3.640e+00   0.876 0.380903
## MOSTYPE.8    1.112e+01  4.121e+00   2.699 0.006964 **
## MOSTYPE.9    2.103e+00  1.876e+00   1.121 0.262317
## MOSTYPE.10   7.769e-02  5.291e+00   0.015 0.988286
## MOSTYPE.11   5.751e+00  3.209e+00   1.792 0.073113 .
## MOSTYPE.12   1.187e+01  4.307e+00   2.756 0.005849 **
## MOSTYPE.13   4.991e+00  3.169e+00   1.575 0.115265
## MOSTYPE.15  -1.290e+02  9.363e+03  -0.014 0.989011
## MOSTYPE.16  -1.389e+02  4.862e+03  -0.029 0.977200
## MOSTYPE.17  -1.477e+02  6.742e+03  -0.022 0.982520
## MOSTYPE.18  -1.382e+02  4.432e+03  -0.031 0.975133
## MOSTYPE.19  -1.429e+02  1.198e+04  -0.012 0.990484
## MOSTYPE.20   6.553e+00  4.117e+00   1.592 0.111443
## MOSTYPE.21  -1.368e+02  5.287e+03  -0.026 0.979356
## MOSTYPE.22  -2.852e+00  3.533e+00  -0.807 0.419610
## MOSTYPE.23  -7.690e+00  2.656e+00  -2.895 0.003786 **
## MOSTYPE.24  -1.134e+00  3.389e+00  -0.335 0.737841
## MOSTYPE.25  -3.325e+00  4.626e+00  -0.719 0.472276
```

```
## MOSTYPE.26  -4.396e+00  5.473e+00  -0.803 0.421891
## MOSTYPE.27  -4.655e+00  5.066e+00  -0.919 0.358133
## MOSTYPE.28  -1.414e+02  3.830e+03  -0.037 0.970548
## MOSTYPE.29  -3.838e+00  2.913e+00  -1.317 0.187692
## MOSTYPE.30   1.248e+00  3.465e+00   0.360 0.718823
## MOSTYPE.31   1.085e+00  4.399e+00   0.247 0.805251
## MOSTYPE.32   4.875e+00  4.368e+00   1.116 0.264414
## MOSTYPE.33   4.514e+00  2.047e+00   2.205 0.027431 *
## MOSTYPE.34   2.725e+00  3.159e+00   0.863 0.388385
## MOSTYPE.35  -2.069e-01  2.280e+00  -0.091 0.927701
## MOSTYPE.36   7.606e+00  2.244e+00   3.390 0.000698 ***
## MOSTYPE.37   7.529e+00  1.951e+00   3.859 0.000114 ***
## MOSTYPE.38   7.328e+00  1.777e+00   4.124 3.72e-05 ***
## MOSTYPE.39   7.084e+00  2.095e+00   3.382 0.000719 ***
## MOSTYPE.40  -1.355e+02  2.255e+03  -0.060 0.952073
## MAANTHUI    -1.875e+00  7.686e-01  -2.439 0.014720 *
## MGEMOMV     -3.843e-01  2.244e-01  -1.713 0.086796 .
## MGEMLEEF     1.186e+00  2.148e-01   5.520 3.38e-08 ***
## MGODRK      -2.828e-01  2.557e-01  -1.106 0.268804
## MGODPR       6.611e-01  1.561e-01   4.235 2.29e-05 ***
## MGODOV       5.679e-01  2.256e-01   2.517 0.011827 *
## MRELGE       7.581e-01  1.959e-01   3.870 0.000109 ***
## MRELSA       7.279e-01  3.381e-01   2.153 0.031327 *
## MFGEKIND    -6.768e-01  1.680e-01  -4.028 5.61e-05 ***
## MOPLHOOG     2.376e+00  2.678e-01   8.872  < 2e-16 ***
## MOPLMIDD     1.117e+00  1.912e-01   5.844 5.11e-09 ***
## MBERHOOG     9.540e-01  3.366e-01   2.834 0.004597 **
## MBERZELF     1.041e+00  3.750e-01   2.775 0.005523 **
## MBERBOER    -1.295e-01  3.831e-01  -0.338 0.735360
## MBERMIDD     8.913e-01  3.223e-01   2.765 0.005688 **
## MBERARBG    -2.133e-02  3.125e-01  -0.068 0.945589
## MBERARBO     4.432e-01  3.134e-01   1.414 0.157330
## MSKA        -7.654e-01  3.729e-01  -2.052 0.040139 *
## MSKB1       -3.766e-01  3.527e-01  -1.068 0.285598
## MSKB2       -1.375e-01  3.193e-01  -0.431 0.666728
## MSKC         1.190e+00  3.446e-01   3.453 0.000555 ***
## MSKD        -3.765e-01  3.338e-01  -1.128 0.259380
## MHHUUR      -2.533e-01  1.004e-01  -2.522 0.011667 *
## MAUT1        4.639e-01  2.144e-01   2.164 0.030491 *
## MAUT2       -9.575e-02  2.558e-01  -0.374 0.708188
## MZFONDS      3.369e-01  1.638e-01   2.057 0.039659 *
## MINKM30      5.020e-01  3.658e-01   1.372 0.169938
## MINK3045     7.467e-01  3.503e-01   2.132 0.033046 *
## MINK4575     4.230e-01  3.519e-01   1.202 0.229377
## MINK7512     6.947e-01  3.639e-01   1.909 0.056250 .
## MINK123M    -2.227e+00  5.191e-01  -4.290 1.79e-05 ***
## MINKGEM      8.550e-01  3.660e-01   2.336 0.019487 *
## MKOOPKLA    -4.631e-01  1.266e+00  -0.366 0.714445
## PWAPART      1.095e+00  2.690e-01   4.071 4.69e-05 ***
## PWABEDR     -4.859e-01  7.166e-01  -0.678 0.497720
## PWALAND     -3.439e+00  6.254e-01  -5.499 3.82e-08 ***
## PPERSAUT     2.078e+00  7.754e-02  26.799  < 2e-16 ***
## PBESAUT      3.010e-01  4.498e-01   0.669 0.503370
## PMOTSCO     -7.989e-01  2.334e-01  -3.422 0.000621 ***
```

```
## PVRAAUT     -3.187e+01  1.100e+03  -0.029 0.976884
## PAANHANG     1.327e+00  1.147e+00   1.157 0.247217
## PTRACTOR     7.668e-02  4.649e-01   0.165 0.868985
## PWERKT      -5.407e+01  1.145e+03  -0.047 0.962325
## PBROM       -2.433e+00  3.738e-01  -6.509 7.57e-11 ***
## PLEVEN      -1.932e-01  2.133e-01  -0.906 0.365029
## PPERSONG    -1.508e+00  1.409e+00  -1.070 0.284664
## PGEZONG      2.144e+00  9.297e-01   2.306 0.021115 *
## PWAOREG      2.770e+00  5.038e-01   5.499 3.83e-08 ***
## PBRAND       1.130e+00  1.456e-01   7.760 8.51e-15 ***
## PZEILPL     -7.506e+00  9.552e+00  -0.786 0.432007
## PPLEZIER     7.329e+00  7.759e-01   9.446  < 2e-16 ***
## PFIETS       6.029e+00  1.190e+00   5.065 4.08e-07 ***
## PINBOED     -1.680e+00  1.103e+00  -1.523 0.127797
## PBYSTAND     1.726e+00  4.259e-01   4.053 5.05e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 15177  on 10947  degrees of freedom
## Residual deviance: 11897  on 10854  degrees of freedom
## AIC: 12085
##
## Number of Fisher Scoring iterations: 15
```
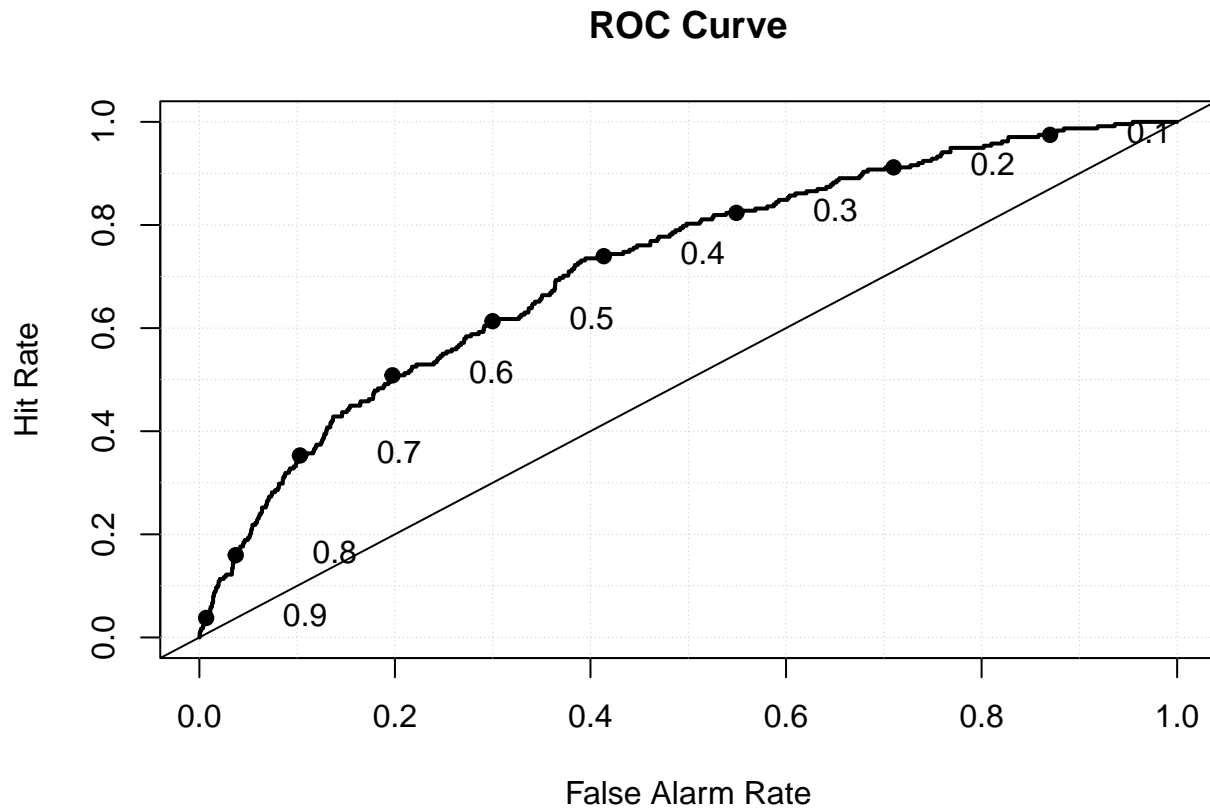
Confusion matrix on the test data set for the complete model

```
pred <- predict(complete_model, newdata = test, type = "response")
```

ROC curve for the complete model

```
roc.plot(test$CARAVAN, pred)
```

```
## Warning in roc.plot.default(test$CARAVAN, pred): Large amount of unique
## predictions used as thresholds. Consider specifying thresholds.
```

## ROC Curve



Confusion matrix for the complete model

```
pred <- ifelse(pred > 0.5, 1, 0)
confusion_matrix <- confusionMatrix(as.factor(test$CARAVAN), as.factor(pred), positive = "1")
print(confusion_matrix)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##          0 2633 1129
##          1   91  147
##
##                Accuracy : 0.695
##                  95% CI : (0.6805, 0.7092)
##     No Information Rate : 0.681
##     P-Value [Acc > NIR] : 0.0295
##
##                   Kappa : 0.1044
##
##  Mcnemar's Test P-Value : <2e-16
##
##             Sensitivity : 0.11520
##             Specificity : 0.96659
##          Pos Pred Value : 0.61765
##          Neg Pred Value : 0.69989
##              Prevalence : 0.31900
##          Detection Rate : 0.03675
##    Detection Prevalence : 0.05950
```

```
##        Balanced Accuracy : 0.54090
##
##          'Positive' Class : 1
##
```
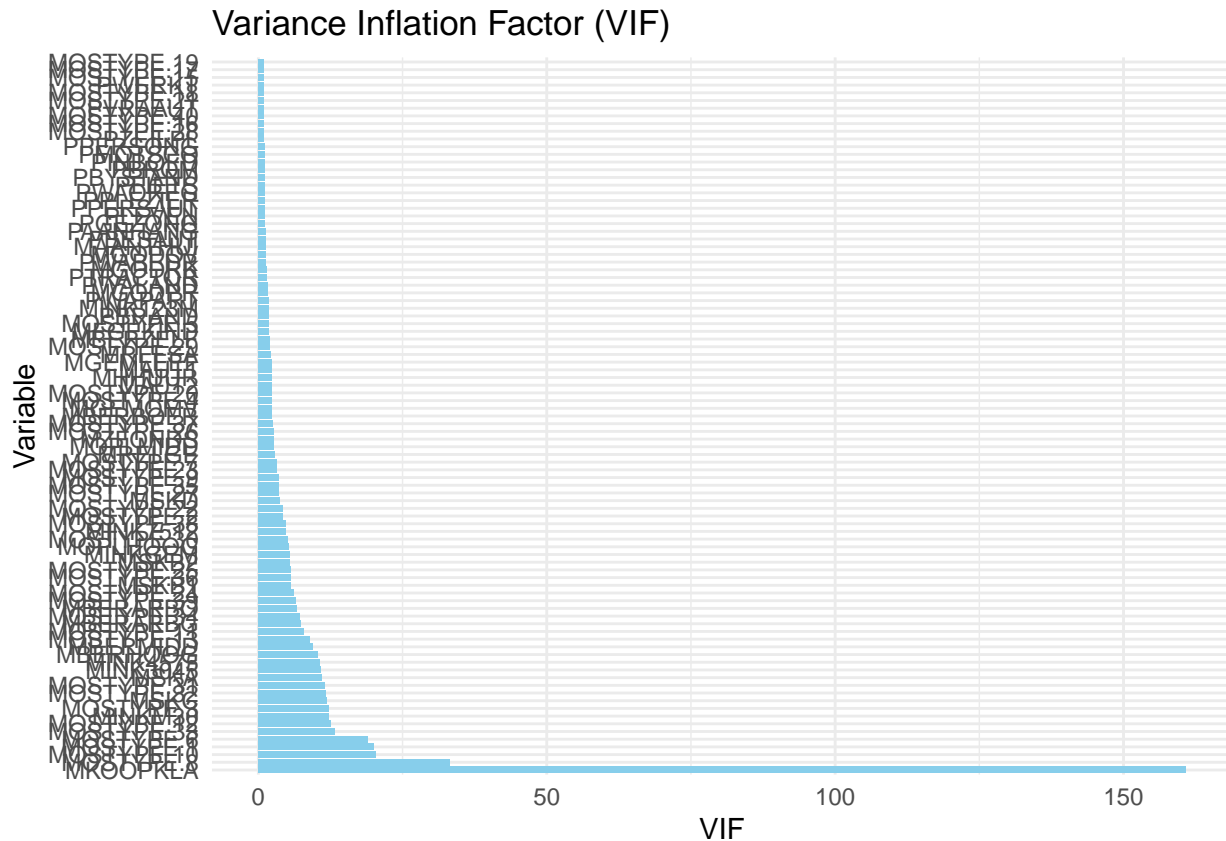
The model uses all the variables, but some of them may not be significant. We can use the Variance Inflation Factor (VIF) to check for multicollinearity.

## Variance Inflation Factor

Calculate the VIF for the complete, reduced models, and the stepwise model. Plot the results using a color map

```
vif_visualizer <- function(model) {
  # Calculate the VIF for the complete model
  vif_values <- vif(model)
  # Create a data frame for plotting
  vif_data <- data.frame(variable = names(vif_values), vif = vif_values)
  # Plot VIF
  ggplot(vif_data, aes(x = reorder(variable, -vif), y = vif)) +
    geom_bar(stat = "identity", fill = "skyblue") +
    coord_flip() +
    labs(title = "Variance Inflation Factor (VIF)",
         x = "Variable",
         y = "VIF") +
    theme_minimal()
}

vif_visualizer(complete_model)
```

**Variance Inflation Factor (VIF)**

The VIF plot shows that are variables that cause multicollinearity. We have to remove some varialbes to reduce the multicollinearity and to make the model model more interpretable.

## Reduced model

We can start by dropping the variables that are likely to be not different from zero and with a very small effect on the response variable. Select the variables that have a p-value less than 0.10 and a coefficient greater than 0.01. This seems to be a good compromise between the number of variables and the model s performance.

```r
vars <- names(which(summary(complete_model)$coefficients[,4] < 0.10 & summary(complete_model)$coefficien
print(vars)
```

```
##  [1] "MOSTYPE.3"  "MOSTYPE.5"  "MOSTYPE.8"  "MOSTYPE.11" "MOSTYPE.12"
##  [6] "MOSTYPE.33" "MOSTYPE.36" "MOSTYPE.37" "MOSTYPE.38" "MOSTYPE.39"
## [11] "MGEMLEEF"   "MGODPR"     "MGODOV"     "MRELGE"     "MRELSA"
## [16] "MOPLHOOG"   "MOPLMIDD"   "MBERHOOG"   "MBERZELF"   "MBERMIDD"
## [21] "MSKC"       "MAUT1"      "MZFONDS"    "MINK3045"   "MINK7512"
## [26] "MINKGEM"    "PWAPART"    "PPERSAUT"   "PGEZONG"    "PWAOREG"
## [31] "PBRAND"     "PPLEZIER"   "PFIETS"     "PBYSTAND"
```

Fit a logistic regression model using only the selected variables

```r
vars <- c(vars, "CARAVAN")
reduced_model <- glm(
    CARAVAN ~ .,
    data = train_rose[, vars],
    family = "binomial"
    )
```

```r
summary(reduced_model)
```

```
##
## Call:
## glm(formula = CARAVAN ~ ., family = "binomial", data = train_rose[,
##     vars])
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) -6.12334    0.28225 -21.695  < 2e-16 ***
## MOSTYPE.3    4.39380    0.91604   4.797 1.61e-06 ***
## MOSTYPE.5    1.98685    2.32739   0.854  0.39328
## MOSTYPE.8    8.46999    0.79666  10.632  < 2e-16 ***
## MOSTYPE.11   3.13478    1.20780   2.595  0.00945 **
## MOSTYPE.12  10.59961    1.29718   8.171 3.05e-16 ***
## MOSTYPE.33   3.74722    0.63948   5.860 4.63e-09 ***
## MOSTYPE.36   7.32108    0.99454   7.361 1.82e-13 ***
## MOSTYPE.37   6.92530    1.25922   5.500 3.80e-08 ***
## MOSTYPE.38   6.77462    0.87957   7.702 1.34e-14 ***
## MOSTYPE.39   5.21889    0.86479   6.035 1.59e-09 ***
## MGEMLEEF     0.93576    0.16701   5.603 2.11e-08 ***
## MGODPR       0.80788    0.13318   6.066 1.31e-09 ***
## MGODOV       0.71904    0.20924   3.436  0.00059 ***
## MRELGE       0.70569    0.15585   4.528 5.96e-06 ***
## MRELSA       0.03064    0.28846   0.106  0.91541
## MOPLHOOG     2.05049    0.19881  10.314  < 2e-16 ***
## MOPLMIDD     0.93640    0.15408   6.077 1.22e-09 ***
## MBERHOOG     0.95475    0.18601   5.133 2.85e-07 ***
## MBERZELF     0.60295    0.28651   2.104  0.03534 *
## MBERMIDD     0.83813    0.14092   5.947 2.72e-09 ***
## MSKC         1.41076    0.15941   8.850  < 2e-16 ***
## MAUT1        0.66519    0.15273   4.355 1.33e-05 ***
## MZFONDS      0.40350    0.14278   2.826  0.00471 **
## MINK3045     0.41073    0.12117   3.390  0.00070 ***
## MINK7512     0.63436    0.22435   2.828  0.00469 **
## MINKGEM      0.47462    0.24551   1.933  0.05321 .
## PWAPART      1.41269    0.24286   5.817 5.99e-09 ***
## PPERSAUT     2.08118    0.07433  27.998  < 2e-16 ***
## PGEZONG      1.34053    0.83395   1.607  0.10796
## PWAOREG      2.74181    0.44145   6.211 5.27e-10 ***
## PBRAND       0.84157    0.12843   6.553 5.65e-11 ***
## PPLEZIER     6.65548    0.68920   9.657  < 2e-16 ***
## PFIETS       6.38188    1.14083   5.594 2.22e-08 ***
## PBYSTAND     2.12510    0.41279   5.148 2.63e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 15177  on 10947  degrees of freedom
## Residual deviance: 12323  on 10913  degrees of freedom
## AIC: 12393
##
## Number of Fisher Scoring iterations: 5
```

Confusion matrix on the test data set for the reduced model

```
pred <- predict(reduced_model, newdata = test, type = "response")
pred <- ifelse(pred > 0.5, 1, 0)
confusion_matrix <- confusionMatrix(as.factor(test$CARAVAN), as.factor(pred), positive = "1")
print(confusion_matrix)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##          0 2642 1120
##          1   96  142
##
##                Accuracy : 0.696
##                  95% CI : (0.6815, 0.7102)
##     No Information Rate : 0.6845
##     P-Value [Acc > NIR] : 0.06044
##
##                   Kappa : 0.0991
##
##  Mcnemar's Test P-Value : < 2e-16
##
##             Sensitivity : 0.1125
##             Specificity : 0.9649
##          Pos Pred Value : 0.5966
##          Neg Pred Value : 0.7023
##              Prevalence : 0.3155
##          Detection Rate : 0.0355
##    Detection Prevalence : 0.0595
##       Balanced Accuracy : 0.5387
##
##        'Positive' Class : 1
##
```

We removed 61 variables from the complete model, and the model performance metrics did not change significantly.

## Stepwise regression

Considering the variables of the reduced model, perform a stepwise regression in order to select the best subset of variables. We use the Bayesian Information Criterion (BIC) to select the best model because we have a large number of variables and the BIC penalizes the number of variables in the model more than Akaike Information Criterion.

```
# Stepwise regression
stepwise_model <- step(
    glm(
        CARAVAN ~ .,
        data = train_rose[, vars],
        family = "binomial",
    ),
    k = log(nrow(train_rose)),
    direction = "both"
)
```

```
## Start:  AIC=12648.22
## CARAVAN ~ MOSTYPE.3 + MOSTYPE.5 + MOSTYPE.8 + MOSTYPE.11 + MOSTYPE.12 +
##     MOSTYPE.33 + MOSTYPE.36 + MOSTYPE.37 + MOSTYPE.38 + MOSTYPE.39 +
##     MGEMLEEF + MGODPR + MGODOV + MRELGE + MRELSA + MOPLHOOG +
##     MOPLMIDD + MBERHOOG + MBERZELF + MBERMIDD + MSKC + MAUT1 +
##     MZFONDS + MINK3045 + MINK7512 + MINKGEM + PWAPART + PPERSAUT +
##     PGEZONG + PWAOREG + PBRAND + PPLEZIER + PFIETS + PBYSTAND
##
##              Df Deviance   AIC
## - MRELSA      1    12323 12639
## - MOSTYPE.5   1    12323 12640
## - PGEZONG     1    12325 12642
## - MINKGEM     1    12326 12643
## - MBERZELF    1    12327 12643
## - MOSTYPE.11  1    12330 12646
## - MINK7512    1    12331 12647
## - MZFONDS     1    12331 12647
## <none>             12323 12648
## - MINK3045    1    12334 12650
## - MGODOV      1    12334 12651
## - MAUT1       1    12342 12658
## - MRELGE      1    12343 12660
## - MOSTYPE.3   1    12346 12662
## - MBERHOOG    1    12349 12665
## - PBYSTAND    1    12352 12668
## - MOSTYPE.37  1    12353 12670
## - MGEMLEEF    1    12354 12670
## - PFIETS      1    12355 12671
## - PWAPART     1    12356 12673
## - MOSTYPE.33  1    12357 12673
## - MBERMIDD    1    12358 12674
## - MOSTYPE.39  1    12359 12675
## - MGODPR      1    12360 12676
## - MOPLMIDD    1    12360 12676
## - PBRAND      1    12366 12682
## - PWAOREG     1    12369 12685
## - MOSTYPE.36  1    12377 12693
## - MOSTYPE.38  1    12382 12698
## - MOSTYPE.12  1    12394 12710
## - MSKC        1    12403 12719
## - MOPLHOOG    1    12432 12749
## - MOSTYPE.8   1    12442 12758
## - PPLEZIER    1    12506 12822
## - PPERSAUT    1    13164 13480
##
## Step:  AIC=12638.93
## CARAVAN ~ MOSTYPE.3 + MOSTYPE.5 + MOSTYPE.8 + MOSTYPE.11 + MOSTYPE.12 +
##     MOSTYPE.33 + MOSTYPE.36 + MOSTYPE.37 + MOSTYPE.38 + MOSTYPE.39 +
##     MGEMLEEF + MGODPR + MGODOV + MRELGE + MOPLHOOG + MOPLMIDD +
##     MBERHOOG + MBERZELF + MBERMIDD + MSKC + MAUT1 + MZFONDS +
##     MINK3045 + MINK7512 + MINKGEM + PWAPART + PPERSAUT + PGEZONG +
##     PWAOREG + PBRAND + PPLEZIER + PFIETS + PBYSTAND
##
##              Df Deviance   AIC
```

```
## - MOSTYPE.5   1    12323 12630
## - PGEZONG     1    12325 12632
## - MINKGEM     1    12326 12633
## - MBERZELF    1    12327 12634
## - MOSTYPE.11  1    12330 12636
## - MZFONDS     1    12331 12638
## - MINK7512    1    12331 12638
## <none>             12323 12639
## - MINK3045    1    12334 12641
## - MGODOV      1    12335 12642
## + MRELSA      1    12323 12648
## - MAUT1       1    12342 12649
## - MOSTYPE.3   1    12346 12653
## - MBERHOOG    1    12349 12656
## - MRELGE      1    12350 12657
## - PBYSTAND    1    12352 12659
## - MOSTYPE.37  1    12354 12660
## - PFIETS      1    12355 12662
## - PWAPART     1    12356 12664
## - MOSTYPE.33  1    12357 12664
## - MGEMLEEF    1    12358 12664
## - MBERMIDD    1    12358 12665
## - MOSTYPE.39  1    12359 12666
## - MGODPR      1    12360 12667
## - MOPLMIDD    1    12361 12668
## - PBRAND      1    12366 12673
## - PWAOREG     1    12369 12676
## - MOSTYPE.36  1    12377 12684
## - MOSTYPE.38  1    12382 12689
## - MOSTYPE.12  1    12394 12701
## - MSKC        1    12403 12710
## - MOPLHOOG    1    12432 12739
## - MOSTYPE.8   1    12442 12749
## - PPLEZIER    1    12506 12813
## - PPERSAUT    1    13164 13471
##
## Step:  AIC=12630.35
## CARAVAN ~ MOSTYPE.3 + MOSTYPE.8 + MOSTYPE.11 + MOSTYPE.12 + MOSTYPE.33 +
##     MOSTYPE.36 + MOSTYPE.37 + MOSTYPE.38 + MOSTYPE.39 + MGEMLEEF +
##     MGODPR + MGODOV + MRELGE + MOPLHOOG + MOPLMIDD + MBERHOOG +
##     MBERZELF + MBERMIDD + MSKC + MAUT1 + MZFONDS + MINK3045 +
##     MINK7512 + MINKGEM + PWAPART + PPERSAUT + PGEZONG + PWAOREG +
##     PBRAND + PPLEZIER + PFIETS + PBYSTAND
##
##             Df Deviance   AIC
## - PGEZONG     1    12326 12624
## - MINKGEM     1    12327 12625
## - MBERZELF    1    12328 12626
## - MOSTYPE.11  1    12330 12628
## - MINK7512    1    12332 12629
## - MZFONDS     1    12332 12629
## <none>             12323 12630
## - MINK3045    1    12335 12633
## - MGODOV      1    12336 12633
```

```
## + MOSTYPE.5    1     12323 12639
## + MRELSA       1     12323 12640
## - MAUT1        1     12342 12640
## - MOSTYPE.3    1     12346 12644
## - MBERHOOG     1     12350 12648
## - MRELGE       1     12351 12648
## - PBYSTAND     1     12352 12650
## - MOSTYPE.37   1     12354 12652
## - PFIETS       1     12356 12653
## - PWAPART      1     12357 12654
## - MOSTYPE.33   1     12357 12655
## - MGEMLEEF     1     12359 12656
## - MOSTYPE.39   1     12359 12657
## - MBERMIDD     1     12360 12658
## - MGODPR       1     12361 12658
## - MOPLMIDD     1     12361 12659
## - PBRAND       1     12367 12665
## - PWAOREG      1     12370 12668
## - MOSTYPE.36   1     12377 12675
## - MOSTYPE.38   1     12383 12680
## - MOSTYPE.12   1     12394 12692
## - MSKC         1     12404 12702
## - MOPLHOOG     1     12433 12731
## - MOSTYPE.8    1     12442 12740
## - PPLEZIER     1     12506 12804
## - PPERSAUT     1     13165 13462
##
## Step:  AIC=12623.78
## CARAVAN ~ MOSTYPE.3 + MOSTYPE.8 + MOSTYPE.11 + MOSTYPE.12 + MOSTYPE.33 +
##     MOSTYPE.36 + MOSTYPE.37 + MOSTYPE.38 + MOSTYPE.39 + MGEMLEEF +
##     MGODPR + MGODOV + MRELGE + MOPLHOOG + MOPLMIDD + MBERHOOG +
##     MBERZELF + MBERMIDD + MSKC + MAUT1 + MZFONDS + MINK3045 +
##     MINK7512 + MINKGEM + PWAPART + PPERSAUT + PWAOREG + PBRAND +
##     PPLEZIER + PFIETS + PBYSTAND
##
##              Df Deviance   AIC
## - MINKGEM     1     12330 12618
## - MBERZELF    1     12330 12619
## - MOSTYPE.11  1     12332 12621
## - MZFONDS     1     12334 12623
## - MINK7512    1     12335 12623
## <none>              12326 12624
## - MINK3045    1     12338 12626
## - MGODOV      1     12338 12626
## + PGEZONG     1     12323 12630
## + MOSTYPE.5   1     12325 12632
## + MRELSA      1     12326 12633
## - MAUT1       1     12346 12634
## - MOSTYPE.3   1     12349 12637
## - MBERHOOG    1     12353 12641
## - MRELGE      1     12354 12642
## - MOSTYPE.37  1     12356 12645
## - PBYSTAND    1     12358 12646
## - PWAPART     1     12359 12647
```

```
## - PFIETS        1    12359 12648
## - MOSTYPE.33 1    12362 12650
## - MGEMLEEF    1    12362 12650
## - MOSTYPE.39 1    12362 12650
## - MBERMIDD    1    12362 12650
## - MGODPR      1    12364 12652
## - MOPLMIDD    1    12364 12653
## - PBRAND      1    12371 12660
## - PWAOREG     1    12373 12661
## - MOSTYPE.36 1    12380 12668
## - MOSTYPE.38 1    12386 12674
## - MOSTYPE.12 1    12397 12685
## - MSKC        1    12406 12695
## - MOPLHOOG    1    12436 12724
## - MOSTYPE.8  1    12445 12733
## - PPLEZIER    1    12509 12797
## - PPERSAUT    1    13174 13462
##
## Step:  AIC=12618.26
## CARAVAN ~ MOSTYPE.3 + MOSTYPE.8 + MOSTYPE.11 + MOSTYPE.12 + MOSTYPE.33 +
##     MOSTYPE.36 + MOSTYPE.37 + MOSTYPE.38 + MOSTYPE.39 + MGEMLEEF +
##     MGODPR + MGODOV + MRELGE + MOPLHOOG + MOPLMIDD + MBERHOOG +
##     MBERZELF + MBERMIDD + MSKC + MAUT1 + MZFONDS + MINK3045 +
##     MINK7512 + PWAPART + PPERSAUT + PWAOREG + PBRAND + PPLEZIER +
##     PFIETS + PBYSTAND
##
##              Df Deviance   AIC
## - MBERZELF  1    12335 12614
## - MOSTYPE.11 1    12337 12616
## - MZFONDS    1    12338 12617
## <none>           12330 12618
## - MINK3045  1    12340 12619
## - MGODOV     1    12342 12621
## + MINKGEM    1    12326 12624
## + PGEZONG    1    12327 12625
## + MOSTYPE.5 1    12329 12627
## + MRELSA     1    12330 12628
## - MAUT1      1    12350 12629
## - MOSTYPE.3 1    12354 12633
## - MINK7512  1    12354 12634
## - MBERHOOG   1    12358 12637
## - MOSTYPE.37 1    12361 12640
## - PBYSTAND   1    12361 12640
## - PFIETS     1    12363 12642
## - PWAPART    1    12363 12642
## - MGEMLEEF   1    12363 12642
## - MRELGE     1    12364 12643
## - MOSTYPE.39 1    12365 12644
## - MOSTYPE.33 1    12366 12645
## - MBERMIDD   1    12367 12646
## - MGODPR     1    12369 12648
## - MOPLMIDD   1    12370 12649
## - PWAOREG    1    12376 12655
## - PBRAND     1    12376 12655
```

```
## - MOSTYPE.36  1    12385 12664
## - MOSTYPE.38  1    12390 12668
## - MOSTYPE.12  1    12403 12682
## - MSKC        1    12409 12688
## - MOPLHOOG    1    12452 12731
## - MOSTYPE.8   1    12455 12734
## - PPLEZIER    1    12513 12792
## - PPERSAUT    1    13179 13458
##
## Step:  AIC=12613.97
## CARAVAN ~ MOSTYPE.3 + MOSTYPE.8 + MOSTYPE.11 + MOSTYPE.12 + MOSTYPE.33 +
##     MOSTYPE.36 + MOSTYPE.37 + MOSTYPE.38 + MOSTYPE.39 + MGEMLEEF +
##     MGODPR + MGODOV + MRELGE + MOPLHOOG + MOPLMIDD + MBERHOOG +
##     MBERMIDD + MSKC + MAUT1 + MZFONDS + MINK3045 + MINK7512 +
##     PWAPART + PPERSAUT + PWAOREG + PBRAND + PPLEZIER + PFIETS +
##     PBYSTAND
##
##               Df Deviance   AIC
## - MZFONDS      1    12341 12611
## - MOSTYPE.11   1    12342 12611
## <none>              12335 12614
## - MINK3045     1    12345 12615
## - MGODOV       1    12348 12617
## + MBERZELF     1    12330 12618
## + MINKGEM      1    12330 12619
## + PGEZONG      1    12333 12621
## + MOSTYPE.5    1    12334 12622
## + MRELSA       1    12335 12623
## - MAUT1        1    12354 12624
## - MBERHOOG     1    12359 12628
## - MOSTYPE.3    1    12359 12628
## - MINK7512     1    12360 12629
## - MOSTYPE.37   1    12365 12635
## - PBYSTAND     1    12366 12635
## - PFIETS       1    12367 12637
## - MBERMIDD     1    12368 12637
## - PWAPART      1    12368 12638
## - MOSTYPE.39   1    12370 12639
## - MOSTYPE.33   1    12370 12640
## - MGEMLEEF     1    12372 12642
## - MRELGE       1    12372 12642
## - MGODPR       1    12374 12644
## - MOPLMIDD     1    12382 12651
## - PBRAND       1    12382 12652
## - PWAOREG      1    12384 12654
## - MOSTYPE.36   1    12389 12659
## - MOSTYPE.38   1    12394 12663
## - MOSTYPE.12   1    12407 12676
## - MSKC         1    12413 12683
## - MOSTYPE.8    1    12458 12728
## - MOPLHOOG     1    12473 12743
## - PPLEZIER     1    12522 12791
## - PPERSAUT     1    13182 13451
##
```

```
## Step:  AIC=12610.86
## CARAVAN ~ MOSTYPE.3 + MOSTYPE.8 + MOSTYPE.11 + MOSTYPE.12 + MOSTYPE.33 +
##     MOSTYPE.36 + MOSTYPE.37 + MOSTYPE.38 + MOSTYPE.39 + MGEMLEEF +
##     MGODPR + MGODOV + MRELGE + MOPLHOOG + MOPLMIDD + MBERHOOG +
##     MBERMIDD + MSKC + MAUT1 + MINK3045 + MINK7512 + PWAPART +
##     PPERSAUT + PWAOREG + PBRAND + PPLEZIER + PFIETS + PBYSTAND
##
##              Df Deviance   AIC
## - MOSTYPE.11  1    12348 12608
## <none>             12341 12611
## - MINK3045    1    12353 12613
## + MZFONDS     1    12335 12614
## - MGODOV      1    12355 12616
## + MINKGEM     1    12337 12616
## + MBERZELF    1    12338 12617
## + PGEZONG     1    12339 12618
## + MOSTYPE.5   1    12340 12619
## - MBERHOOG    1    12360 12620
## + MRELSA      1    12341 12620
## - MAUT1       1    12362 12622
## - MOSTYPE.3   1    12364 12624
## - MINK7512    1    12365 12625
## - MOSTYPE.37  1    12371 12632
## - PBYSTAND    1    12372 12632
## - MBERMIDD    1    12372 12633
## - PFIETS      1    12374 12634
## - MRELGE      1    12374 12634
## - MOSTYPE.39  1    12375 12635
## - MGEMLEEF    1    12376 12637
## - PWAPART     1    12377 12637
## - MOSTYPE.33  1    12378 12638
## - MGODPR      1    12380 12640
## - MOPLMIDD    1    12384 12644
## - PBRAND      1    12386 12646
## - PWAOREG     1    12390 12651
## - MOSTYPE.36  1    12395 12655
## - MOSTYPE.38  1    12400 12660
## - MOSTYPE.12  1    12413 12674
## - MSKC        1    12421 12682
## - MOSTYPE.8   1    12462 12722
## - MOPLHOOG    1    12473 12734
## - PPLEZIER    1    12524 12785
## - PPERSAUT    1    13193 13454
##
## Step:  AIC=12608.27
## CARAVAN ~ MOSTYPE.3 + MOSTYPE.8 + MOSTYPE.12 + MOSTYPE.33 + MOSTYPE.36 +
##     MOSTYPE.37 + MOSTYPE.38 + MOSTYPE.39 + MGEMLEEF + MGODPR +
##     MGODOV + MRELGE + MOPLHOOG + MOPLMIDD + MBERHOOG + MBERMIDD +
##     MSKC + MAUT1 + MINK3045 + MINK7512 + PWAPART + PPERSAUT +
##     PWAOREG + PBRAND + PPLEZIER + PFIETS + PBYSTAND
##
##              Df Deviance   AIC
## <none>             12348 12608
## - MINK3045    1    12359 12610
```

```
## + MOSTYPE.11   1    12341 12611
## + MZFONDS      1    12342 12611
## + MINKGEM      1    12343 12613
## - MGODOV       1    12362 12614
## + MBERZELF     1    12345 12614
## + PGEZONG      1    12346 12615
## + MOSTYPE.5    1    12347 12617
## + MRELSA       1    12348 12617
## - MAUT1        1    12367 12618
## - MBERHOOG     1    12368 12619
## - MOSTYPE.3    1    12368 12620
## - MINK7512     1    12372 12623
## - MOSTYPE.37   1    12376 12627
## - MOSTYPE.39   1    12379 12630
## - PBYSTAND     1    12381 12632
## - MOSTYPE.33   1    12381 12632
## - MGEMLEEF     1    12381 12632
## - PFIETS       1    12381 12633
## - MRELGE       1    12382 12633
## - MBERMIDD     1    12382 12633
## - PWAPART      1    12383 12634
## - MOPLMIDD     1    12390 12641
## - MGODPR       1    12390 12641
## - PWAOREG      1    12396 12648
## - PBRAND       1    12397 12648
## - MOSTYPE.36   1    12399 12650
## - MOSTYPE.38   1    12404 12655
## - MOSTYPE.12   1    12417 12668
## - MSKC         1    12433 12684
## - MOSTYPE.8    1    12463 12714
## - MOPLHOOG     1    12478 12729
## - PPLEZIER     1    12530 12781
## - PPERSAUT     1    13196 13447
```

```
summary(stepwise_model)
```

```
##
## Call:
## glm(formula = CARAVAN ~ MOSTYPE.3 + MOSTYPE.8 + MOSTYPE.12 +
##     MOSTYPE.33 + MOSTYPE.36 + MOSTYPE.37 + MOSTYPE.38 + MOSTYPE.39 +
##     MGEMLEEF + MGODPR + MGODOV + MRELGE + MOPLHOOG + MOPLMIDD +
##     MBERHOOG + MBERMIDD + MSKC + MAUT1 + MINK3045 + MINK7512 +
##     PWAPART + PPERSAUT + PWAOREG + PBRAND + PPLEZIER + PFIETS +
##     PBYSTAND, family = "binomial", data = train_rose[, vars])
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -5.61098    0.21554 -26.032  < 2e-16 ***
## MOSTYPE.3    4.09539    0.90863   4.507 6.57e-06 ***
## MOSTYPE.8    8.19330    0.78281  10.466  < 2e-16 ***
## MOSTYPE.12  10.38805    1.28984   8.054 8.03e-16 ***
## MOSTYPE.33   3.61839    0.62928   5.750 8.92e-09 ***
## MOSTYPE.36   7.03967    0.98589   7.140 9.31e-13 ***
## MOSTYPE.37   6.61406    1.25247   5.281 1.29e-07 ***
## MOSTYPE.38   6.49172    0.86878   7.472 7.88e-14 ***
```

```
## MOSTYPE.39   4.78557     0.85764    5.580 2.41e-08 ***
## MGEMLEEF     0.89089     0.15464    5.761 8.37e-09 ***
## MGODPR       0.85004     0.13104    6.487 8.76e-11 ***
## MGODOV       0.79302     0.20815    3.810 0.000139 ***
## MRELGE       0.73313     0.12636    5.802 6.56e-09 ***
## MOPLHOOG     2.09887     0.18679   11.237  < 2e-16 ***
## MOPLMIDD     0.95751     0.14869    6.439 1.20e-10 ***
## MBERHOOG     0.75627     0.16806    4.500 6.80e-06 ***
## MBERMIDD     0.78734     0.13506    5.830 5.55e-09 ***
## MSKC         1.43977     0.15759    9.136  < 2e-16 ***
## MAUT1        0.66533     0.15115    4.402 1.07e-05 ***
## MINK3045     0.39705     0.11955    3.321 0.000897 ***
## MINK7512     0.89665     0.18275    4.906 9.28e-07 ***
## PWAPART      1.42191     0.24062    5.909 3.43e-09 ***
## PPERSAUT     2.08064     0.07403   28.107  < 2e-16 ***
## PWAOREG      2.83725     0.44612    6.360 2.02e-10 ***
## PBRAND       0.88786     0.12640    7.024 2.15e-12 ***
## PPLEZIER     6.61833     0.69063    9.583  < 2e-16 ***
## PFIETS       6.46924     1.13351    5.707 1.15e-08 ***
## PBYSTAND     2.24460     0.41196    5.449 5.08e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 15177  on 10947  degrees of freedom
## Residual deviance: 12348  on 10920  degrees of freedom
## AIC: 12404
##
## Number of Fisher Scoring iterations: 5
```
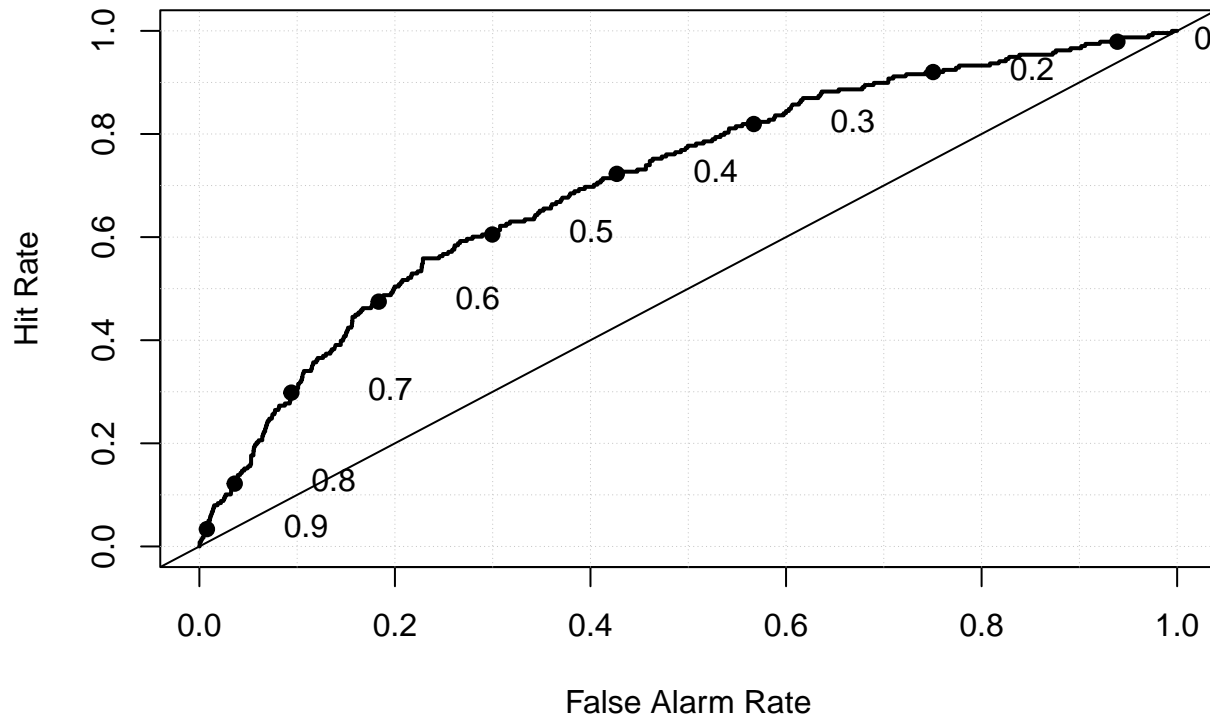
Confusion matrix on the test data set for the stepwise model

```
pred <- predict(stepwise_model, newdata = test, type = "response")
```

ROC curve for the stepwise model

```
roc.plot(test$CARAVAN, pred)
```

```
## Warning in roc.plot.default(test$CARAVAN, pred): Large amount of unique
## predictions used as thresholds. Consider specifying thresholds.
```

**ROC Curve**



```
pred <- ifelse(pred > 0.5, 1, 0)
confusion_matrix <- confusionMatrix(as.factor(test$CARAVAN), as.factor(pred), positive = "1")
print(confusion_matrix)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##          0 2634 1128
##          1   94  144
##
##                Accuracy : 0.6945
##                  95% CI : (0.68, 0.7088)
##     No Information Rate : 0.682
##     P-Value [Acc > NIR] : 0.04604
##
##                   Kappa : 0.1006
##
##  Mcnemar's Test P-Value : < 2e-16
##
##             Sensitivity : 0.1132
##             Specificity : 0.9655
##          Pos Pred Value : 0.6050
##          Neg Pred Value : 0.7002
##              Prevalence : 0.3180
##          Detection Rate : 0.0360
##    Detection Prevalence : 0.0595
##       Balanced Accuracy : 0.5394
```

```
##
##          'Positive' Class : 1
##
```

The stepwise model has 28 variables, and the model performance metrics are similar to the reduced model.

## Bayesian Information Criterion Summary

Calculate the BIC for the complete, reduced models, and the stepwise model

```
BIC(complete_model, reduced_model, stepwise_model)
```

```
##                df      BIC
## complete_model 94 12771.04
## reduced_model  35 12648.22
## stepwise_model 28 12608.27
```

## Interpretation and Conclusion

We reached the goal to reduce the number of variables without worsening the model performance.

The stepwise model, with 28 variables, reaches the best balanced accuracy on the test set, 0.5430. The factors that result to influence more the target variable are:

```
-MOSTYPE.12 Affluent young families,
if this variable is 1 (True), the log odds of
the target variable to be 1 increases by 10.38.


-PPLEZIER Contribution boat policies


-PFIETS Contribution bicycle policies


-MOSTYPE.36 Couples with teens, Married with children,
if this variable is 1 (True), the log odds of the
target variable to be 1 increases by 7.03.


-MOSTYPE.38 Traditional families,  if this variable
is 1 (True), the log odds of the target variable
to be 1 increases by 6.49
```

In particular these vairables have a positive effect on the target variable, meaning that the higher the value of the variable, the higher the probability of the target variable to be 1. The interpretation of the coefficient is note very straightforward because of the encoding and the normalization process applied to the data.