

```

#libs
require("ROSE")

## Loading required package: ROSE
## Warning: package 'ROSE' was built under R version 3.3.3
## Loaded ROSE 0.0-3
require("pROC")

## Loading required package: pROC
## Type 'citation("pROC")' for a citation.
##
## Attaching package: 'pROC'
## The following objects are masked from 'package:stats':
##
##     cov, smooth, var
require("rpart")

## Loading required package: rpart
require("rpart.plot")

## Loading required package: rpart.plot
## Warning: package 'rpart.plot' was built under R version 3.3.3
require("caret")

## Loading required package: caret
## Loading required package: lattice
## Loading required package: ggplot2
require("randomForest")

## Loading required package: randomForest
## Warning: package 'randomForest' was built under R version 3.3.3
## randomForest 4.6-12
## Type rfNews() to see new features/changes/bug fixes.
##
## Attaching package: 'randomForest'
## The following object is masked from 'package:ggplot2':
##
##     margin
require("e1071")

## Loading required package: e1071
#main file
cc <- data.frame(read.csv("data/cc.csv"))
cc <- cc[,c(2:31)]

```

```

#check balance
fraud <- nrow(cc[cc$Class == 1,])
notFraud <- nrow(cc) - fraud
paste("fraud: ", fraud, "|| not fraud: ", notFraud)

## [1] "fraud: 492 || not fraud: 284315"

#make the size smaller for easier use
set.seed(65)
cc_simp <- cc[sample(nrow(cc), 25000), ]
fraud <- nrow(cc_simp[cc_simp$Class == 1,])
notFraud <- nrow(cc_simp) - fraud
paste("fraud: ", fraud, "|| not fraud: ", notFraud)

## [1] "fraud: 43 || not fraud: 24957"

#split data for testing models
trainLength <- floor(.7*nrow(cc_simp))
testLength <- nrow(cc_simp) - trainLength

train_model <- cc_simp[1:trainLength,]
train_eval <- cc_simp[(trainLength + 1):nrow(cc_simp),]

#oversample
fraud <- nrow(train_model[train_model$Class == 1,])
notFraud <- nrow(train_model) - fraud
paste("fraud: ", fraud, "|| not fraud: ", notFraud)

## [1] "fraud: 29 || not fraud: 17471"

train_model <- ovun.sample(Class ~ ., data = train_model, method = "over", N = notFraud*2, seed = 1)$data

#functions for sd and se
mysd <- function(predict, target) {
  diff_sq <- (predict - mean(target))^2
  return(mean(sqrt(diff_sq)))
}

myse <- function(predict, target) {
  diff_sq <- (predict - target)^2
  return(mean(sqrt(diff_sq)))
}

#Model1 - Multiple Linear Regression - Base Line
mlr1 <- glm(Class~., data = train_model)
BIC(mlr1)

## [1] 5837.002

predict_ml1 <- predict(mlr1, train_eval, type = 'response')
table(train_eval$Class, predict_ml1 > 0.5)

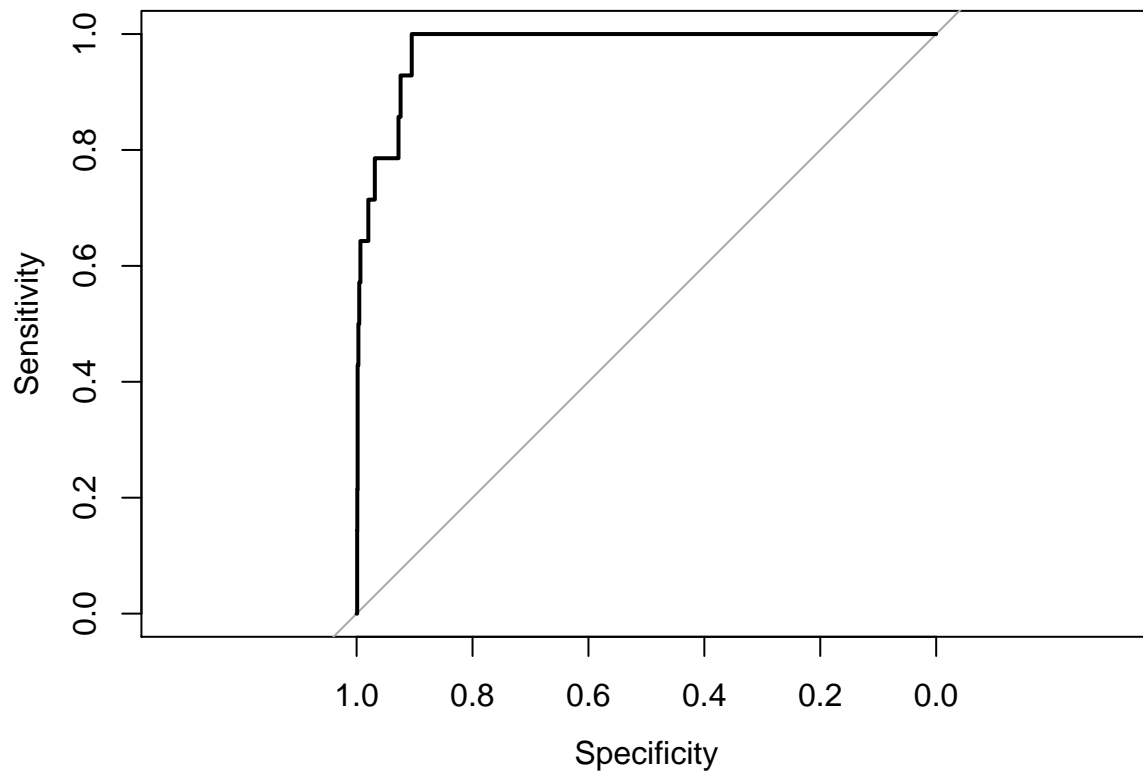
##
## FALSE TRUE
## 0 7278 208

```

```
##      1      4     10
mysd(predict_mlr1, train_eval$Class)

## [1] 0.1943988
myse(predict_mlr1, train_eval$Class)

## [1] 0.1945416
auc_mlr1 <- roc(train_eval$Class, predict_mlr1)
plot(auc_mlr1)
```



```
##
## Call:
## roc.default(response = train_eval$Class, predictor = predict_mlr1)
##
## Data: predict_mlr1 in 7486 controls (train_eval$Class 0) < 14 cases (train_eval$Class 1).
## Area under the curve: 0.9773
#Model2 - Poisson Model
poisson1 <- glm(Class ~ ., family = "poisson", data = train_model)

## Warning: glm.fit: fitted rates numerically 0 occurred
BIC(poison1)

## [1] 43690.02
```

```
predict_poisson1 <- predict(poisson1, train_eval, type = 'response')
table(train_eval$Class, predict_poisson1 > 0.5)
```

```
##
##      FALSE TRUE
##  0  7274  212
##  1     5    9
```

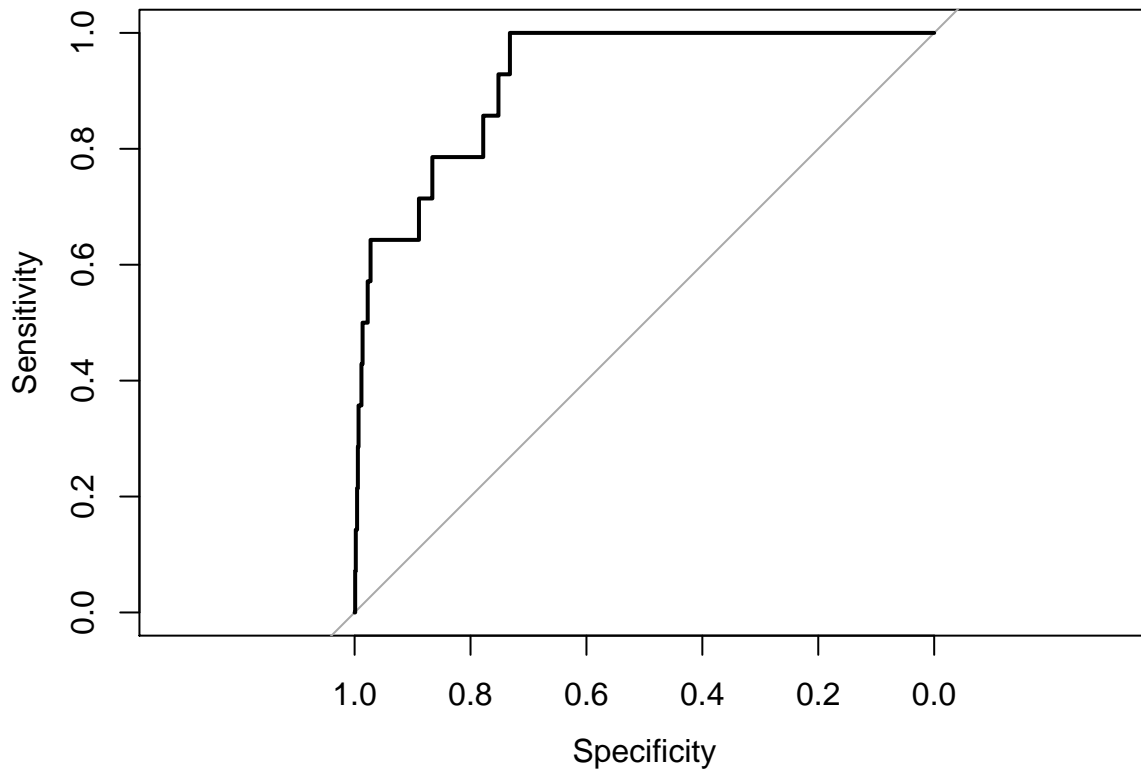
```
mysd(predict_poisson1, train_eval$Class)
```

```
## [1] 0.2760893
```

```
myse(predict_poisson1, train_eval$Class)
```

```
## [1] 0.2767727
```

```
auc_poisson <- roc(train_eval$Class, predict_poisson1)
plot(auc_poisson)
```



```
##
## Call:
## roc.default(response = train_eval$Class, predictor = predict_poisson1)
##
## Data: predict_poisson1 in 7486 controls (train_eval$Class 0) < 14 cases (train_eval$Class 1).
## Area under the curve: 0.923
```

```
#logit model
```

```
logit1 <- glm(Class ~., family = binomial(link='logit'), data = train_model)
```

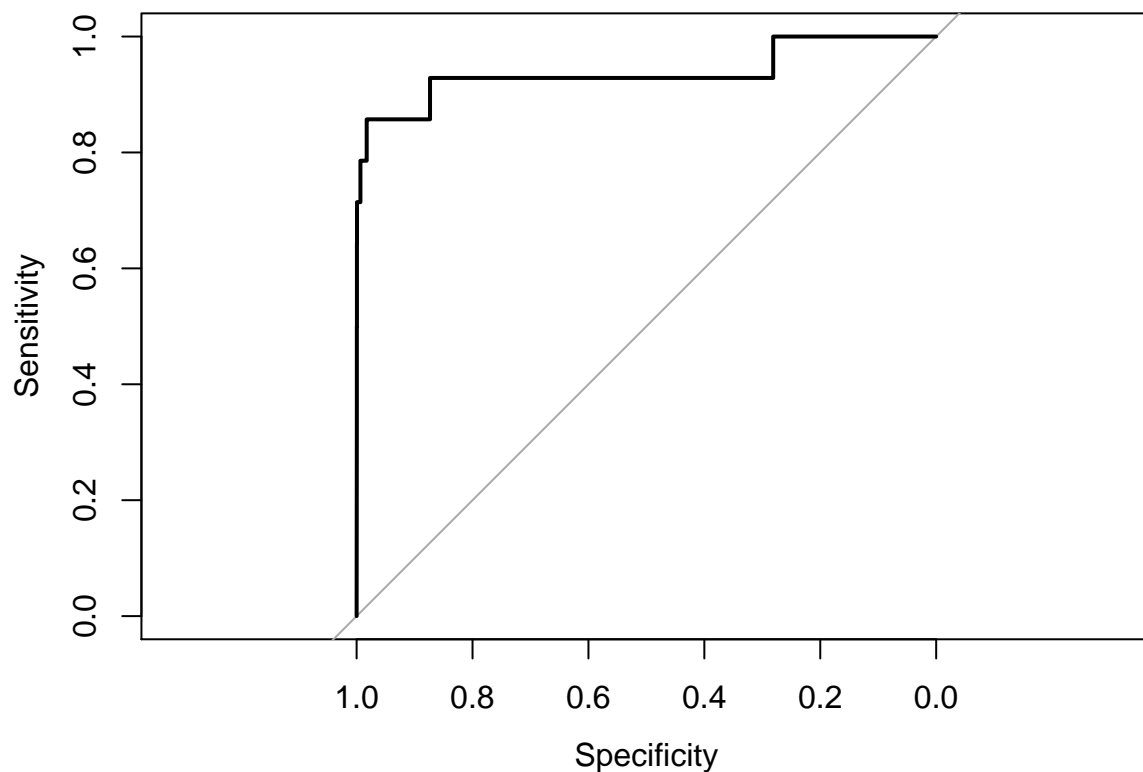
```
## Warning: glm.fit: algorithm did not converge
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
BIC(logit1)
```

```
## [1] 3954.386
```

```
predict_logit1 <- predict(logit1, train_eval, type = 'response')
table(train_eval$Class, predict_logit1 > 0.5)
```

```
##
##      FALSE TRUE
## 0  7347  139
## 1     2   12
```

```
auc_logit1 <- roc(train_eval$Class, predict_logit1)
plot(auc_logit1)
```



```
##
## Call:
## roc.default(response = train_eval$Class, predictor = predict_logit1)
##
## Data: predict_logit1 in 7486 controls (train_eval$Class 0) < 14 cases (train_eval$Class 1).
## Area under the curve: 0.9377
# backward stepwise
stepwisel <- glm(Class ~ ., data = train_model)
backward <- step(stepwisel, trace = 0)
```

```
BIC(backward)
```

```
## [1] 5826.854
```

```
predict_backward <- predict(backward, train_eval, type = 'response')  
table(train_eval$Class, predict_backward > 0.5)
```

```
##
```

```
##      FALSE TRUE
```

```
##    0  7279  207
```

```
##    1     4   10
```

```
mysd(predict_backward, train_eval$Class)
```

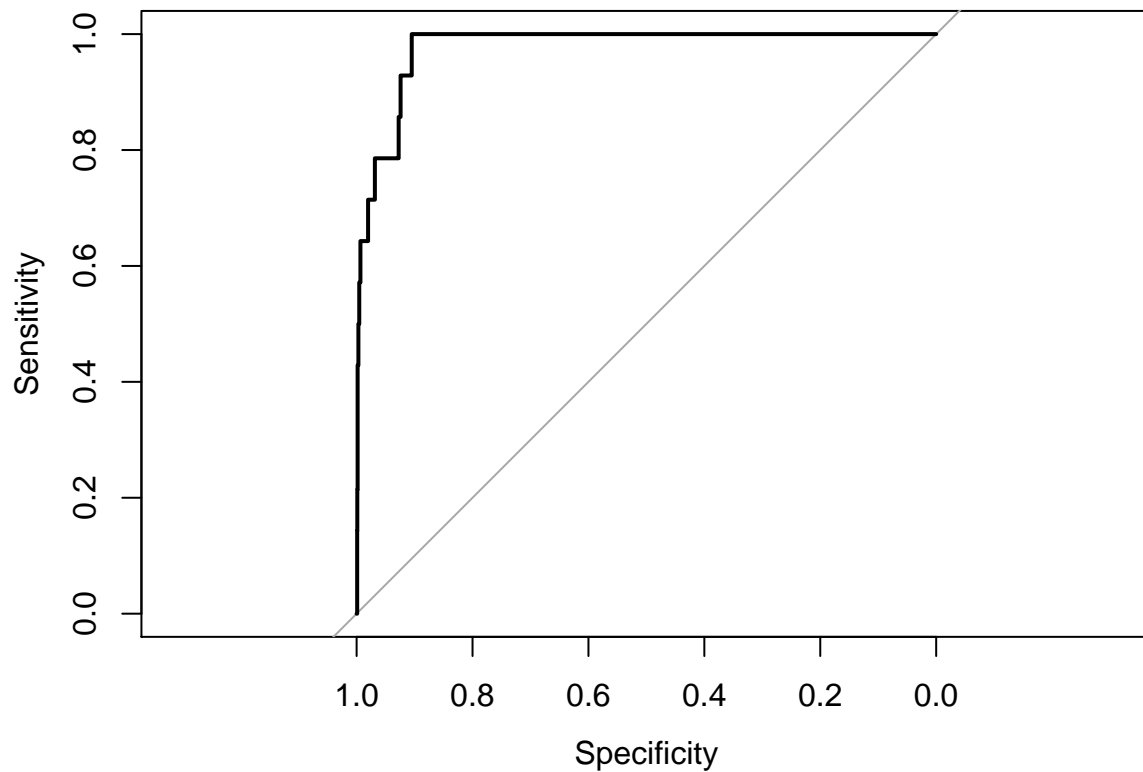
```
## [1] 0.1943553
```

```
myse(predict_backward, train_eval$Class)
```

```
## [1] 0.1945002
```

```
auc_backward <- roc(train_eval$Class, predict_backward)
```

```
plot(auc_backward)
```



```
##
```

```
## Call:
```

```
## roc.default(response = train_eval$Class, predictor = predict_backward)
```

```
##
```

```
## Data: predict_backward in 7486 controls (train_eval$Class 0) < 14 cases (train_eval$Class 1).
```

```
## Area under the curve: 0.9773
```

```

#forward stepwise
stepwise2 <- glm(Class ~ 1,data = train_model)
forward <- step(stepwise2, scope = list(lower=formula(stepwise2), upper=formula(stepwise1)), direction = "BIC",
BIC(forward))

## [1] 5826.854

predict_forward <- predict(forward, train_eval, type = 'response')
table(train_eval$Class, predict_forward > 0.5)

##
##      FALSE TRUE
##    0  7279  207
##    1     4   10

mysd(predict_forward, train_eval$Class)

## [1] 0.1943553

myse(predict_forward, train_eval$Class)

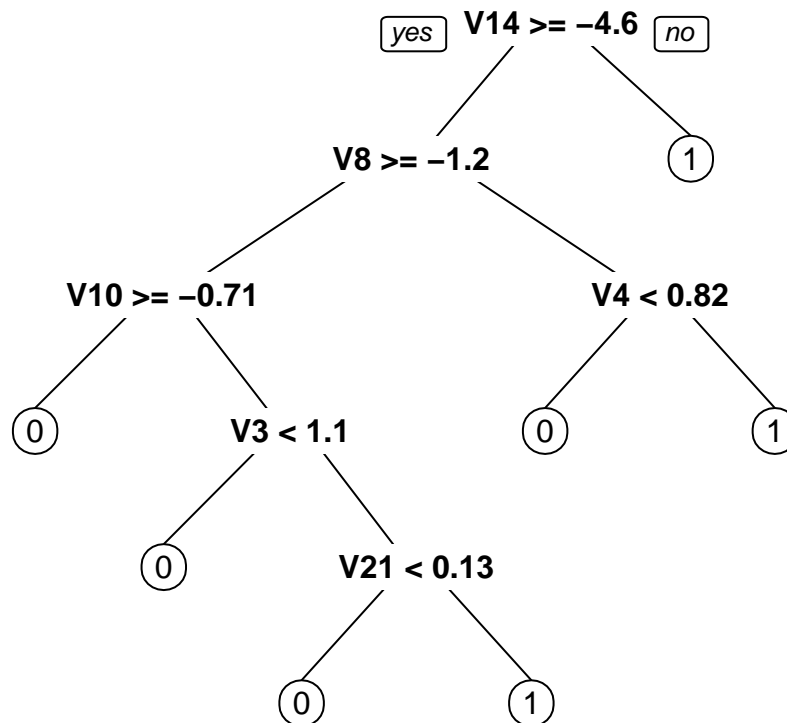
## [1] 0.1945002

auc_forward <- roc(train_eval$Class, predict_forward)
plot(auc_forward)

##
## Call:
## roc.default(response = train_eval$Class, predictor = predict_forward)
##
## Data: predict_forward in 7486 controls (train_eval$Class 0) < 14 cases (train_eval$Class 1).
## Area under the curve: 0.9773

# decision tree
decision <- rpart(Class ~ ., data = train_model, method = "class")
prp(decision)

```



```
predict_decision <- predict(decision, train_eval, type = "class")
confusionMatrix(train_eval$Class, predict_decision)
```

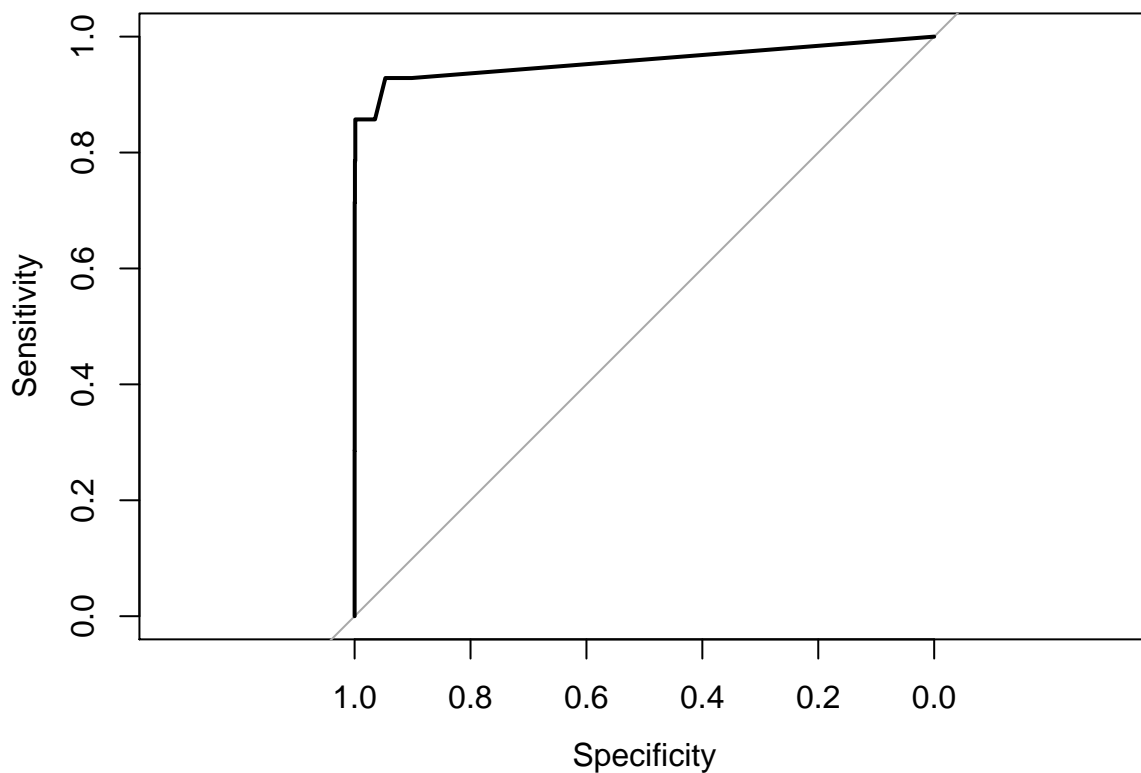
```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##           0 7292 194
##           1    3   11
##
##           Accuracy : 0.9737
##           95% CI : (0.9699, 0.9772)
##           No Information Rate : 0.9727
##           P-Value [Acc > NIR] : 0.3005
##
##           Kappa : 0.0973
##           Mcnemar's Test P-Value : <2e-16
##
##           Sensitivity : 0.99959
##           Specificity : 0.05366
##           Pos Pred Value : 0.97408
##           Neg Pred Value : 0.78571
##           Prevalence : 0.97267
##           Detection Rate : 0.97227
##           Detection Prevalence : 0.99813
##           Balanced Accuracy : 0.52662
```



```
##
##      'Positive' Class : 0
##
#decision tree random forest (kind of broke with raw data)
rforest <- randomForest(Class ~ ., data = train_model)

## Warning in randomForest.default(m, y, ...): The response has five or fewer
## unique values. Are you sure you want to do regression?
predict_rforest <- predict(rforest, train_eval)
table(train_eval$Class, predict_rforest > 0.5)

##
##      FALSE TRUE
##      0 7485    1
##      1     6    8
auc_rforest <- roc(train_eval$Class, predict_rforest)
plot(auc_rforest)
```

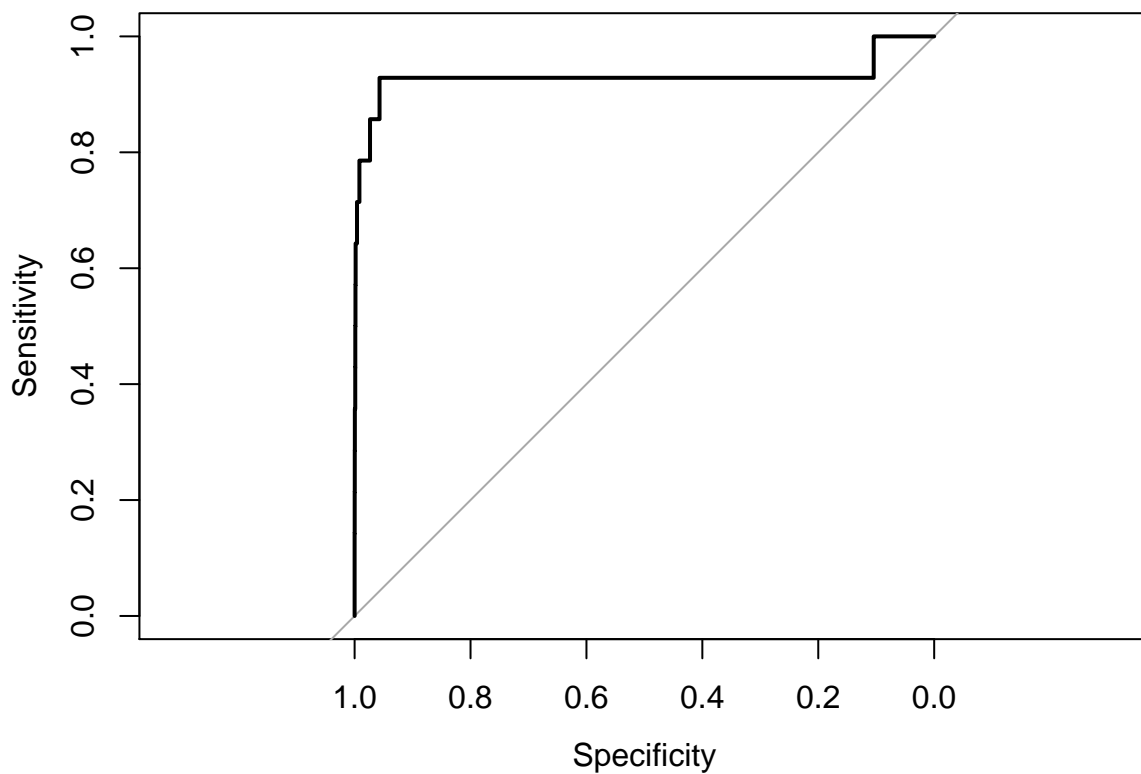


```
##
## Call:
## roc.default(response = train_eval$Class, predictor = predict_rforest)
##
## Data: predict_rforest in 7486 controls (train_eval$Class 0) < 14 cases (train_eval$Class 1).
## Area under the curve: 0.9574
```

```
#svm (kind of broken with raw data)
svm <- svm(Class ~ ., data = train_model)
predict_svm <- predict(svm, train_eval)
table(train_eval$Class, predict_svm > 0.5)
```

```
##
##      FALSE TRUE
##  0  7460   26
##  1     5    9
```

```
auc_svm <- roc(train_eval$Class, predict_svm)
plot(auc_svm)
```



```
##
## Call:
## roc.default(response = train_eval$Class, predictor = predict_svm)
##
## Data: predict_svm in 7486 controls (train_eval$Class 0) < 14 cases (train_eval$Class 1).
## Area under the curve: 0.9297
```

```
#tables only
table(train_eval$Class, predict_mlr1 > 0.5)
```

```
##
##      FALSE TRUE
##  0  7278   208
##  1     4    10
```

```
table(train_eval$Class, predict_poisson1 > 0.5)
```

```
##
##      FALSE TRUE
##  0  7274  212
##  1     5    9
```

```
table(train_eval$Class, predict_logit1 > 0.5)
```

```
##
##      FALSE TRUE
##  0  7347  139
##  1     2   12
```

```
table(train_eval$Class, predict_backward > 0.5)
```

```
##
##      FALSE TRUE
##  0  7279  207
##  1     4   10
```

```
table(train_eval$Class, predict_forward > 0.5)
```

```
##
##      FALSE TRUE
##  0  7279  207
##  1     4   10
```

```
confusionMatrix(train_eval$Class, predict_decision)
```

```
## Confusion Matrix and Statistics
```

```
##
##              Reference
## Prediction    0    1
##              0 7292 194
##              1    3   11
##
##              Accuracy : 0.9737
##              95% CI : (0.9699, 0.9772)
##              No Information Rate : 0.9727
##              P-Value [Acc > NIR] : 0.3005
##
##              Kappa : 0.0973
##              Mcnemar's Test P-Value : <2e-16
##
##              Sensitivity : 0.99959
##              Specificity : 0.05366
##              Pos Pred Value : 0.97408
##              Neg Pred Value : 0.78571
##              Prevalence : 0.97267
##              Detection Rate : 0.97227
##              Detection Prevalence : 0.99813
##              Balanced Accuracy : 0.52662
##
##              'Positive' Class : 0
##
```

```
table(train_eval$Class, predict_rforest > 0.5)
```

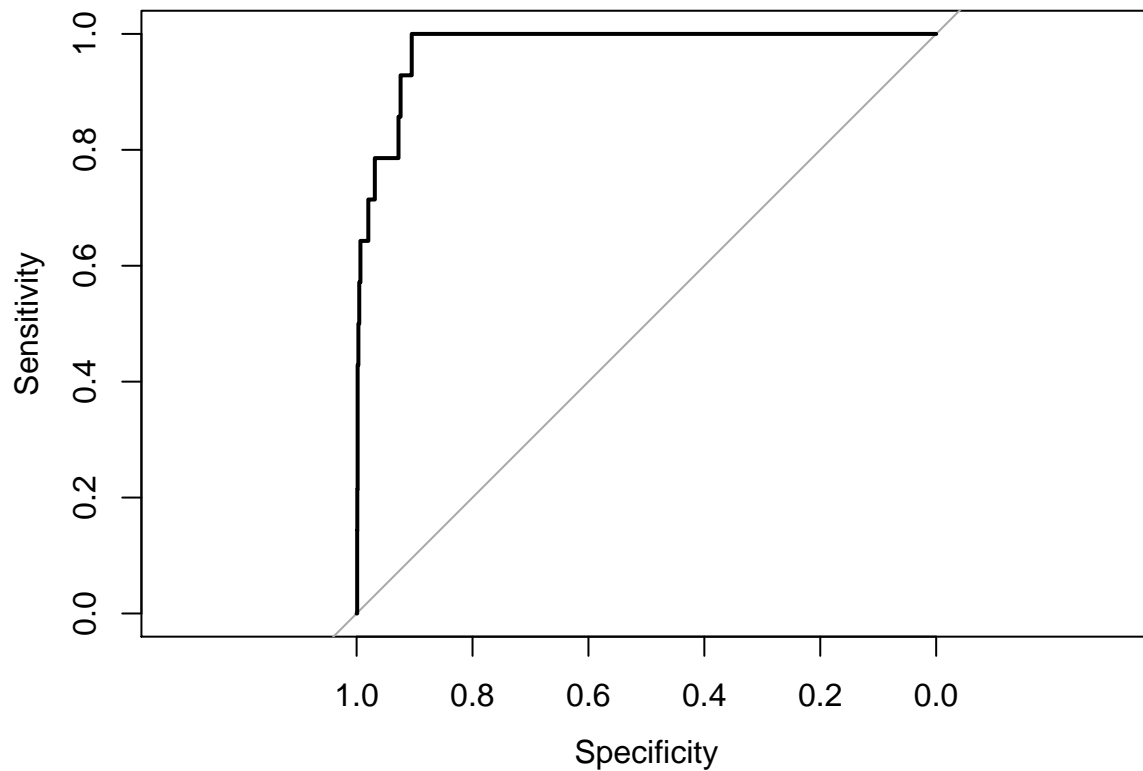
```
##  
##      FALSE TRUE  
##    0  7485    1  
##    1     6    8
```

```
table(train_eval$Class, predict_svm > 0.5)
```

```
##  
##      FALSE TRUE  
##    0  7460    26  
##    1     5     9
```

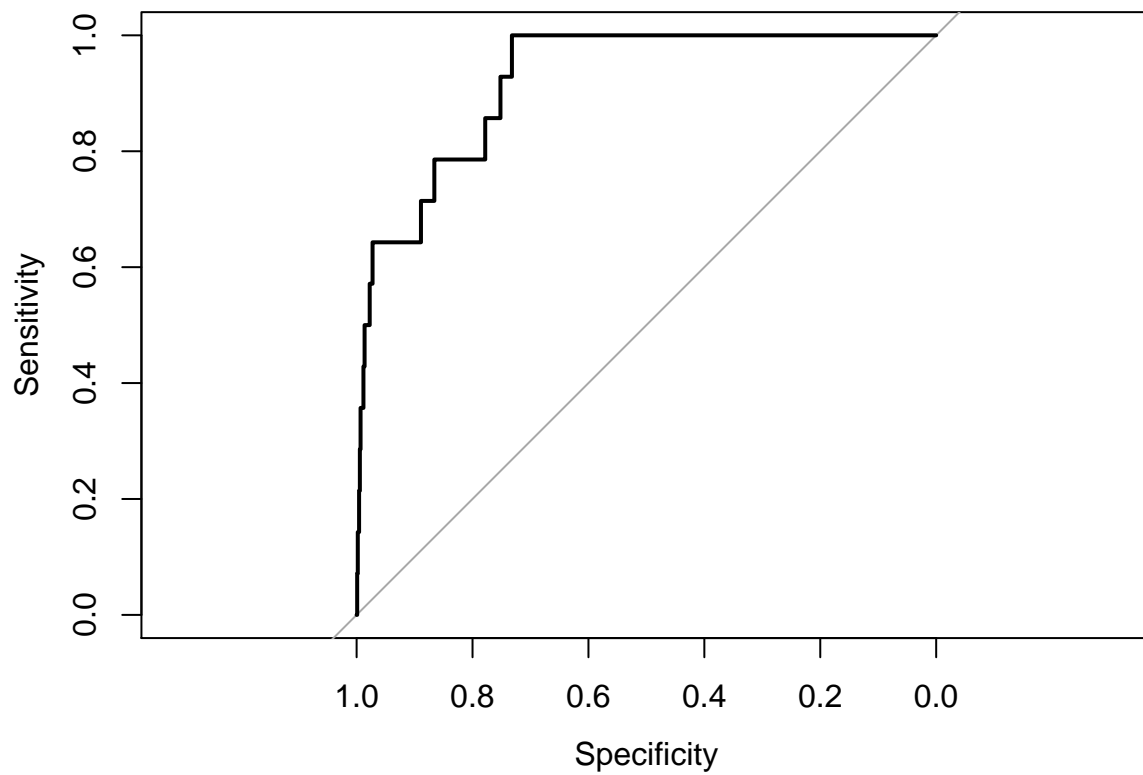
```
#AUC plots only
```

```
plot(auc_mlr1)
```

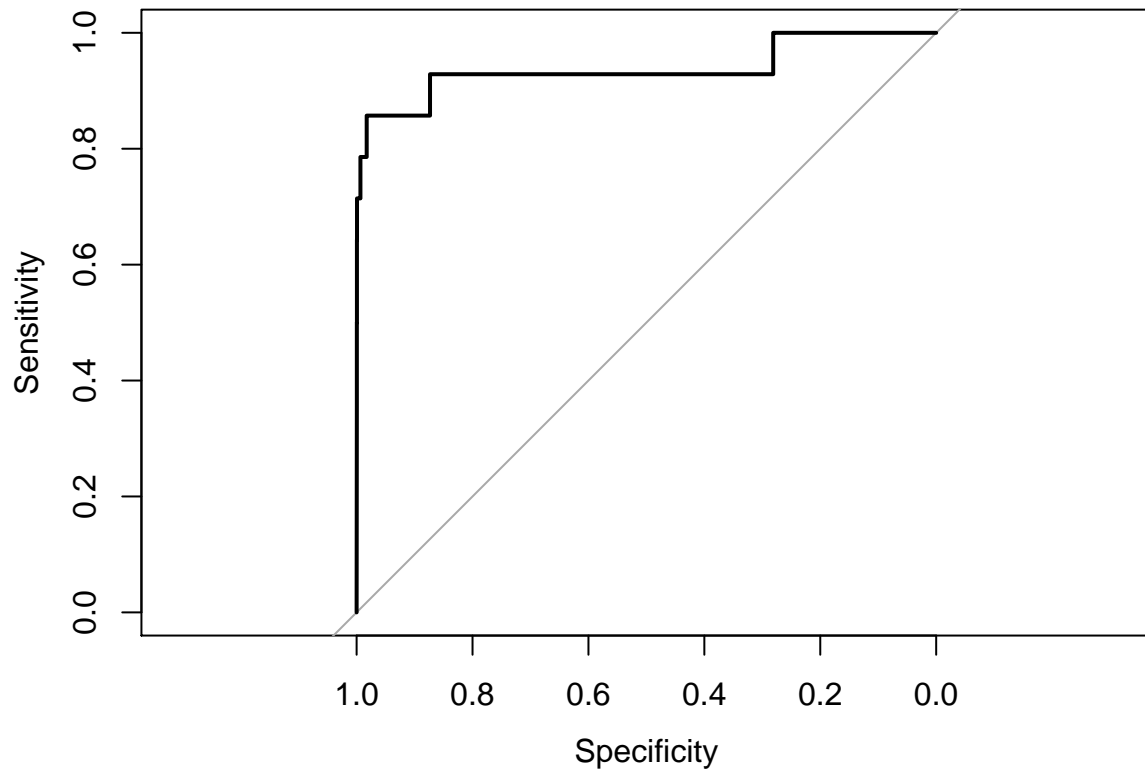


```
##  
## Call:  
## roc.default(response = train_eval$Class, predictor = predict_mlr1)  
##  
## Data: predict_mlr1 in 7486 controls (train_eval$Class 0) < 14 cases (train_eval$Class 1).  
## Area under the curve: 0.9773
```

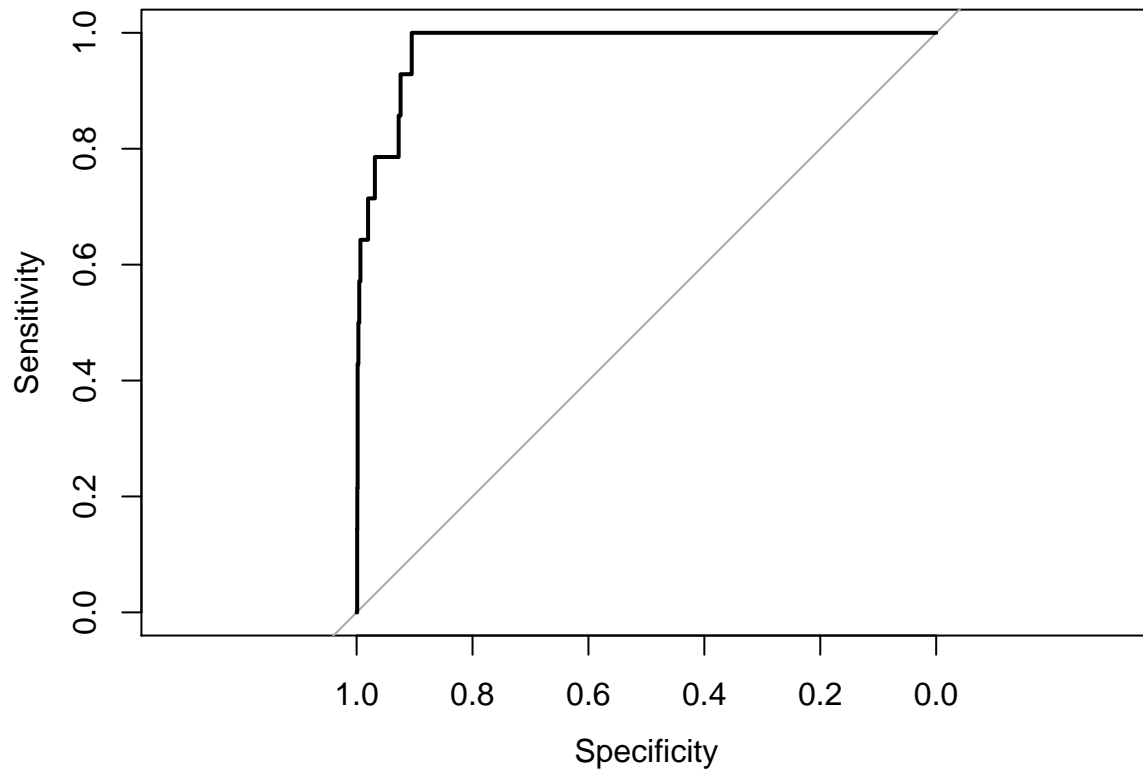
```
plot(auc_poisson)
```



```
##  
## Call:  
## roc.default(response = train_eval$Class, predictor = predict_poisson1)  
##  
## Data: predict_poisson1 in 7486 controls (train_eval$Class 0) < 14 cases (train_eval$Class 1).  
## Area under the curve: 0.923  
plot(auc_logit1)
```

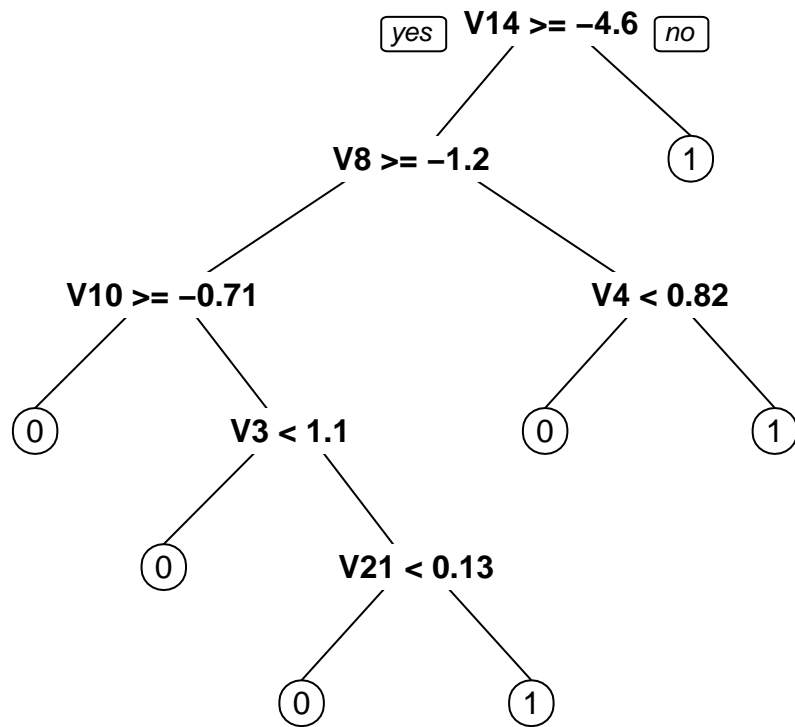


```
##  
## Call:  
## roc.default(response = train_eval$Class, predictor = predict_logit1)  
##  
## Data: predict_logit1 in 7486 controls (train_eval$Class 0) < 14 cases (train_eval$Class 1).  
## Area under the curve: 0.9377  
plot(auc_backward)
```

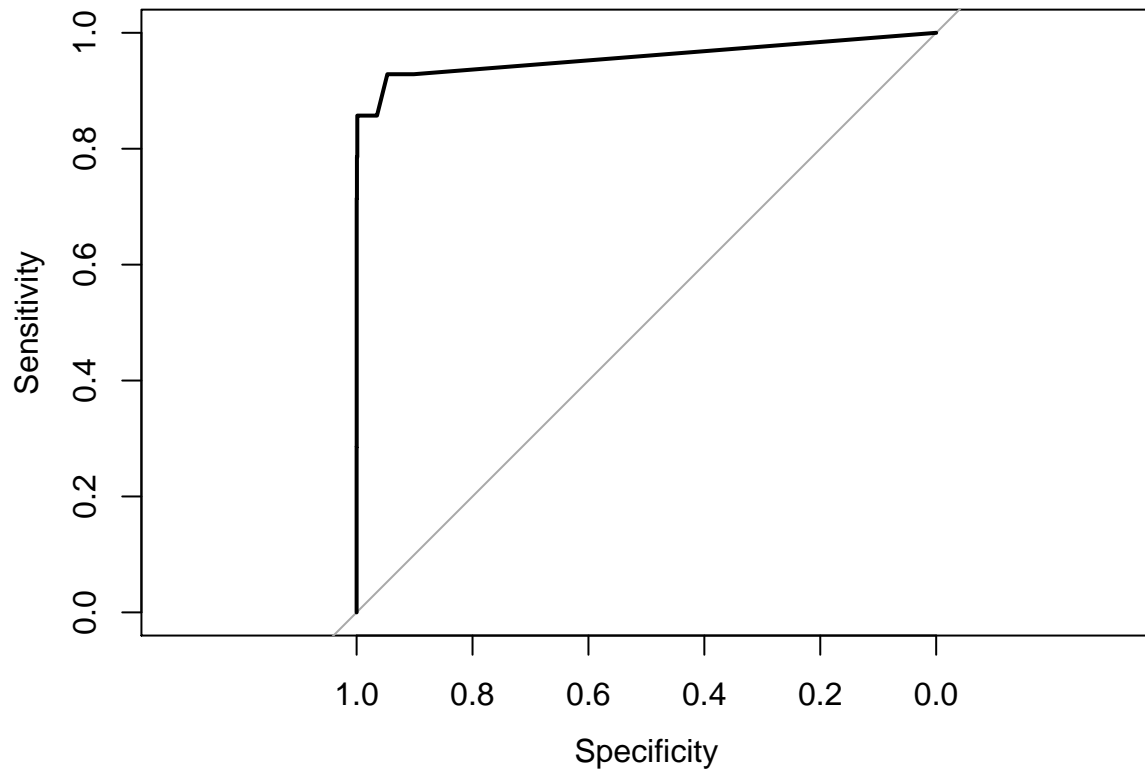


```
##
## Call:
## roc.default(response = train_eval$Class, predictor = predict_backward)
##
## Data: predict_backward in 7486 controls (train_eval$Class 0) < 14 cases (train_eval$Class 1).
## Area under the curve: 0.9773
plot(auc_forward)

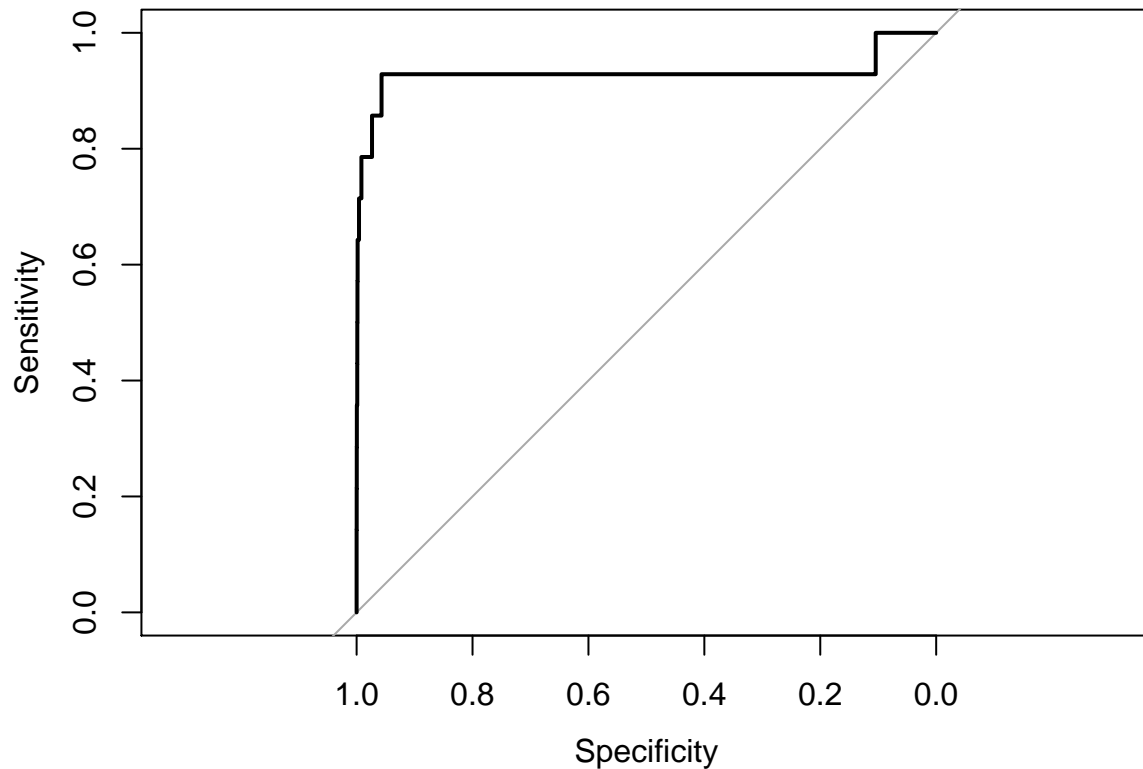
##
## Call:
## roc.default(response = train_eval$Class, predictor = predict_forward)
##
## Data: predict_forward in 7486 controls (train_eval$Class 0) < 14 cases (train_eval$Class 1).
## Area under the curve: 0.9773
prp(decision)
```



```
plot(auc_rforest)
```

```
##  
## Call:  
## roc.default(response = train_eval$Class, predictor = predict_rforest)  
##  
## Data: predict_rforest in 7486 controls (train_eval$Class 0) < 14 cases (train_eval$Class 1).  
## Area under the curve: 0.9574  
plot(auc_svm)
```



```
##
## Call:
## roc.default(response = train_eval$Class, predictor = predict_svm)
##
## Data: predict_svm in 7486 controls (train_eval$Class 0) < 14 cases (train_eval$Class 1).
## Area under the curve: 0.9297
accuracy.meas(train_eval$Class, predict_mlr1)

##
## Call:
## accuracy.meas(response = train_eval$Class, predicted = predict_mlr1)
##
## Examples are labelled as positive when predicted is greater than 0.5
##
## precision: 0.046
## recall: 0.714
## F: 0.043
accuracy.meas(train_eval$Class, predict_poisson1)

##
## Call:
## accuracy.meas(response = train_eval$Class, predicted = predict_poisson1)
##
## Examples are labelled as positive when predicted is greater than 0.5
##
## precision: 0.041
```

```

## recall: 0.643
## F: 0.038
accuracy.meas(train_eval$Class, predict_logit1)

##
## Call:
## accuracy.meas(response = train_eval$Class, predicted = predict_logit1)
##
## Examples are labelled as positive when predicted is greater than 0.5
##
## precision: 0.079
## recall: 0.857
## F: 0.073
accuracy.meas(train_eval$Class, predict_backward)

##
## Call:
## accuracy.meas(response = train_eval$Class, predicted = predict_backward)
##
## Examples are labelled as positive when predicted is greater than 0.5
##
## precision: 0.046
## recall: 0.714
## F: 0.043
accuracy.meas(train_eval$Class, predict_forward)

##
## Call:
## accuracy.meas(response = train_eval$Class, predicted = predict_forward)
##
## Examples are labelled as positive when predicted is greater than 0.5
##
## precision: 0.046
## recall: 0.714
## F: 0.043
accuracy.meas(train_eval$Class, predict_decision)

##
## Call:
## accuracy.meas(response = train_eval$Class, predicted = predict_decision)
##
## Examples are labelled as positive when predicted is greater than 0.5
##
## precision: 0.002
## recall: 1.000
## F: 0.002
accuracy.meas(train_eval$Class, predict_rforest)

##
## Call:
## accuracy.meas(response = train_eval$Class, predicted = predict_rforest)
##
## Examples are labelled as positive when predicted is greater than 0.5

```

```
##
## precision: 0.889
## recall: 0.571
## F: 0.348
accuracy.meas(train_eval$Class, predict_svm)

##
## Call:
## accuracy.meas(response = train_eval$Class, predicted = predict_svm)
##
## Examples are labelled as positive when predicted is greater than 0.5
##
## precision: 0.257
## recall: 0.643
## F: 0.184
```