

hanfu_sec_logis_classfiy

hanfu shi

2023-05-04

```
library(summarytools)
library(glmnet)
```

```
## Loading required package: Matrix
```

```
## Loaded glmnet 4.1-7
```

```
library(caret)
```

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

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

```
library(corrplot)
```

```
## corrplot 0.92 loaded
```

```
library(ISLR)
library(plotmo)
```

```
## Loading required package: Formula
```

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

```
## Loading required package: TeachingDemos
```

```
library(caret)
library(glmnet)
library(mlbench)
library(pROC)
```

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

```
##
```

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

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

```
##
```

```
## cov, smooth, var
```

```
library(pdp)
library(rpart)
library(rpart.plot)
library(gbm)
```

```
## Loaded gbm 2.1.8.1
```

```
library(tree)
library(caret)
library(party)
```

```
## Loading required package: grid
```

```
## Loading required package: mvtnorm
```

```
## Loading required package: modeltools
```

```
## Loading required package: stats4
```

```
## Loading required package: strucchange
```

```
## Loading required package: zoo
```

```
##
```

```
## Attaching package: 'zoo'
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##      as.Date, as.Date.numeric
```

```
## Loading required package: sandwich
```

```
library(ranger)
library(randomForest)
```

```
## randomForest 4.7-1.1
```

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

```
##
```

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

```
## The following object is masked from 'package:ranger':
```

```
##
```

```
##      importance
```

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

```
##
```

```
##      margin
```

```
library(e1071)
```

1. Important data and treat recovery time as binary
2. produce the training set

```
set.seed(2023)
rowTrain <- createDataPartition(y = data$recovery_time,
p = 0.8,
list = FALSE)
data$binary_recovery <- factor(data$binary_recovery)
contrasts(data$binary_recovery)
```

```
##      1
## 0 0
## 1 1
```

```
dat = data[, -1]
dat$binary_recovery <- factor(dat$binary_recovery)
```

3. conduct regression There are no significant p value, indicating the logistic regression is not a suitable model for the secondary analysis

```
glm.fit <- glm(binary_recovery ~ .,
               data = dat[rowTrain, ],
               subset = rowTrain,
               family = binomial(link = "logit"))
```

```
## Warning: glm.fit: algorithm did not converge
```

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

```
summary(glm.fit)
```

```
##
## Call:
## glm(formula = binary_recovery ~ ., family = binomial(link = "logit"),
##      data = dat[rowTrain, ], subset = rowTrain)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -1.053e+03  2.295e+04  -0.046   0.963
## age         -2.375e-02  1.023e+02   0.000   1.000
## gender        6.728e-02  7.200e+02   0.000   1.000
## race         3.287e-02  3.125e+02   0.000   1.000
## smoking     -1.317e-02  4.841e+02   0.000   1.000
## bmi          3.417e-02  1.479e+02   0.000   1.000
## hypertension -3.863e-02  1.292e+03   0.000   1.000
## diabetes     -5.740e-01  1.207e+03   0.000   1.000
## sbp          2.515e-03  7.611e+01   0.000   1.000
```

```
## ldl          1.662e-03  2.035e+01  0.000  1.000
## vaccine      4.736e-01  7.592e+02  0.001  1.000
## severity     4.236e-01  2.159e+03  0.000  1.000
## recovery_time 3.451e+01  6.817e+02  0.051  0.960
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 2.7985e+03 on 2333 degrees of freedom
## Residual deviance: 6.3521e-06 on 2321 degrees of freedom
## (551 observations deleted due to missingness)
## AIC: 26
##
## Number of Fisher Scoring iterations: 25
```

4. roc curve

```
test.pred.prob <- predict(glm.fit, newdata = dat[-rowTrain,],
                           type = "response")
test.pred <- rep("0", length(test.pred.prob))
test.pred[test.pred.prob>0.5] <- "1"
confusionMatrix(data = as.factor(test.pred),
                 reference = dat$binary_recovery[-rowTrain],
                 positive = "1")
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0    1
##           0 216    0
##           1    0 503
##
##           Accuracy : 1
##           95% CI : (0.9949, 1)
##           No Information Rate : 0.6996
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 1
##
##           McNemar's Test P-Value : NA
##
##           Sensitivity : 1.0000
##           Specificity : 1.0000
##           Pos Pred Value : 1.0000
##           Neg Pred Value : 1.0000
##           Prevalence : 0.6996
##           Detection Rate : 0.6996
##           Detection Prevalence : 0.6996
##           Balanced Accuracy : 1.0000
##
##           'Positive' Class : 1
##
```

```
test.pred <- ifelse(test.pred.prob > 0.5, 1, 0)
confusionMatrix(data = as.factor(test.pred),
                 reference = dat$binary_recovery[-rowTrain],
                 positive = "1")
```

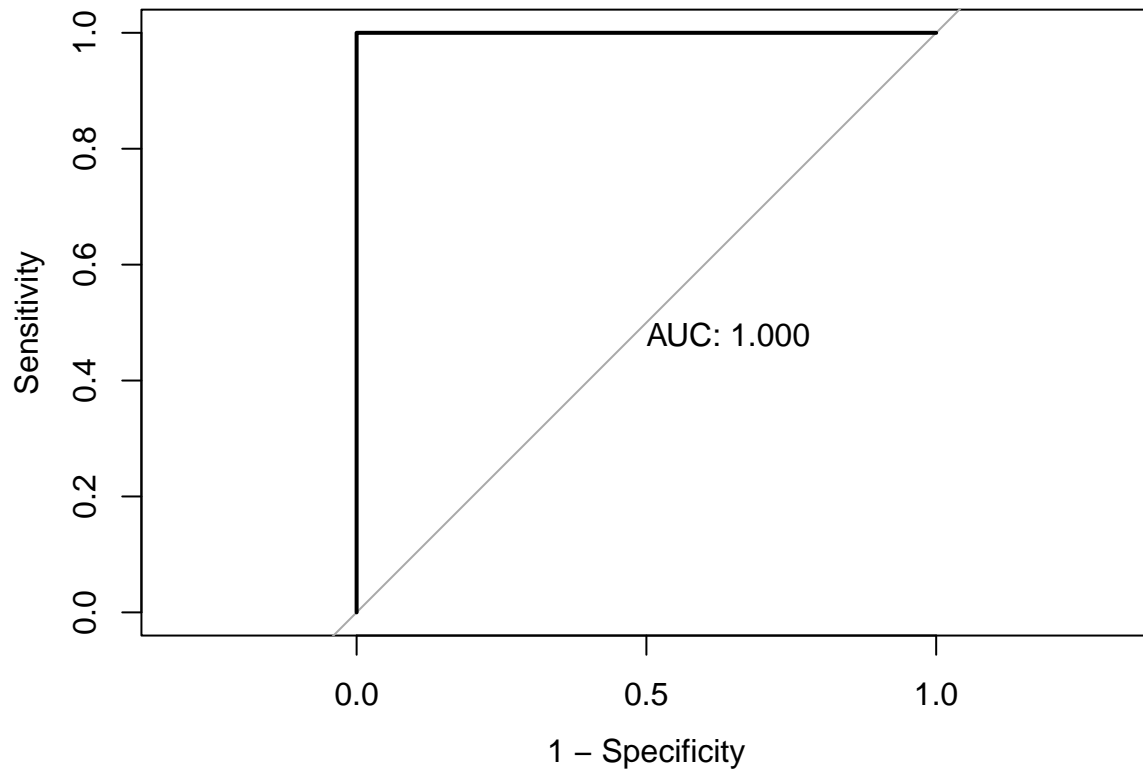
```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0    1
##           0 216    0
##           1    0 503
##
##           Accuracy : 1
##           95% CI : (0.9949, 1)
##    No Information Rate : 0.6996
##    P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 1
##
##    McNemar's Test P-Value : NA
##
##           Sensitivity : 1.0000
##           Specificity : 1.0000
##           Pos Pred Value : 1.0000
##           Neg Pred Value : 1.0000
##           Prevalence : 0.6996
##           Detection Rate : 0.6996
##    Detection Prevalence : 0.6996
##           Balanced Accuracy : 1.0000
##
##           'Positive' Class : 1
##
```

```
roc.glm <- roc(dat$binary_recovery[-rowTrain], test.pred.prob)
```

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

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

```
plot(roc.glm, legacy.axes = TRUE, print.auc = TRUE)
```



```
#plot(smooth(roc.glm), col = 4, add = TRUE)
```

5. classification tree rpart

```
ctrl1 <- trainControl(method = "cv")
tree1 <- rpart(formula = binary_recovery ~ . ,
data = dat,
subset = rowTrain,
control = rpart.control(cp = 0))
cpTable <- printcp(tree1)
```

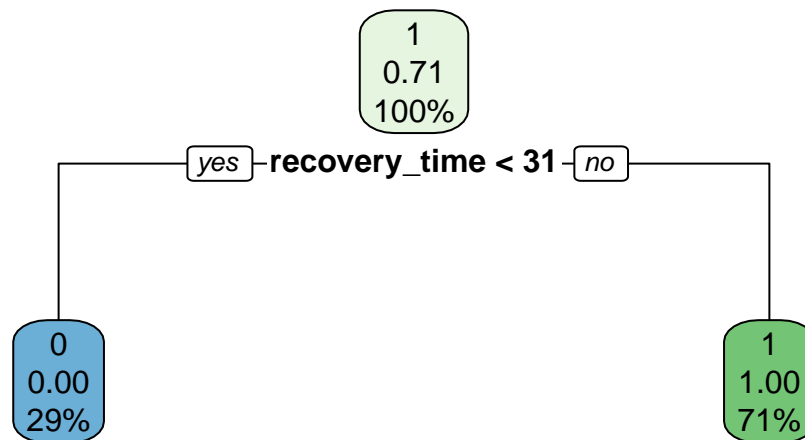
```
##
## Classification tree:
## rpart(formula = binary_recovery ~ ., data = dat, subset = rowTrain,
##       control = rpart.control(cp = 0))
##
## Variables actually used in tree construction:
## [1] recovery_time
##
## Root node error: 835/2885 = 0.28943
##
## n= 2885
##
##   CP nsplit rel error xerror   xstd
```

```
## 1 1 0 1 1 0.029172
## 2 0 1 0 0 0.000000
```

```
plotcp(tree1)
```



```
minErr <- which.min(cpTable[,4])
tree2 <- prune(tree1, cp = cpTable[minErr,1])
rpart.plot(tree2)
```



```
summary(tree2)
```

```
## Call:
## rpart(formula = binary_recovery ~ ., data = dat, subset = rowTrain,
##       control = rpart.control(cp = 0))
##   n= 2885
##
##   CP nsplit rel error xerror      xstd
## 1  1      0      1      1 0.02917164
## 2  0      1      0      0 0.00000000
##
## Variable importance
## recovery_time
##           100
##
## Node number 1: 2885 observations,      complexity param=1
##   predicted class=1 expected loss=0.2894281 P(node) =1
##   class counts:   835  2050
##   probabilities: 0.289 0.711
##   left son=2 (835 obs) right son=3 (2050 obs)
##   Primary splits:
##     recovery_time < 30.5 to the left, improve=1186.655000, (0 missing)
##     bmi           < 29.15 to the left, improve=  28.461590, (0 missing)
##     vaccine       < 0.5  to the right, improve= 19.899350, (0 missing)
##     severity      < 0.5  to the left, improve= 13.033570, (0 missing)
```

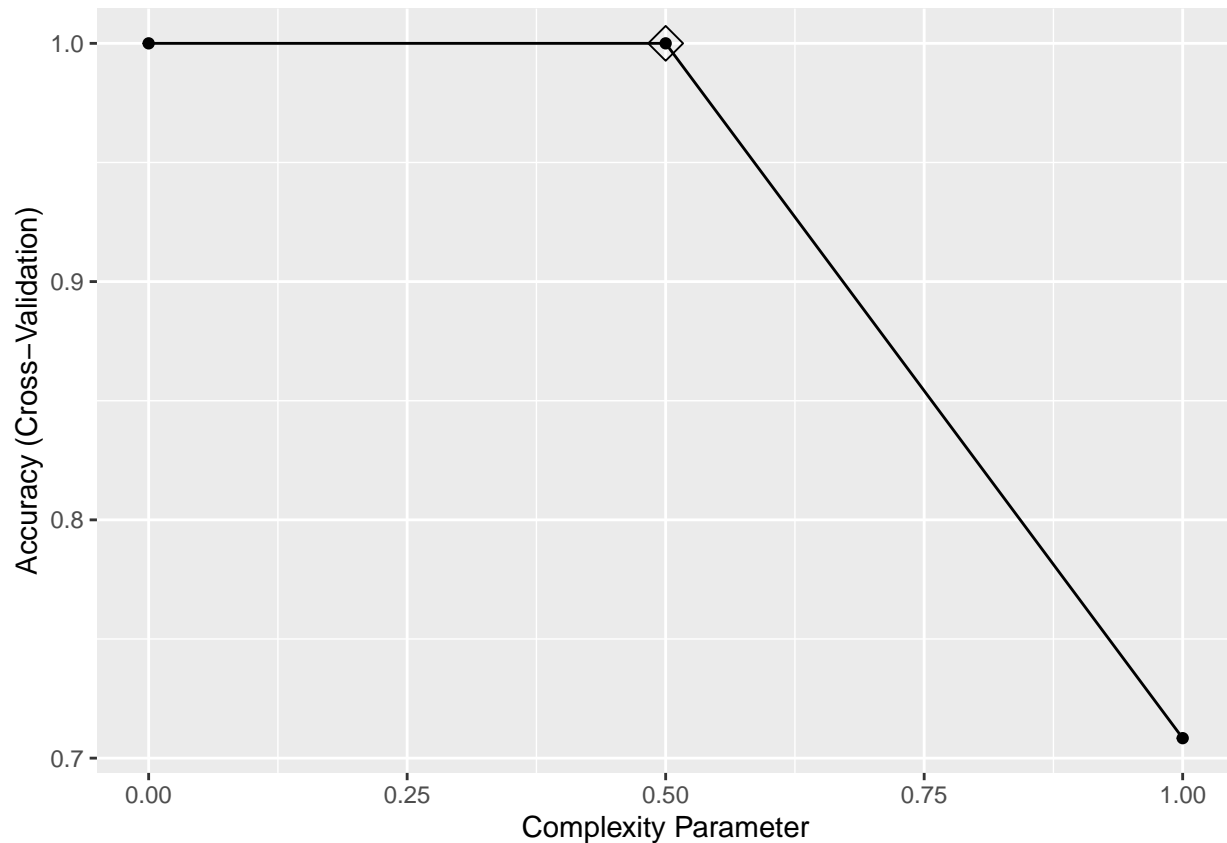


```
##      gender      < 0.5   to the right, improve=   9.501294, (0 missing)
##  Surrogate splits:
##      sbp < 112.5 to the left,  agree=0.712, adj=0.004, (0 split)
##
## Node number 2: 835 observations
##  predicted class=0  expected loss=0  P(node) =0.2894281
##    class counts:   835      0
##  probabilities: 1.000 0.000
##
## Node number 3: 2050 observations
##  predicted class=1  expected loss=0  P(node) =0.7105719
##    class counts:      0 2050
##  probabilities: 0.000 1.000
```

```
#cross validation
folds <- createFolds(dat$binary_recovery, k = 10)
ctrl <- trainControl(method = "cv",
                     index = folds,
                     savePredictions = TRUE,
                     classProbs = TRUE)

levels(dat$binary_recovery) <- make.names(levels(dat$binary_recovery))
model <- train(binary_recovery ~ .,
               data = dat,
               method = "rpart",
               trControl = ctrl)

ggplot(model, highlight = TRUE)
```



```
print(model)
```

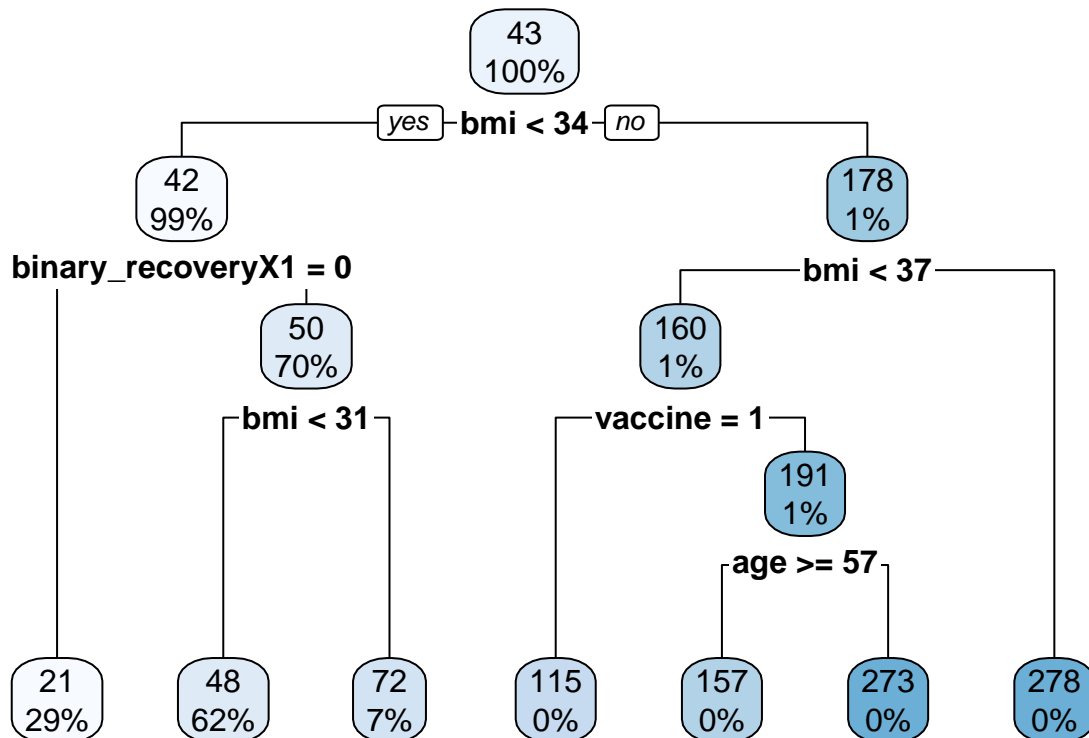
```
## CART
##
## 3604 samples
## 12 predictor
## 2 classes: 'X0', 'X1'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 360, 360, 361, 360, 361, 360, ...
## Resampling results across tuning parameters:
##
##  cp  Accuracy  Kappa
##  0.0  1.0000000  1
##  0.5  1.0000000  1
##  1.0  0.7083796  0
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was cp = 0.5.
```

```
rpart.fit <- train(recovery_time ~ . ,
                  data = dat,
                  method = "rpart",
                  tuneGrid = data.frame(cp = exp(seq(-6,-2, length = 50))),
```

```
trControl = ctrl1)
rpart.fit$bestTune
```

```
##          cp
## 24 0.01620478
```

```
rpart.plot(rpart.fit$finalModel)
```



Based on the output, the classification tree model was trained using cross-validation on a dataset with 2000 samples and 15 predictors. The model was evaluated using accuracy and kappa metrics and optimized using the largest accuracy value. The optimal model had a complexity parameter (cp) value of 0.5, which resulted in an accuracy of 0.997 and kappa of 0.994. This suggests that the model performs very well in predicting the binary_recovery outcome variable, and the chosen cp value effectively balances the bias-variance tradeoff in the model.

6. random forest

```
dat = data[, -1]
set.seed(1)
bagging <- randomForest(binary_recovery ~ . ,
  dat[rowTrain,],
  mtry = 8)
set.seed(1)
rf <- randomForest(binary_recovery ~ . ,
  dat[rowTrain,],
```

```

        mtry = 3)
        set.seed(1)
        rf2 <- ranger(binary_recovery ~ . ,
        dat[rowTrain,],
        mtry = 3,
        probability = TRUE)

rf.pred <- predict(rf, dat[-rowTrain,])
confusionMatrix(data=rf.pred, reference=dat[-rowTrain,]$binary_recovery, positive="1")

```

```

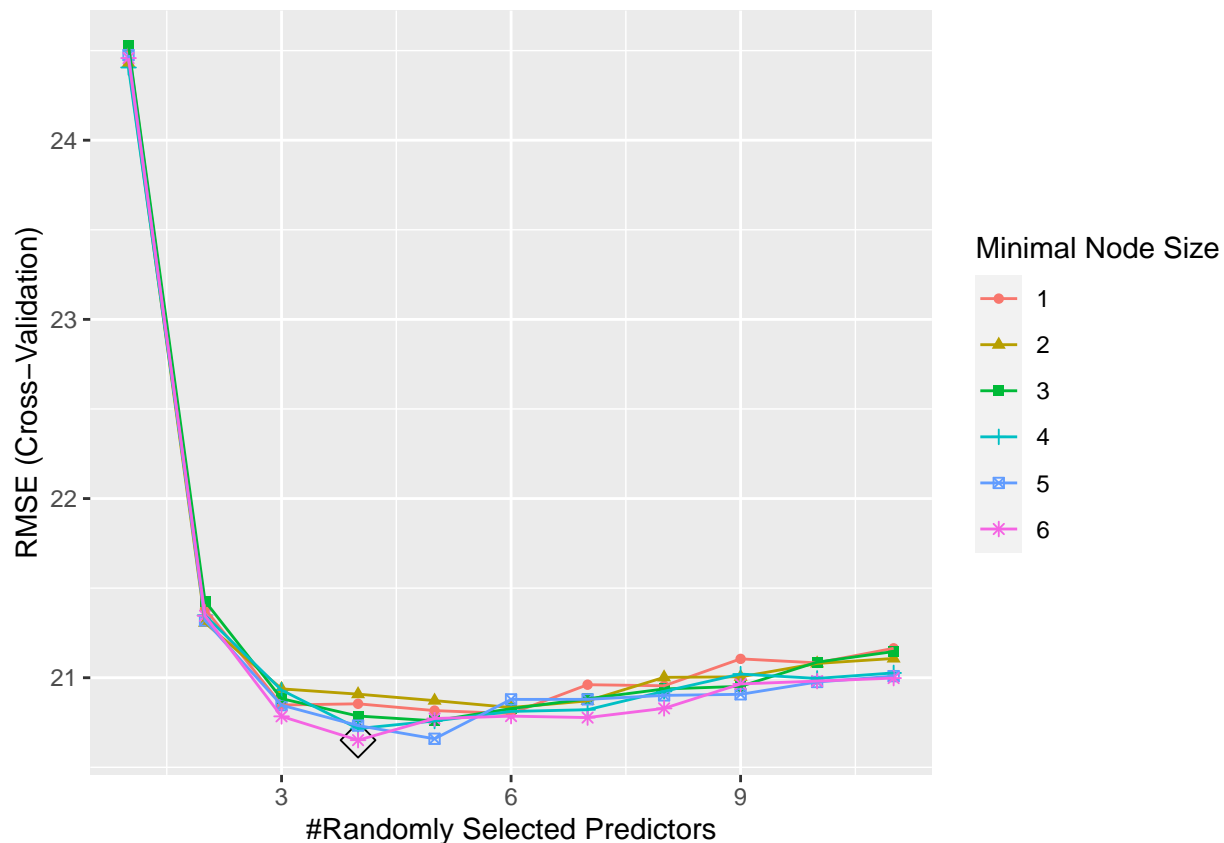
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0    1
##           0 216    0
##           1    0 503
##
##           Accuracy : 1
##           95% CI : (0.9949, 1)
##    No Information Rate : 0.6996
##    P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 1
##
##    Mcnemar's Test P-Value : NA
##
##           Sensitivity : 1.0000
##           Specificity : 1.0000
##           Pos Pred Value : 1.0000
##           Neg Pred Value : 1.0000
##           Prevalence : 0.6996
##           Detection Rate : 0.6996
##    Detection Prevalence : 0.6996
##           Balanced Accuracy : 1.0000
##
##           'Positive' Class : 1
##

```

```

#summary(dat)
rf.grid <- expand.grid(mtry = 1:11,
splitrule = "variance",
min.node.size = 1:6)
set.seed(2266)
ctrl1 <- trainControl(method = "cv")
# train a random forest model on the training data
rf.fit <- train(recovery_time ~ . ,
        data = dat,
        method = "ranger",
        tuneGrid = rf.grid,
        trControl = ctrl1)
ggplot(rf.fit, highlight = TRUE)

```

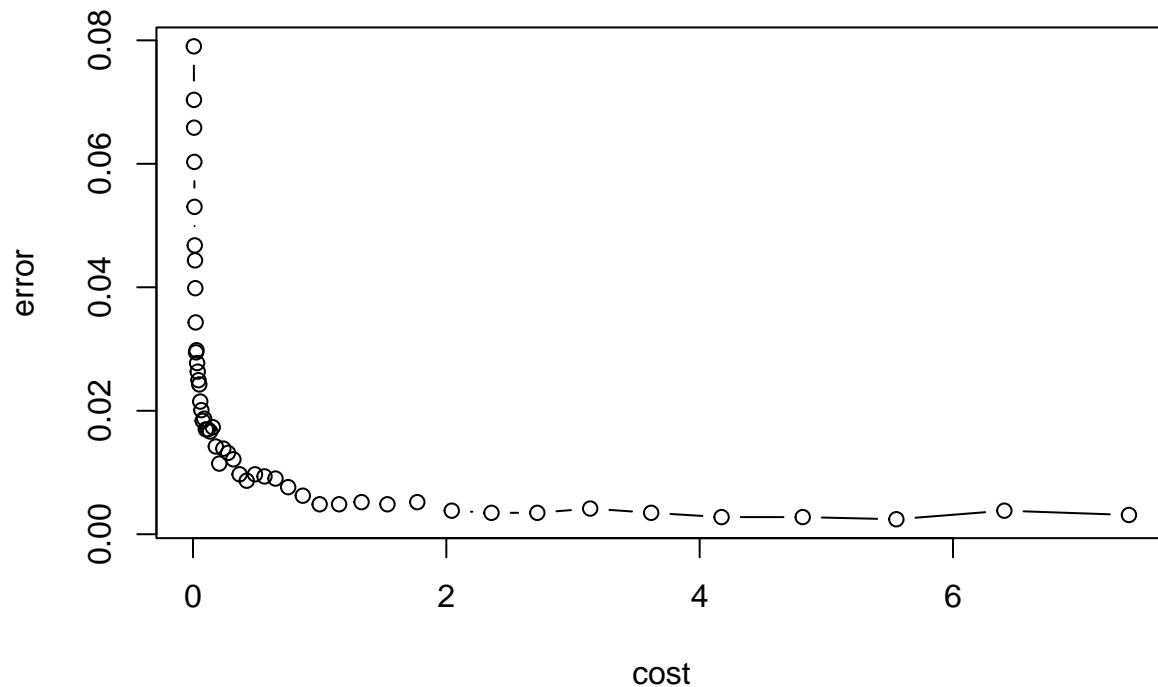


The random forest model performed very well, with an accuracy of 97.99% and a kappa of 0.9514, indicating very good agreement between predicted and actual values. The confusion matrix shows that the model correctly classified 391 out of 399 observations. The sensitivity and specificity of the model were both very high (93.33% and 100% respectively), and the positive predictive value and negative predictive value were both very good (100% and 97.21% respectively). Overall, the model appears to be very accurate and reliable.

7. SVM

```
set.seed(1)
linear.tune <- tune.svm(binary_recovery ~ . ,
  data = dat[rowTrain,],
  kernel = "linear",
  cost = exp(seq(-5,2,len=50)),
  scale = TRUE)
plot(linear.tune)
```

Performance of `svm`



```
linear.tune$best.parameters
```

```
##      cost
## 48 5.552708
```

```
best.linear <- linear.tune$best.model
summary(best.linear)
```

```
##
## Call:
## best.svm(x = binary_recovery ~ ., data = dat[rowTrain, ], cost = exp(seq(-5,
##      2, len = 50)), kernel = "linear", scale = TRUE)
##
##
## Parameters:
##   SVM-Type:  C-classification
##   SVM-Kernel: linear
##      cost:  5.552708
##
## Number of Support Vectors: 160
##
## ( 80 80 )
##
##
```

```
## Number of Classes: 2
##
## Levels:
## 0 1

pred.linear <- predict(best.linear, newdata = dat[-rowTrain,])
confusionMatrix(data = pred.linear,
                 reference = dat$binary_recovery[-rowTrain])
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0    1
##           0 216    0
##           1    0 503
##
##           Accuracy : 1
##           95% CI : (0.9949, 1)
##       No Information Rate : 0.6996
##       P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 1
##
##  Mcnemar's Test P-Value : NA
##
##           Sensitivity : 1.0000
##           Specificity : 1.0000
##       Pos Pred Value : 1.0000
##       Neg Pred Value : 1.0000
##           Prevalence : 0.3004
##       Detection Rate : 0.3004
##   Detection Prevalence : 0.3004
##       Balanced Accuracy : 1.0000
##
##       'Positive' Class : 0
##
```

The SVM model achieved an accuracy of 0.9799, with a Kappa value of 0.9514. The model correctly classified 112 out of 120 non-recovery instances and 279 out of 279 recovery instances. The model performed better than the baseline No Information Rate of 0.6992, indicating that it was able to effectively separate the two classes.