# ActSc 632 Assignment 2 Solutions

**Spring 2023**

**Department of Statistics and Actuarial Science, University of Waterloo**

## Question 1

Determine the number of predictors in this data set, and whether they are quantitative of qualita- tive. How many observations are there? How many observations are classified as a "bad credit"? "good credit"?

```
library(tree)
library(CASdatasets)
```

```
## Loading required package: xts
```

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

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

```
## The following objects are masked from 'package:base':
##
##     as.Date, as.Date.numeric
```

```
## Loading required package: sp
```

```
## The legacy packages maptools, rgdal, and rgeos, underpinning the sp package,
## which was just loaded, will retire in October 2023.
## Please refer to R-spatial evolution reports for details, especially
## https://r-spatial.org/r/2023/05/15/evolution4.html.
## It may be desirable to make the sf package available;
## package maintainers should consider adding sf to Suggests:.
## The sp package is now running under evolution status 2
##     (status 2 uses the sf package in place of rgdal)
```

```
library(pROC)
```

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

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

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

```
library(randomForest)
```

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

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

```
# Load the data
data(credit)

within(credit,{
installment_rate <-factor(installment_rate)
residence_since <- factor(residence_since)
existing_credits <- factor(existing_credits)
num_dependents <- factor(num_dependents)
})
```

```
# Summary of the data
summary(credit)
```

```
##  checking_status    duration     credit_history    purpose     credit_amount
##  A11:274         Min.   : 4.0   A30: 40       A43    :280   Min.   :  250
##  A12:269         1st Qu.:12.0   A31: 49       A40    :234   1st Qu.: 1366
##  A13: 63         Median :18.0   A32:530       A42    :181   Median : 2320
##  A14:394         Mean   :20.9   A33: 88       A41    :103   Mean   : 3271
##                  3rd Qu.:24.0   A34:293       A49    : 97   3rd Qu.: 3972
##                  Max.   :72.0                 A46    : 50   Max.   :18424
##                                               (Other): 55
##  savings    employment installment_rate personal_status other_parties
##  A61:603   A71: 62   Min.   :1.000    A91: 50       A101:907
##  A62:103   A72:172   1st Qu.:2.000    A92:310       A102: 41
##  A63: 63   A73:339   Median :3.000    A93:548       A103: 52
##  A64: 48   A74:174   Mean   :2.973    A94: 92
##  A65:183   A75:253   3rd Qu.:4.000
##                      Max.   :4.000
##
##  residence_since property_magnitude     age       other_payment_plans
##  Min.   :1.000   A121:282           Min.   :19.00   A141:139
##  1st Qu.:2.000   A122:232           1st Qu.:27.00   A142: 47
##  Median :3.000   A123:332           Median :33.00   A143:814
##  Mean   :2.845   A124:154           Mean   :35.55
##  3rd Qu.:4.000                      3rd Qu.:42.00
##  Max.   :4.000                      Max.   :75.00
##
##  housing    existing_credits    job      num_dependents  telephone
##  A151:179   Min.   :1.000   A171: 22   Min.   :1.000   A191:596
##  A152:713   1st Qu.:1.000   A172:200   1st Qu.:1.000   A192:404
##  A153:108   Median :1.000   A173:630   Median :1.000
##             Mean   :1.407   A174:148   Mean   :1.155
##             3rd Qu.:2.000              3rd Qu.:1.000
##             Max.   :4.000              Max.   :2.000
##
##  foreign_worker     class
##  A201:963      Min.   :0.0
##  A202: 37      1st Qu.:0.0
##                Median :0.0
##                Mean   :0.3
##                3rd Qu.:1.0
##                Max.   :1.0
##
```

```
print(c(sum(credit$class == 1),
        sum(credit$class == 0)))
```

```
## [1] 300 700
```

There are 20 predictors. Duration, credit amount, age are quantitative, and the other are qualitative. Depending on the encoding used for the data set (which differs from source to source) the predictors installment rate, residence since, existing credits and num dependents may be encoded as integers and thus be a priori be considered quantitative. But when looking at the meaning of the different values these predictors can take, it seems best to treat them as categorical. It won't matter in our analysis though, as these 4 predictors do not end up being considered in our models. and the others are categorical. There are 300 observations out of 1000 classified as bad credit.

# Question 2

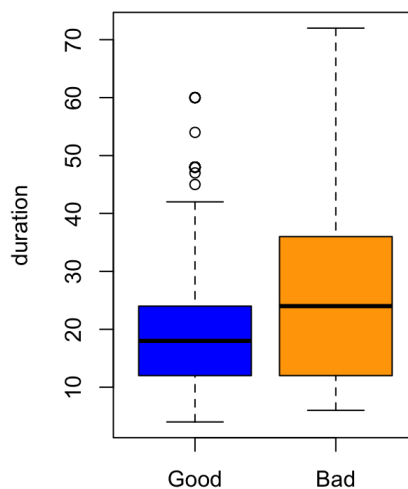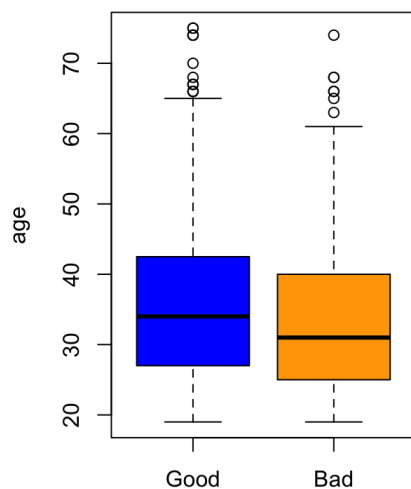To explore the data further, produce the following plots:

- for each of the predictors age and duration, make a box plot showing the distribution of the observations, separately for the "good" and "bad" observations;
- for the pairs duration & savings and duration & credit history, plot the observations (using duration on the x axis) and use different symbols for the "good" and "bad" observations.

Comments on the plots you obtained.

```
par(mfrow=c(1,2))

# Create a box plot of 'age' for the two classes, displayed inline
boxplot(credit$age[credit$class==0], credit$age[credit$class==1],
        names=c("Good","Bad"), ylab = "age", col=c("blue","orange"))

# Create a box plot of 'duration' for the two classes, displayed inline
boxplot(credit$duration[credit$class==0], credit$duration[credit$class==1],
        names=c("Good","Bad"), ylab = "duration", col=c("blue","orange"))
```
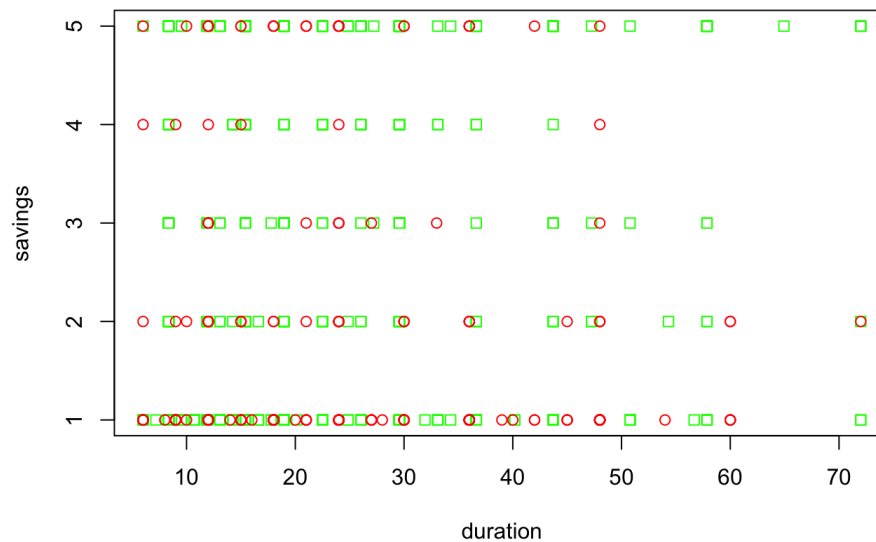
Clearly the bad credit observations

seems to have a longer duration, and the distribution of the duration has a larger variance; they also seem to be slightly younger compared to the good credit.

```
# Set up the plotting area
par(mfrow=c(1,1))

# Subset 'savings' and 'duration' for the two classes
csG<-credit$savings[credit$class==0]
csB<-credit$savings[credit$class==1]
dG<-credit$duration[credit$class==0]
dB<-credit$duration[credit$class==1]

# Plot 'duration' vs 'savings' for the two classes, displayed inline
plot(dG, csG, col=c("green"), pch=0, axes = F, ylab="", xlab="")
par(new=T)
plot(dB, csB, col=c("red"), pch=1, xlab="duration", ylab="savings", axes=T)
```
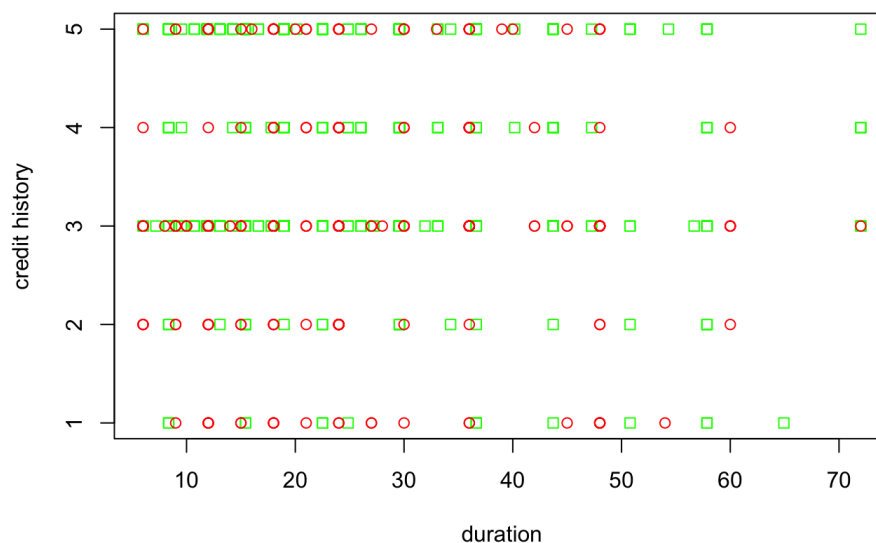
```
# Set up the plotting area
par(mfrow=c(1,1))

# Subset 'credit_history' for the two classes
chG<-credit$credit_history[credit$class==0]
chB<-credit$credit_history[credit$class==1]

# Plot 'duration' vs 'credit_history' for the two classes, displayed inline
plot(dG, chG, col=c("green"), pch=0, axes = F, ylab="", xlab="")
par(new=T)
plot(dB, chB, col=c("red"), pch=1, xlab="duration", ylab="credit history", axes=T)
```



It is harder to see a clear trend in these two graphs, in that the good and bad credits are not easily classified according to duration for a given level of credit history or savings.

# Question 3

Randomly split your data in 70% of the observations for training and 30% for testing.

```
set.seed(1)
train <- sample(1:nrow(credit), 0.7 * nrow(credit))

credit$class <- as.factor(credit$class)
```

# Question 4

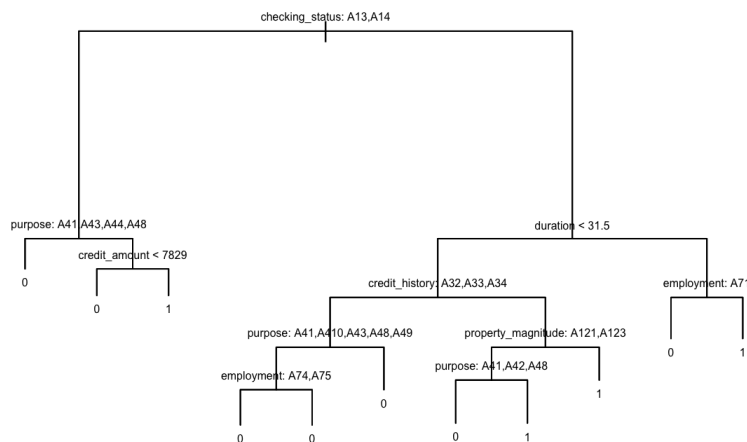Simply using recursive binary partitioning, obtain a tree for this classifcation problem.

- How many leaves does your tree have?
- How many factors were used to build this tree?
- What is the deviance for this tree? (If you used something else than the default definition of deviance in R, please specify how is deviance determined).
- Plot the tree you obtained using R. There should be enough information that given an observation, one could determine in which leaf it ends up.
- Use your tree to make predictions for the test data set. Produce the confusion matrix corresponding to your tree and plot the ROC curve.

```
tree.credit<-tree(class~.,data=credit,subset=train)
summary(tree.credit)
```

```
##
## Classification tree:
## tree(formula = class ~ ., data = credit, subset = train)
## Variables actually used in tree construction:
## [1] "checking_status"    "purpose"              "credit_amount"
## [4] "duration"           "credit_history"       "employment"
## [7] "property_magnitude"
## Number of terminal nodes:  11
## Residual mean deviance:  0.9321 = 642.2 / 689
## Misclassification error rate: 0.2143 = 150 / 700
```

So there are 11 leaves and 7 factors were used. The deviance is 0.9321. Note that since the tree will depend on the training set, which is randomly chosen, you may have obtained a tree that has a quite different structure, with a different subset of factors used, and a different number of terminal nodes. As mentioned in class, this method has a high variance, which is why the results can be so different.

```
plot(tree.credit)
text(tree.credit,pretty=0,cex=0.5)
```



We recall that the categories listed on a node are those used to determine which observations go in the left child.

```
pred.tree.cred<-predict(tree.credit,newdata=credit[-train,],type="class")

credit.test<-credit$class[-train]
tree.tab<-table(pred.tree.cred,credit.test)

tree.tab
```

```
##               credit.test
## pred.tree.cred   0    1
##              0 188   64
##              1  20   28
```

```
print(c("Overall error:", (tree.tab[1,2]+tree.tab[2,1])/sum(tree.tab[,])))
```

```
## [1] "Overall error:" "0.28"
```

```
print(c("Type I error:", tree.tab[2,1]/sum(tree.tab[,1])))
```

```
## [1] "Type I error:"       "0.0961538461538462"
```

```
print(c("Type II error:", tree.tab[1,2]/sum(tree.tab[,2])))
```

```
## [1] "Type II error:"     "0.695652173913043"
```

```
tree.dev.pred <- predict(tree.credit,newdata=credit[-train,],type="tree")

print(c("Deviance:", deviance(tree.dev.pred)))
```
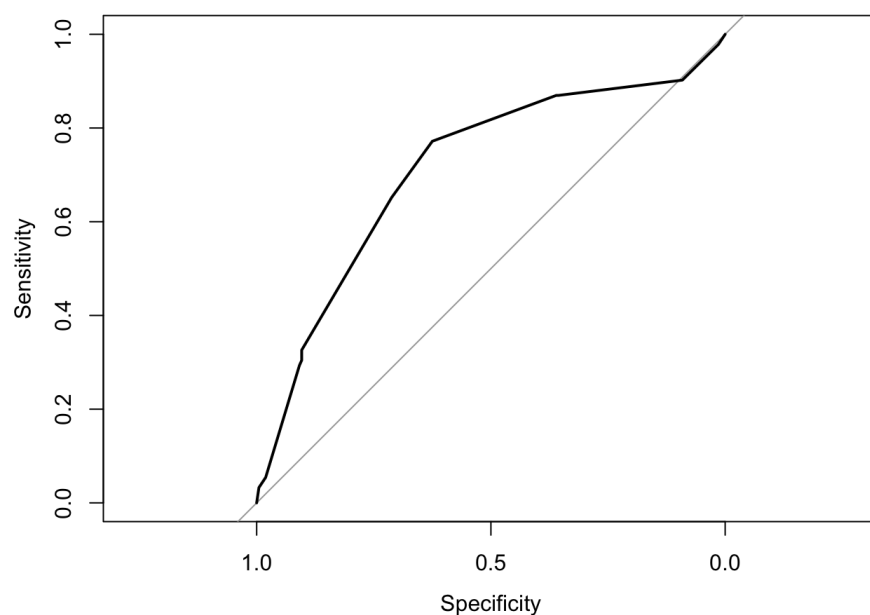
```
## [1] "Deviance:"          "356.791311123743"
```

```
roc(credit.test,predict(tree.credit,newdata=credit[-train,],type="vector")[,2],plot=TRUE)
```

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

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

```
##
## Call:
## roc.default(response = credit.test, predictor = predict(tree.credit,     newdata = credit[-train, ], type = "v
ector")[, 2], plot = TRUE)
##
## Data: predict(tree.credit, newdata = credit[-train, ], type = "vector")[, 2] in 208 controls (credit.test 0) <
92 cases (credit.test 1).
## Area under the curve: 0.7124
```
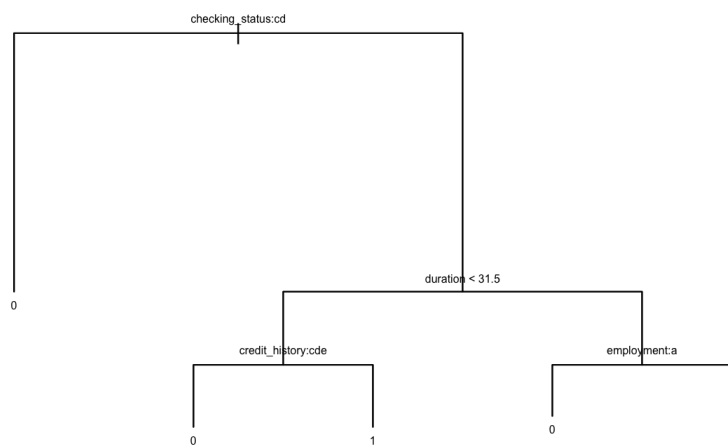
The ROC curve has AUC of 0.7124.

# Question 5

Now try to use pruning to see if you can improve your results.

```
cv.credit<-cv.tree(tree.credit,FUN=prune.misclass)
print(cv.credit)
```

```
## $size
## [1] 11  9  7  5  4  1
##
## $dev
## [1] 193 191 186 184 190 208
##
## $k
## [1]     -Inf  0.00000  1.50000  4.00000  6.00000 13.66667
##
## $method
## [1] "misclass"
##
## attr(,"class")
## [1] "prune"          "tree.sequence"
```

```
prune.credit<-prune.misclass(tree.credit,best=5)
plot(prune.credit)
text(prune.credit,cex=0.5)
```



```
summary(prune.credit)
```

```
##
## Classification tree:
## snip.tree(tree = tree.credit, nodes = c(12L, 2L, 13L))
## Variables actually used in tree construction:
## [1] "checking_status" "duration"        "credit_history"  "employment"
## Number of terminal nodes:  5
## Residual mean deviance:  1.036 = 719.9 / 695
## Misclassification error rate: 0.23 = 161 / 700
```

```
pred.prune.cred<-predict(prune.credit,credit[-train,],type="class")
prune.tab<-table(pred.prune.cred,credit.test)

print(c("Overall error:", (prune.tab[1,2]+prune.tab[2,1])/sum(prune.tab[,])))
```
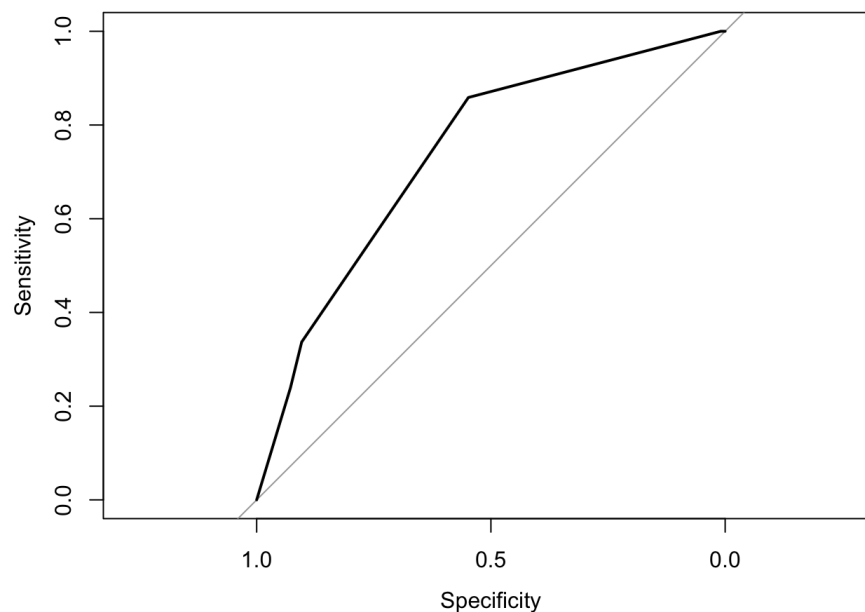
```
## [1] "Overall error:" "0.27"
```

```
prtree.dev.pred<- predict(prune.credit,newdata=credit[-train,],type="tree")
print(c("Deviance:", deviance(prtree.dev.pred)))
```

```
## [1] "Deviance:"       "323.24622569088"
```

```
pred.prune.cred.prob<-predict(prune.credit,credit[-train,])
roc.prune.cred<-roc(credit.test,pred.prune.cred.prob[,2],plot=TRUE)
```

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

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



```
print(c("AUC", roc.prune.cred$auc))
```

```
## [1] "AUC"                 "0.738268185618729"
```

The overall error has slightly improved. The deviance has improved to 323.3462, and the ROC curve now has AUC of 0.7383.

# Question 6

Now try bagging and random forests to see if you can improve your results.

```
set.seed(5)
fullbag.credit<-randomForest(class~.,data=credit,subset=train,mtry=20,importance=TRUE)
yhat.bag<-predict(fullbag.credit,newdata=credit[-train,])
bag.tab<-table(yhat.bag,credit.test)

bag.tab
```

```
##          credit.test
## yhat.bag   0   1
##        0 181  49
##        1  27  43
```

```
rf.credit<-randomForest(class~.,data=credit,subset=train,mtry=20,importance=TRUE)

yhat.rf<-predict(rf.credit,newdata=credit[-train,])
rf.tab<-table(yhat.rf,credit.test)

rf.tab
```
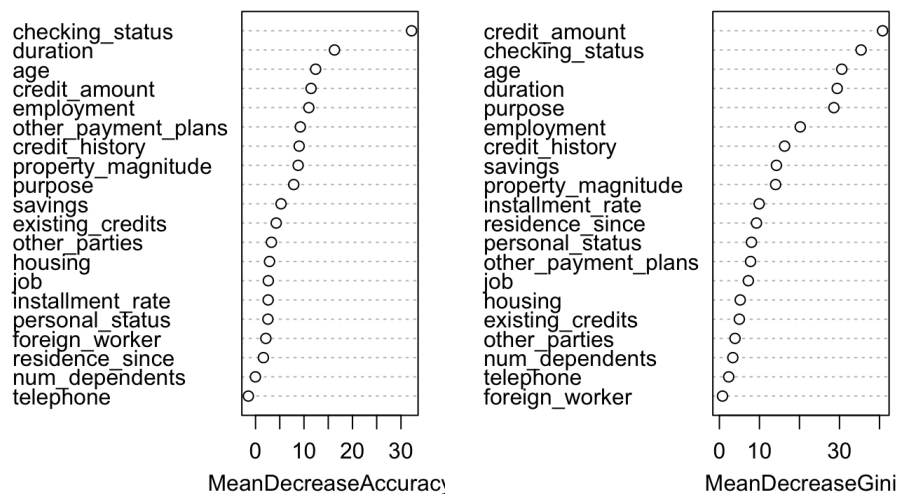
```
##         credit.test
## yhat.rf   0   1
##       0 185  47
##       1  23  45
```

```
varImpPlot(rf.credit)
```

### rf.credit



We see that whether we use the accuracy (prediction error) or Gini index, the predictor checking status seems very important, as well as duration. Other important predictors are credit history and savings, purpose, age and credit amount. These are pretty consistent with the (limited set of ) predictors that were used to construct the pruned tree.

# Question 7

Conclude by making a suggestion as to which of the above methods is the best for this problem.

```
res.table <- data.frame(matrix(ncol = 4, nrow = 4))

colnames(res.table) <- c("Method", "Overall error", "Type I error", "Type II error")

res.table[1,1] <- "Single Tree"
res.table[1,2] <- (tree.tab[1,2]+tree.tab[2,1])/sum(tree.tab[,])
res.table[1,3] <- tree.tab[2,1]/sum(tree.tab[,1])
res.table[1,4] <- tree.tab[1,2]/sum(tree.tab[,2])

res.table[2,1] <- "Pruned Tree"
res.table[2,2] <- (prune.tab[1,2]+prune.tab[2,1])/sum(prune.tab[,])
res.table[2,3] <- prune.tab[2,1]/sum(prune.tab[,1])
res.table[2,4] <- prune.tab[1,2]/sum(prune.tab[,2])

res.table[3,1] <- "Bagging"
res.table[3,2] <- (bag.tab[1,2]+bag.tab[2,1])/sum(bag.tab[,])
res.table[3,3] <- bag.tab[2,1]/sum(bag.tab[,1])
res.table[3,4] <- bag.tab[1,2]/sum(bag.tab[,2])

res.table[4,1] <- "Random Forest"
res.table[4,2] <- (rf.tab[1,2]+rf.tab[2,1])/sum(rf.tab[,])
res.table[4,3] <- rf.tab[2,1]/sum(rf.tab[,1])
res.table[4,4] <- rf.tab[1,2]/sum(rf.tab[,2])

res.table
```

```
##           Method Overall error Type I error Type II error
## 1    Single Tree     0.2800000   0.09615385     0.6956522
## 2    Pruned Tree     0.2700000   0.09615385     0.6630435
## 3        Bagging     0.2533333   0.12980769     0.5326087
## 4  Random Forest     0.2333333   0.11057692     0.5108696
```

Based on our results, the single tree is definitely not a good choice. Among the three other methods, the random forest has the best overall prediction error and type-1 error, but bagging has the lowest type-2 error. Given that bagging and random forest have much less variance, our assessment of their performance is much more reliable than for the single tree and pruned tree and as such, we would recommend using one of those two methods.