

1. Importing the data and defining helper functions

```
library(readxl)
data <- read_excel("~/RStudio/CS_859_Team_Project/NOV23/Drug_Consumption_data_month.xlsx")

f_cols = c("Gender", "Country", "Ethnicity", "Alcohol", "Amphet", "Benzos", "Cannabis", "Ecstasy", "Legalh", "Nicotine")

#drugs = list("Alcohol", "Amphet", "Benzos", "Cannabis", "Ecstasy", "Legalh", "Nicotine")

data[f_cols] <- lapply(data[f_cols], factor)

data
```

```
## # A tibble: 1,885 × 31
##       ID      Age Gender  Educa...1 Country Ethni...2 Nscore   Escore   Oscore Ascore
##   <dbl>   <dbl> <fct>    <dbl> <fct>    <fct>    <dbl>    <dbl>    <dbl>  <dbl>
## 1     1    0.498  0.48246 -0.0592 0.96082 0.126     0.313   -0.575   -0.583  -0.917
## 2     2   -0.0785 -0.48246  1.98    0.96082 -0.316... -0.678    1.94     1.44    0.761
## 3     3    0.498  -0.48246 -0.0592 0.96082 -0.316... -0.467    0.805   -0.847  -1.62
## 4     4   -0.952  0.48246  1.16    0.96082 -0.316... -0.149   -0.806   -0.0193  0.590
## 5     5    0.498  0.48246  1.98    0.96082 -0.316...  0.735   -1.63   -0.452  -0.302
## 6     6    2.59   0.48246 -1.23    0.24923 -0.316... -0.678   -0.300   -1.56    2.04
## 7     7    1.09   -0.48246  1.16   -0.570... -0.316... -0.467   -1.09   -0.452  -0.302
## 8     8    0.498  -0.48246 -1.74    0.96082 -0.316... -1.33     1.94   -0.847  -0.302
## 9     9    0.498  0.48246 -0.0592 0.24923 -0.316...  0.630    2.57   -0.976    0.761
## 10    10    1.82   -0.48246  1.16    0.96082 -0.316... -0.246    0.00332 -1.42    0.590
## # ... with 1,875 more rows, 21 more variables: Cscore <dbl>, Impulsiveness <dbl>,
## # SS <dbl>, Alcohol <fct>, Amphet <fct>, Amyl <dbl>, Benzos <fct>,
## # Caff <dbl>, Cannabis <fct>, Choc <dbl>, Coke <dbl>, Crack <dbl>,
## # Ecstasy <fct>, Heroin <dbl>, Ketamine <dbl>, Legalh <fct>, LSD <dbl>,
## # Meth <dbl>, Mushrooms <dbl>, Nicotine <fct>, VSA <dbl>, and abbreviated
## # variable names 1Education, 2Ethnicity
```

```

alcohol = data[c("Age", "Gender", "Education", "Country", "Ethnicity", "Nscore", "Escore", "Oscore", "Ascore", "Cscore", "Impulsiveness", "SS", "Alcohol")]

amphet = data[c("Age", "Gender", "Education", "Country", "Ethnicity", "Nscore", "Escore", "Oscore", "Ascore", "Cscore", "Impulsiveness", "SS", "Amphet")]

benzos = data[c("Age", "Gender", "Education", "Country", "Ethnicity", "Nscore", "Escore", "Oscore", "Ascore", "Cscore", "Impulsiveness", "SS", "Benzos")]

cannabis = data[c("Age", "Gender", "Education", "Country", "Ethnicity", "Nscore", "Escore", "Oscore", "Ascore", "Cscore", "Impulsiveness", "SS", "Cannabis")]

ecstasy = data[c("Age", "Gender", "Education", "Country", "Ethnicity", "Nscore", "Escore", "Oscore", "Ascore", "Cscore", "Impulsiveness", "SS", "Ecstasy")]

legalh = data[c("Age", "Gender", "Education", "Country", "Ethnicity", "Nscore", "Escore", "Oscore", "Ascore", "Cscore", "Impulsiveness", "SS", "Legalh")]

nicotine = data[c("Age", "Gender", "Education", "Country", "Ethnicity", "Nscore", "Escore", "Oscore", "Ascore", "Cscore", "Impulsiveness", "SS", "Nicotine")]

```

```

f1 <- function(rf_model)

{
p = rf_model$confusion[4]/(rf_model$confusion[3]+rf_model$confusion[4])
r = rf_model$confusion[4]/(rf_model$confusion[2]+rf_model$confusion[4])
#sp = rf_model$confusion[4]/(rf_model$confusion[2]+rf_model$confusion[4])
#se = rf_model$confusion[1]/(rf_model$confusion[2]+rf_model$confusion[4])

score = (2*p*r)/(r+p)

#sprintf('Precision: %.4f', p)
#sprintf('Recall: %.4f', r)
#sprintf('Sensitivity: %.4f', se)
#sprintf('Specificity: %.4f', sp)
sprintf('F1 Score: %.4f', score)
}

```

2. Running a Random Forest

```
library(randomForest)
```

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

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

```
set.seed(365)

rf_oh <- randomForest(Alcohol~.,
                      data = alcohol,
                      importance = TRUE,
                      mtry = 4,
                      ntree = 1000,
                      CUTOFF = .6,
                      verbose = TRUE)

print(rf_oh)
```

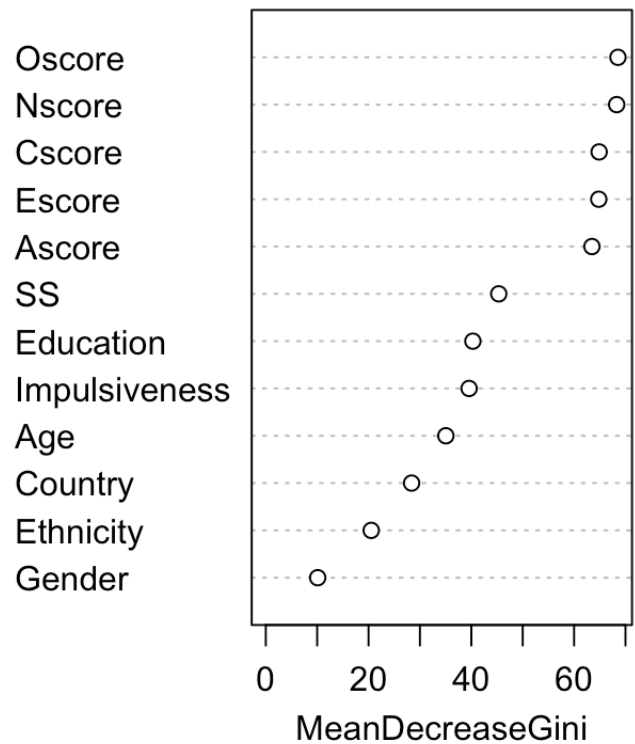
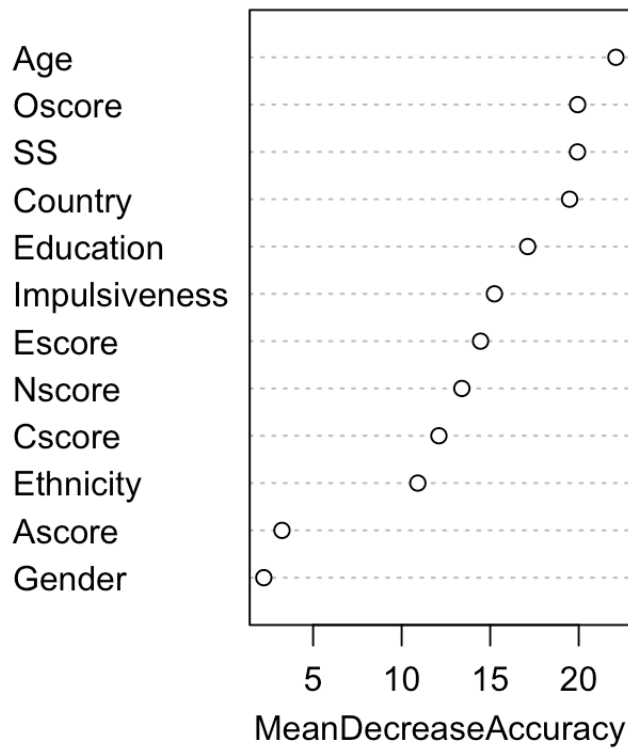
```
##
## Call:
## randomForest(formula = Alcohol ~ ., data = alcohol, importance = TRUE,      mtry
= 4, ntree = 1000, CUTOFF = 0.6, verbose = TRUE)
##              Type of random forest: classification
##              Number of trees: 1000
## No. of variables tried at each split: 4
##
##              OOB estimate of  error rate: 18.14%
## Confusion matrix:
##      0      1 class.error
## 0 11   323   0.96706587
## 1 19  1532   0.01225016
```

```
f1(rf_oh) #Positive Class F1 Score function
```

```
## [1] "F1 Score: 0.8996"
```

```
varImpPlot(rf_oh)
```

rf_oh



```
set.seed(365)

rf_amp <- randomForest(Amphet~.,
                        data = amphet,
                        importance = TRUE,
                        mtry = 4,
                        ntree = 1000,
                        CUTOFF = .6,
                        verbose = TRUE)

print(rf_amp)
```

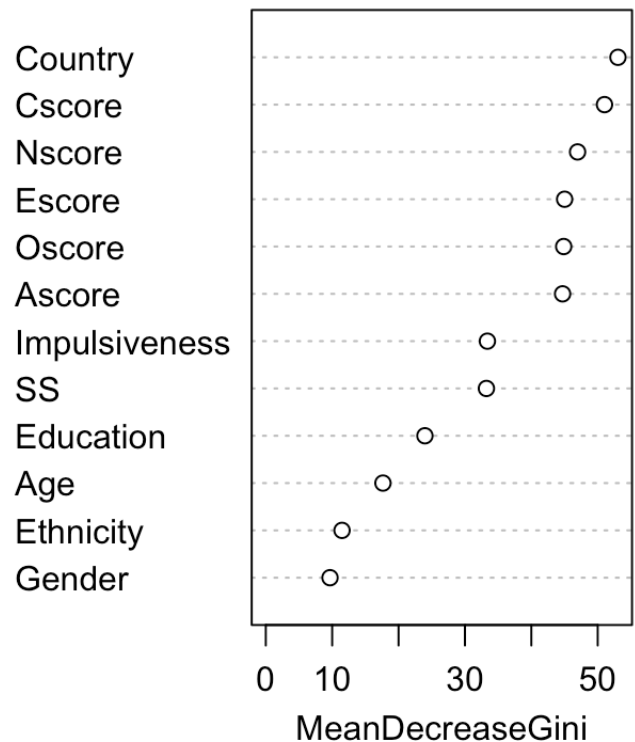
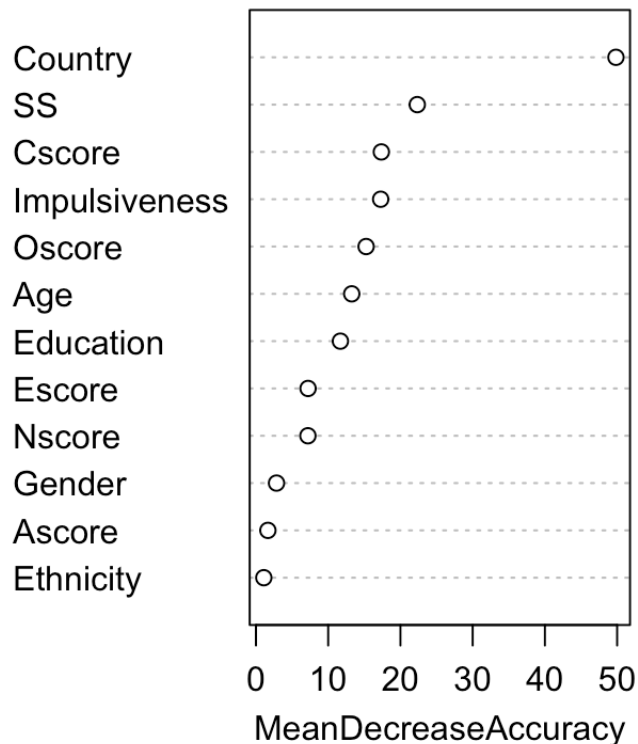
```
##
## Call:
## randomForest(formula = Amphet ~ ., data = amphet, importance = TRUE,          mtry =
4, ntree = 1000, CUTOFF = 0.6, verbose = TRUE)
##
##           Type of random forest: classification
##           Number of trees: 1000
## No. of variables tried at each split: 4
##
##           OOB estimate of  error rate: 13.58%
## Confusion matrix:
##      0  1 class.error
## 0 1616 31    0.0188221
## 1  225 13    0.9453782
```

```
f1(rf_amp) #Positive Class F1 Score function
```

```
## [1] "F1 Score: 0.0922"
```

```
varImpPlot(rf_amp)
```

rf_amp



```
set.seed(365)

rf_ben <- randomForest(Benzos~.,
                        data = benzos,
                        importance = TRUE,
                        mtry = 4,
                        ntree = 1000,
                        CUTOFF = .6,
                        verbose = TRUE)

print(rf_ben)
```

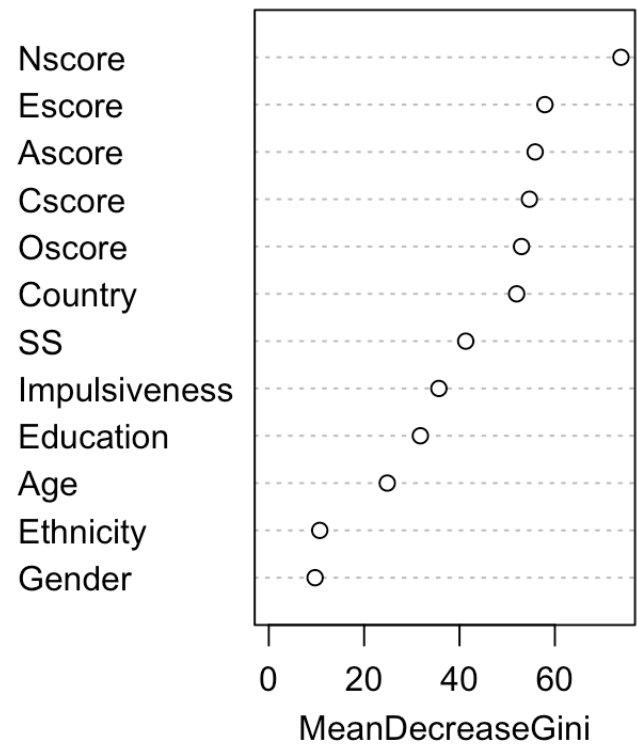
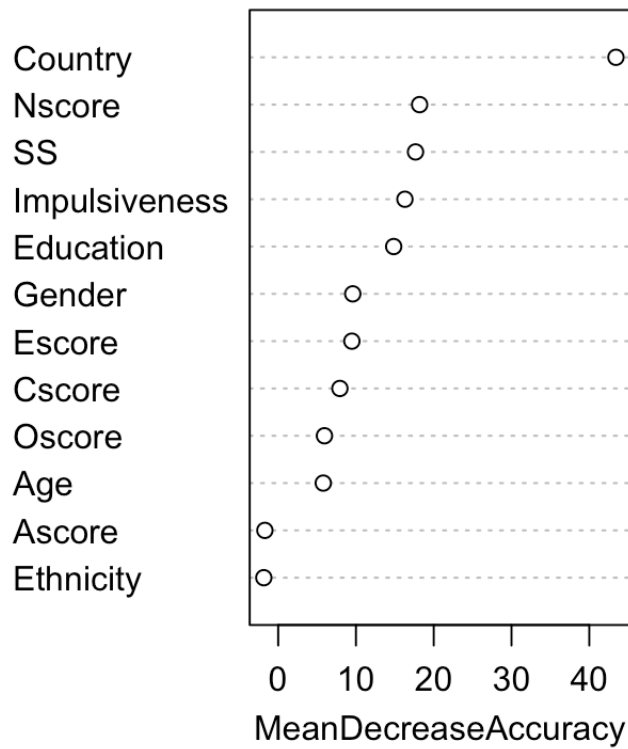
```
##
## Call:
## randomForest(formula = Benzos ~ ., data = benzos, importance = TRUE,      mtry =
4, ntree = 1000, CUTOFF = 0.6, verbose = TRUE)
##              Type of random forest: classification
##              Number of trees: 1000
## No. of variables tried at each split: 4
##
##              OOB estimate of  error rate: 15.97%
## Confusion matrix:
##           0   1 class.error
## 0 1551 35    0.0220681
## 1  266 33    0.8896321
```

```
f1(rf_ben) #Positive Class F1 Score function
```

```
## [1] "F1 Score: 0.1798"
```

```
varImpPlot(rf_ben)
```

rf_ben



```
set.seed(365)

rf_can <- randomForest(Cannabis~.,
                        data = cannabis,
                        importance = TRUE,
                        mtry = 4,
                        ntree = 1000,
                        CUTOFF = .6,
                        verbose = TRUE)

print(rf_can)
```

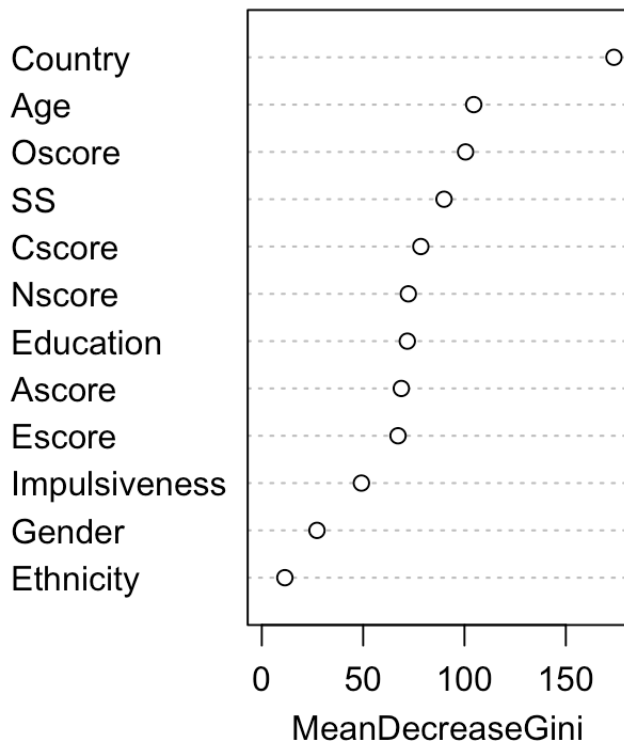
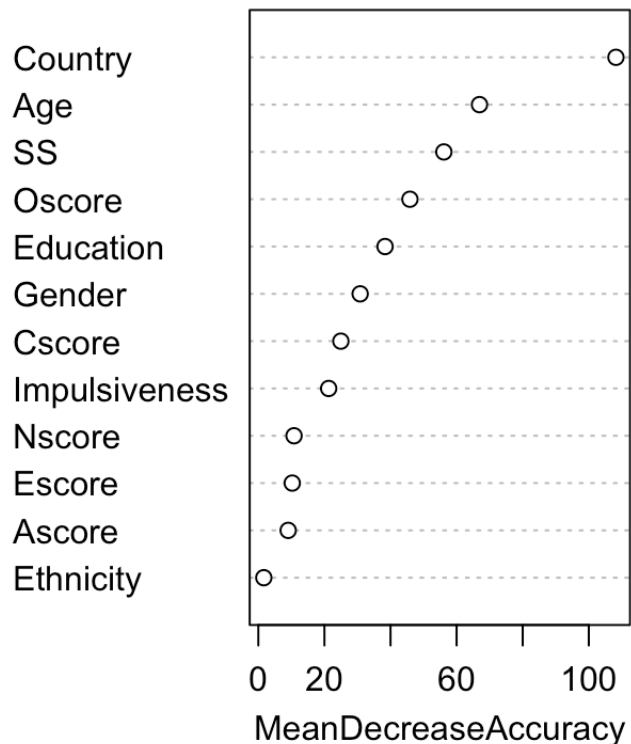
```
##
## Call:
## randomForest(formula = Cannabis ~ ., data = cannabis, importance = TRUE,      mtr
y = 4, ntree = 1000, CUTOFF = 0.6, verbose = TRUE)
##
##           Type of random forest: classification
##
##           Number of trees: 1000
## No. of variables tried at each split: 4
##
##           OOB estimate of  error rate: 20.48%
## Confusion matrix:
##      0    1 class.error
## 0 900 197    0.1795807
## 1 189 599    0.2398477
```

```
f1(rf_can) #Positive Class F1 Score function
```

```
## [1] "F1 Score: 0.7563"
```

```
varImpPlot(rf_can)
```

rf_can




```
set.seed(365)

rf_xt <- randomForest(Ecstasy~.,
                      data = ecstasy,
                      importance = TRUE,
                      mtry = 4,
                      ntree = 1000,
                      CUTOFF = .6,
                      verbose = TRUE)

print(rf_xt)
```

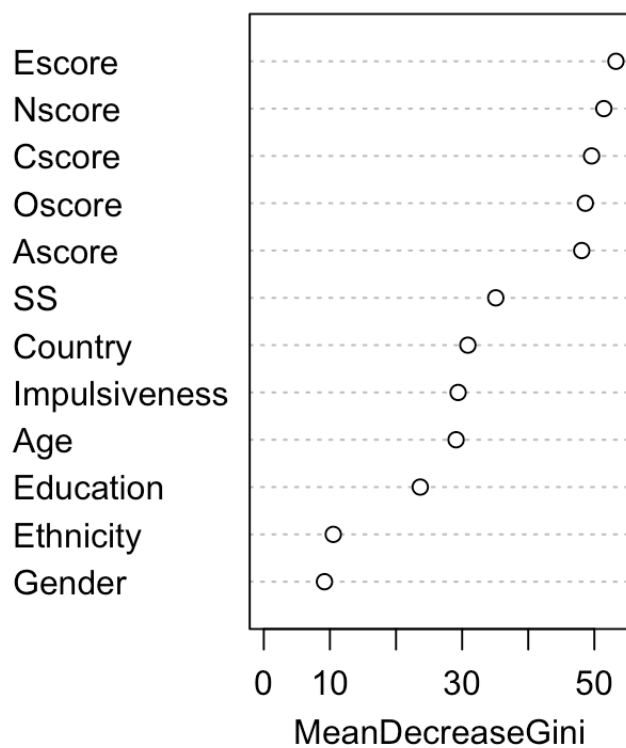
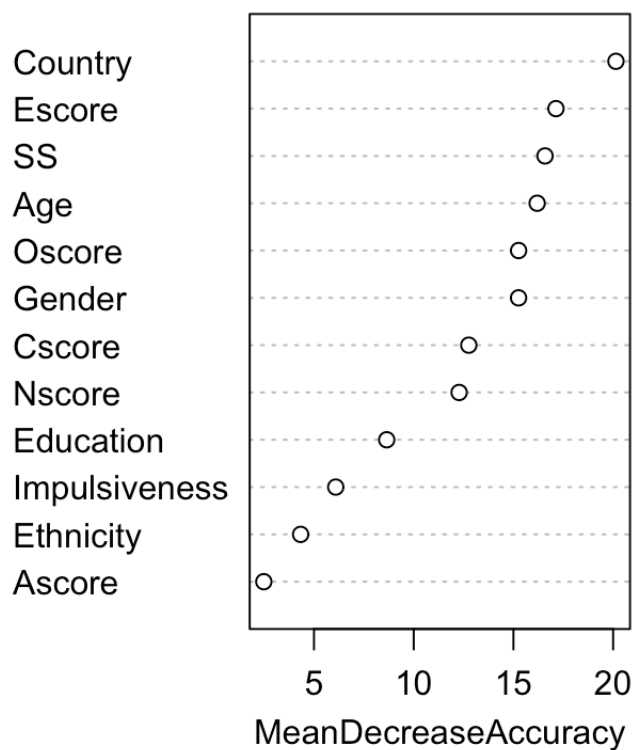
```
##
## Call:
## randomForest(formula = Ecstasy ~ ., data = ecstasy, importance = TRUE,      mtry
= 4, ntree = 1000, CUTOFF = 0.6, verbose = TRUE)
##              Type of random forest: classification
##              Number of trees: 1000
## No. of variables tried at each split: 4
##
##              OOB estimate of  error rate: 12.73%
## Confusion matrix:
##           0   1 class.error
## 0 1635 10 0.006079027
## 1   230 10 0.958333333
```

```
f1(rf_xt) #Positive Class F1 Score function
```

```
## [1] "F1 Score: 0.0769"
```

```
varImpPlot(rf_xt)
```

rf_xt



```
set.seed(365)

rf_lh <- randomForest(Legalh~.,
                      data = legalh,
                      importance = TRUE,
                      mtry = 4,
                      ntree = 1000,
                      CUTOFF = .6,
                      verbose = TRUE)

print(rf_lh)
```

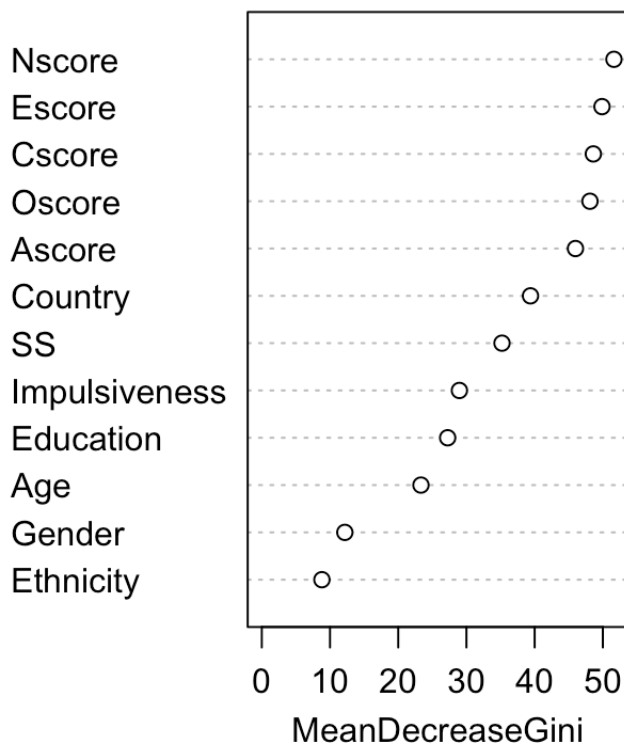
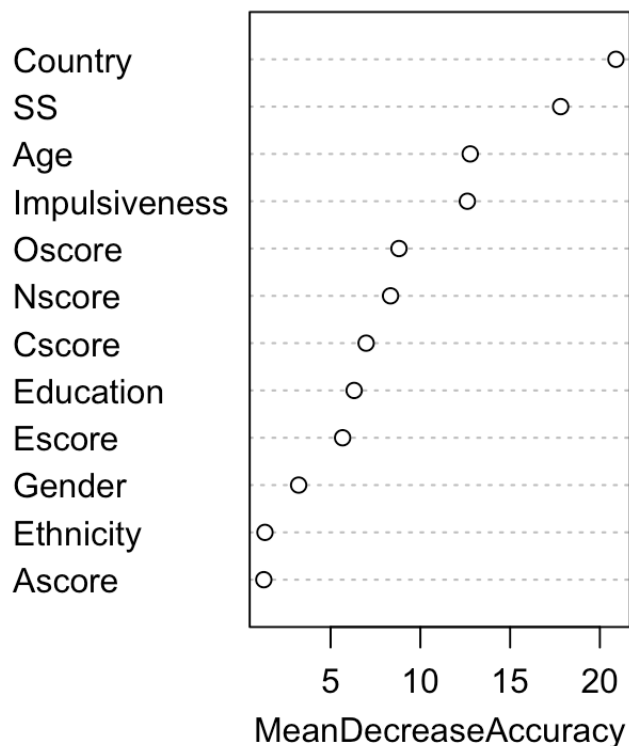
```
##
## Call:
## randomForest(formula = Legalh ~ ., data = legalh, importance = TRUE,      mtry =
4, ntree = 1000, CUTOFF = 0.6, verbose = TRUE)
##
##           Type of random forest: classification
##           Number of trees: 1000
## No. of variables tried at each split: 4
##
##           OOB estimate of  error rate: 13.47%
## Confusion matrix:
##      0  1 class.error
## 0 1623 21  0.01277372
## 1  233  8  0.96680498
```

```
f1(rf_lh) #Positive Class F1 Score function
```

```
## [1] "F1 Score: 0.0593"
```

```
varImpPlot(rf_lh)
```

rf_lh



```
set.seed(365)

rf_nic <- randomForest(Nicotine~.,
                        data = nicotine,
                        importance = TRUE,
                        mtry = 4,
                        ntree = 1000,
                        CUTOFF = .6,
                        verbose = TRUE)

print(rf_nic)
```

```
##
## Call:
## randomForest(formula = Nicotine ~ ., data = nicotine, importance = TRUE,      mtr
y = 4, ntree = 1000, CUTOFF = 0.6, verbose = TRUE)
##              Type of random forest: classification
##              Number of trees: 1000
## No. of variables tried at each split: 4
##
##              OOB estimate of  error rate: 33.85%
## Confusion matrix:
##      0    1 class.error
## 0 689 321    0.3178218
## 1 317 558    0.3622857
```

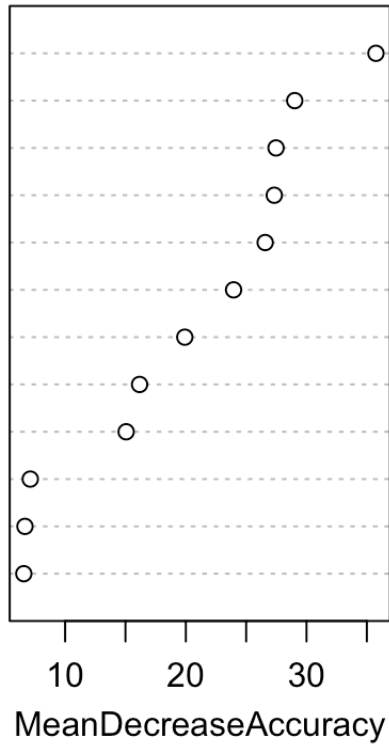
```
f1(rf_nic) #Positive Class F1 Score function
```

```
## [1] "F1 Score: 0.6363"
```

```
varImpPlot(rf_nic)
```

rf_nic

SS
Age
Country
Education
Impulsiveness
Cscore
Oscore
Gender
Nscore
Ascore
Ethnicity
Escore



Cscore
Nscore
Oscore
Ascore
Escore
SS
Education
Impulsiveness
Country
Age
Gender
Ethnicity

