

2. Gene expression analysis: Creating random forest models on lab data

Fay Webster

2022-05-27

Aim:

- Predicting health impact of infections utilizing immune parameters as predictors
- Predicted variable: WL as a proxy of health
- To do that we are using immune data from experimental lab infections.
- We are training random forest models on the immune data from experimental lab infections
- And we test them on the field.
- We then compare the differences in the predicted health impact among non-hybrid and hybrid mice.

In this document I am preparing the models using the lab data only.

Load necessary libraries:

```
#install.packages("optimx", version = "2021-10.12") # this package is required for  
#the parasite load package to work  
library(tidyverse)  
library(tidyr)  
library(dplyr)  
library(cowplot)  
library(randomForest)  
library(ggplot2)  
library(caret)  
library(ggpubr)  
library(rfUtilities) # Implements a permutation test cross-validation for  
# Random Forests models
```

Laboratory data

Importing the data

We start with the data from experimental lab infections.

```
# Here we import the cleaned data set from the previous script derived from the  
# data set challenge infections  
g <-  
  read.csv("https://raw.githubusercontent.com/fayweb/Eimeria_mouse_immunity/main/output_data/gene_expression_data.csv")  
  
# vectors for selecting gene columns  
Genes <- c("IFNy", "CXCR3_bio", "IL.6", "IL.10", "IL.13", "IL.10", "IL.13",  
           "IL1RN", "CASP1", "CXCL9", "IDO1", "IRGM1", "MPO", "MUC2", "MUC5AC",  
           "MYD88", "NCR1", "PRF1", "RETNLB", "SOCS1", "TICAM1", "TNF")
```

Data cleaning / preparation

```
# we need to change the in challenge infections to a factor
g$Parasite_challenge <- as.factor(g$Parasite_challenge)
g$Eim_MC <- as.factor(g$Eim_MC)

# Here I create a new column, where we get the actual infection status
# According to the melting curve for eimeria
g <- g %>%
  dplyr::mutate(current_infection = case_when(
    Parasite_challenge == "E_ferrisi" & Eim_MC == "TRUE" ~ "E_ferrisi",
    Parasite_challenge == "E_ferrisi" & Eim_MC == "FALSE" ~ "uninfected",
    Parasite_challenge == "E_falciformis" & Eim_MC == "TRUE" ~ "E_falciformis",
    Parasite_challenge == "E_falciformis" & Eim_MC == "FALSE" ~ "uninfected",
    Parasite_challenge == "uninfected" & Eim_MC == "TRUE" ~ "infected_eimeria",
    Parasite_challenge == "uninfected" & Eim_MC == "FALSE" ~ "uninfected",
    TRUE ~ ""
  ))

# create variable maximum weight loss instead of maximum relative weight loss
g <- g %>% dplyr::mutate(max_WL = 100 - max_WL)
```

Imputation of missing data

Imputing missing data + cleaning

Here, I am using a function from the random forest package, `rfImpute` which utilizes random forests to impute missing data in the other variables.

The variables used for imputing mainly the immune gene expression are the current infection, the state of *Eimeria* infection, oocysts and the non-missing genes.

```
#Start by selecting only the genes and the maximum weight loss for each mouse
# Apparently the relative end weight doesn't work so well for predictions

g.1 <- g %>% dplyr::select(c(all_of(Genes), max_WL,
                             mouse_strain, CD4, Treg, Div_Treg,
                             Treg17, Th1, Div_Th1, Th17, Div_Th17, CD8, Act_CD8,
                             Div_Act_CD8, IFNy_CD4, IFNy_CD8, Treg_prop,
                             IL17A_CD4))

sapply(g.1, function(x) sum(is.na(x)))
```

##	IFNy	CXCR3_bio	IL.6	IL.10	IL.13	IL1RN
##	27	0	10	10	86	0
##	CASP1	CXCL9	ID01	IRGM1	MPO	MUC2
##	2	0	0	0	15	0
##	MUC5AC	MYD88	NCR1	PRF1	RETNLB	SOCS1
##	0	0	10	23	0	0
##	TICAM1	TNF	max_WL	mouse_strain	CD4	Treg
##	1	2	0	0	36	36
##	Div_Treg	Treg17	Th1	Div_Th1	Th17	Div_Th17
##	36	36	36	36	36	36
##	CD8	Act_CD8	Div_Act_CD8	IFNy_CD4	IFNy_CD8	Treg_prop
##	36	36	36	36	36	68
##	IL17A_CD4					

```
##          68
g.1 <- g.1 %>% mutate_if(is.character, as.factor)
g.1 <- g.1 %>% mutate_if(is.integer, as.numeric)

# to get reproducible results we use a seed
set.seed(42)

# We want the maximum weight loss to be predicted by the data in all of
# the other columns

# iter = how many random forests are needed, in theory 6 are enough
g.imputed <- rfImpute(max_WL ~ ., data = g.1, iter = 6)

##          |      Out-of-bag      |
## Tree |      MSE   %Var(y) |
## 300 |      27.59   64.60 |
##          |      Out-of-bag      |
## Tree |      MSE   %Var(y) |
## 300 |      28.21   66.06 |
##          |      Out-of-bag      |
## Tree |      MSE   %Var(y) |
## 300 |      29.31   68.63 |
##          |      Out-of-bag      |
## Tree |      MSE   %Var(y) |
## 300 |      28.1    65.80 |
##          |      Out-of-bag      |
## Tree |      MSE   %Var(y) |
## 300 |      28.2    66.02 |
##          |      Out-of-bag      |
## Tree |      MSE   %Var(y) |
## 300 |      29.02   67.94 |

g_minus <- g %>%
  dplyr::select(-c(all_of(Genes), max_WL,
                        mouse_strain, CD4, Treg, Div_Treg,
                        Treg17, Th1, Div_Th1, Th17, Div_Th17, CD8, Act_CD8,
                        Div_Act_CD8, IFNy_CD4, IFNy_CD8, Treg_prop,
                        IL17A_CD4))

#full data set containing the imputed gene expression data
g.imputed <- cbind(g_minus, g.imputed)

How many mice are in the infection planning?

g.imputed %>%
  group_by(Parasite_challenge) %>%
  summarize(length(EH_ID))

## # A tibble: 3 x 2
##   Parasite_challenge `length(EH_ID)`
##   <fct>              <int>
## 1 E_falciformis      22
## 2 E_ferrisi          47
## 3 uninfected         47

How many mice are indeed infected?
```

```
g.imputed %>%
  filter(infection == "challenge") %>%
  group_by(current_infection) %>%
  summarize(length(EH_ID))
```

```
## # A tibble: 4 x 2
##   current_infection `length(EH_ID)`
##   <chr>              <int>
## 1 E_falciformis      22
## 2 E_ferrisi          39
## 3 infected_eimeria   9
## 4 uninfected        46
```

I guess mice got mixed up here?

Splitting data into training and testing sets

Splitting between training and testing: - Assess model performance on unseen data - Avoid over-fitting

Random forest for predicting percentage of maximum weight loss

Dividing data into training and testing

```
Genes <- c("IFNy", "IL.6", "IL.10", "IL.13", "IL.10", "IL.13", "IL1RN",
           "CASP1", "CXCL9", "IDO1", "IRGM1", "MPO", "MUC2", "MUC5AC", "MYD88",
           "NCR1", "PRF1", "RETNLB", "SOCS1", "TICAM1", "TNF")

g.imputed_full <- g.imputed

write.csv(g.imputed_full,
          "output_data/gene_expression/data_products/lab_imputed_gene_expression.csv",
          row.names = FALSE)

#select the relevant columns:
g.imputed <- g.imputed %>%
  dplyr::select(c(max_WL, all_of(Genes)))

# split data into training and test

set.seed(123) # this will help us reproduce this random assignment

# in this way we can pick the random numbers

training.samples <- g.imputed$max_WL%>%
  createDataPartition(p = .7,
                      list = FALSE)

# this is the partiicipation! In this case 0.7 = training data and 0.3 = testing
# we don't want to get a list in return

train.data <- g.imputed[training.samples, ]
test.data <- g.imputed[-training.samples, ]
```

Building the model

```
#train the model
weight_loss_predict <- randomForest(max_WL ~., data = train.data,
                                     proximity = TRUE, ntree = 1000)

# ntree = number of trees

# save the model
save(weight_loss_predict, file = "r_scripts/models/predict_weight_loss.RData")

print(weight_loss_predict)

##
## Call:
## randomForest(formula = max_WL ~ ., data = train.data, proximity = TRUE,      ntree = 1000)
##              Type of random forest: regression
##              Number of trees: 1000
## No. of variables tried at each split: 6
##
##              Mean of squared residuals: 30.55973
##              % Var explained: 26.91
```

Plotting the `weight_loss_predict` will illustrate the error rate as we average across more trees and shows that our error rate stabilizes with around 200 trees.

Model - quality testing

Cross-validation

MSE: As a brief explanation, mean squared error (MSE) is the average of the summation of the squared difference between the actual output value and the predicted output value. Our goal is to reduce the MSE as much as possible.

Variance explained: %explained variance is a measure of how well out-of-bag predictions explain the target variance of the training set.

```
predict_WL_cv <- rf.crossValidation(x = weight_loss_predict, xdata = train.data,
                                   p = 0.10, n = 99, ntree = 501)
```

```
## running: regression cross-validation with 99 iterations
```

```
predict_WL_cv$fit.var.exp
```

```
## [1] 26.91
```

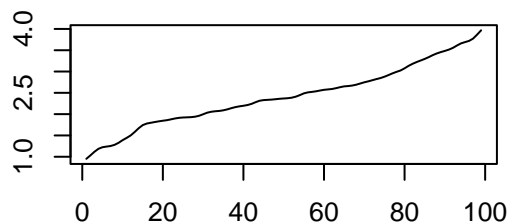
```
par(mfrow=c(2,2))
```

```
plot(predict_WL_cv)
```

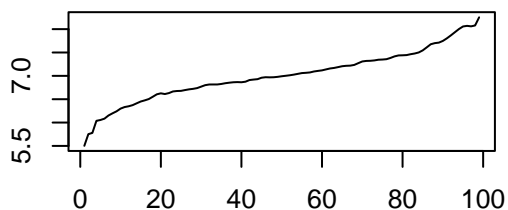
```
# Root Mean Squared Error (observed vs. predicted) from each Bootstrap
# iteration (cross-validation)
plot(predict_WL_cv, stat = "mse")
```

```
#Percent variance explained from specified fit model
plot(predict_WL_cv, stat = "var.exp")
```

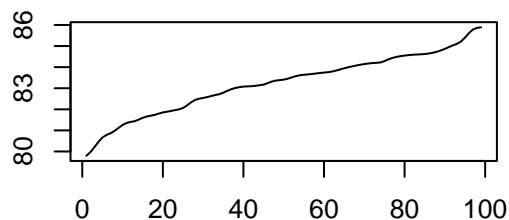
```
#Mean Absolute Error from each Bootstrapped model
plot(predict_WL_cv, stat = "mae")
```

Cross-validated Root Mean Squared Error

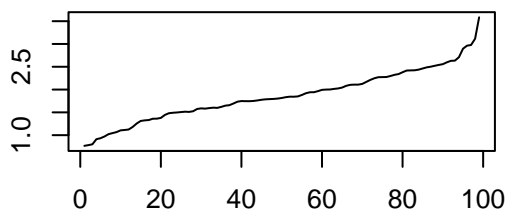
rmse

Model Mean Square Error

mse

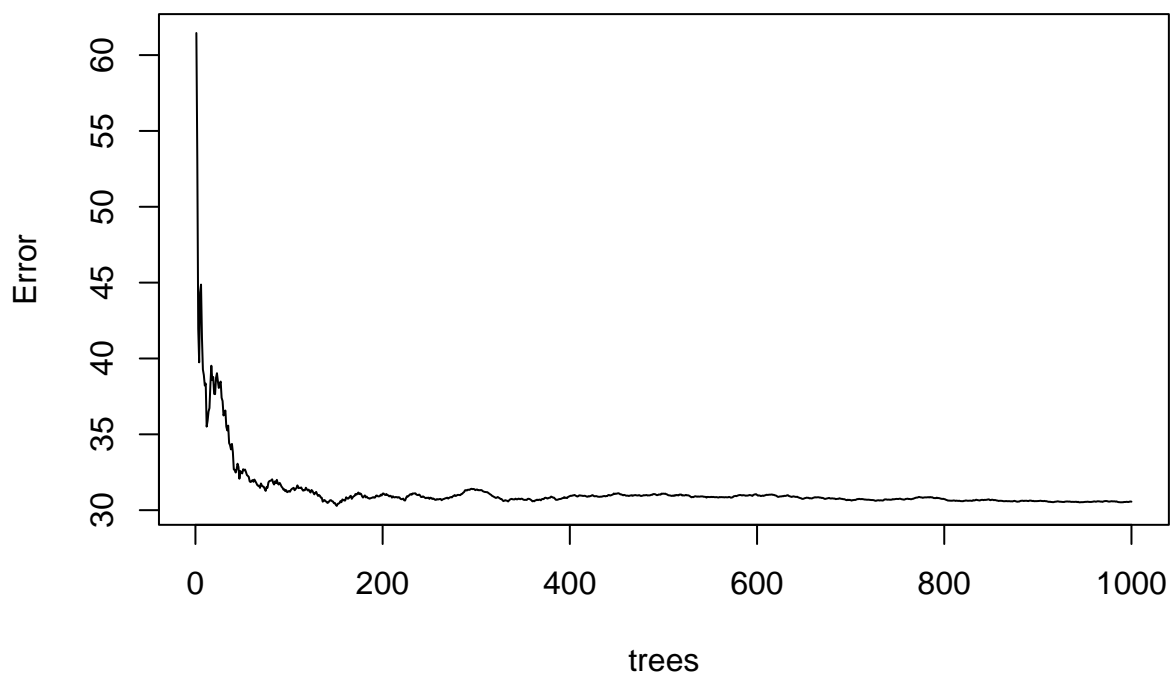
Model percent variance explained

var.exp

Cross-validated Mean Absolute Error

mae

```
plot(weight_loss_predict)
```

weight_loss_predict

The plotted error rate above is based on the OOB sample error and can be accessed directly at `m1$mse`. Thus, we can find which number of trees providing the lowest error rate, which is 257 trees providing an weight error of 5.024738.

```
# number of trees with lowest MSE
which.min(weight_loss_predict$mse)
```

```
## [1] 151
## [1] 257

# RMSE of this optimal random forest
sqrt(weight_loss_predict$mse[which.min(weight_loss_predict$mse)])

## [1] 5.502798
## [1] 5.024738
```

<https://uc-r.github.io/s>

RandomForest also allows us to use a validation set to measure predictive accuracy if we did not want to use the OOB samples.

Tutorial: <https://hackernoon.com/random-forest-regression-in-r-code-and-interpretation>

Random forest regression in R provides two outputs: decrease in mean square error (MSE) and node purity. Prediction error described as MSE is based on permuting out-of-bag sections of the data per individual tree and predictor, and the errors are then averaged. In the regression context, Node purity is the total decrease in residual sum of squares when splitting on a variable averaged over all trees (i.e. how well a predictor decreases variance). MSE is a more reliable measure of variable importance. If the two importance metrics show different results, listen to MSE. If all of your predictors are numerical, then it shouldn't be too much of an issue

Mean Decrease Gini (IncNodePurity) - This is a measure of variable importance based on the Gini impurity index used for the calculating the splits in trees.

Improving Your Model Your model depends on the quality of your dataset and the type of Machine Learning algorithm used. Therefore, to improve the accuracy of your model, you should:

Check what attributes affect our model the most and what variables to leave out in future analysis Find out what other attributes affect a person's wage; we can use as predictors in future analysis Tweak the algorithm (e.g. change the ntree value) Use a different machine learning algorithm If any of these reduces the RMSE significantly, you have succeeded in improving your model!

Application of weight_loss_predict

Using the testing data

Let's now make some predictions using our test data.

```
#The predict() function in R is used to predict the values based on the
# input data.
predictions <- predict(weight_loss_predict, test.data)

# assign test.data to a new object, so that we can make changes
result <- test.data

#add the new variable of predictions to the result object
result <- cbind(result, predictions)

#add the results to a data frame containing test data and the prediction
result <- cbind(g[row.names(result), ], predictions)

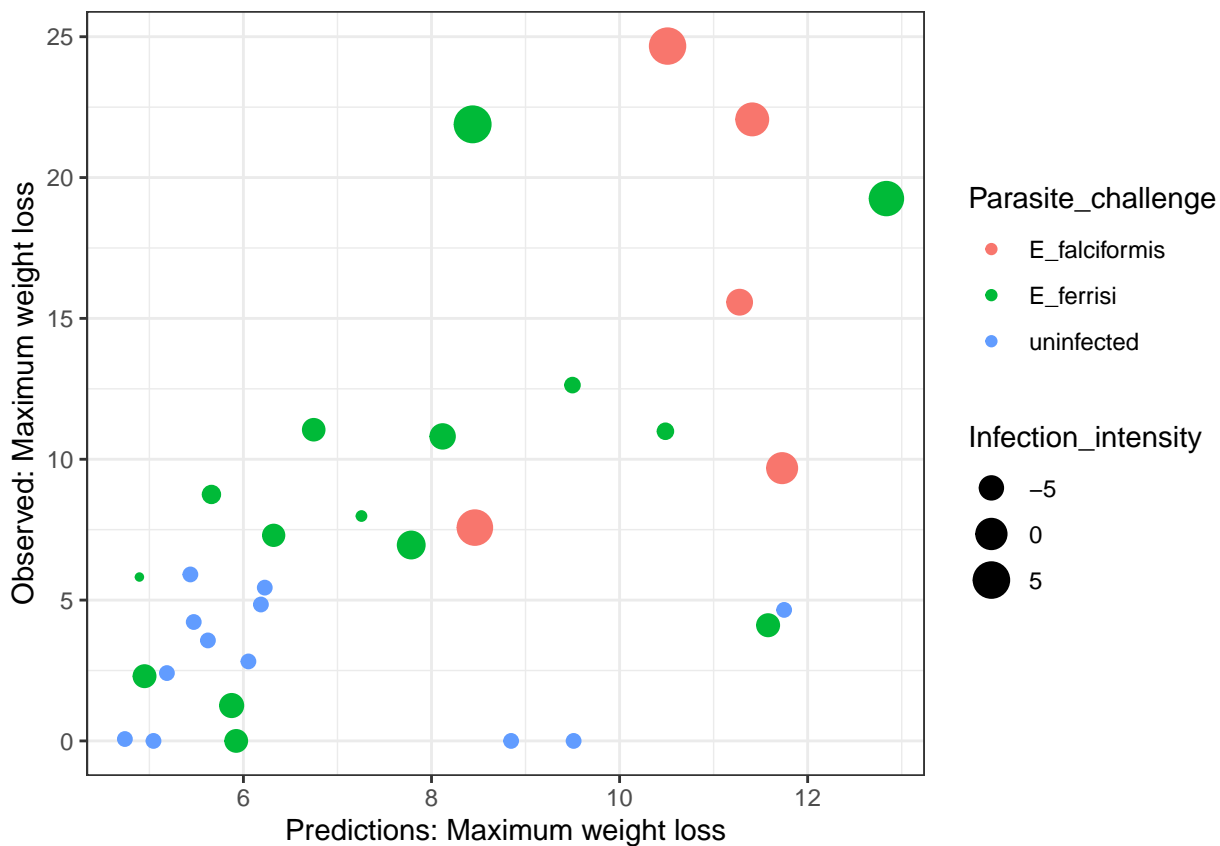
# what is the correlation between predicted and actual data?
cor(result$max_WL, result$predictions,
    method = c("pearson", "kendall", "spearman"))

## [1] 0.5810315
```

Visualizing the predictions

```
# trying to find a way to represent the delta ct for the negative ones
# please find a better way to do this
result <- result %>%
  dplyr::mutate(Infection_intensity = case_when(
    Parasite_challenge == "uninfected" ~ -9,
    TRUE ~ delta
  ))

result %>%
  ggplot() +
  geom_point(aes(x = predictions, y = max_WL,
                 color = Parasite_challenge, size = Infection_intensity)) +
  labs(x = "Predictions: Maximum weight loss",
       y = "Observed: Maximum weight loss") +
  theme_bw()
```



Using the same method to predict either Melting curve or infecting parasite

(2nd validation)

As a second part I am using the same method to predict either infection with Eimeria in general or the species of eimeria.

Predicting eimeria species

Predicting parasite: splitting into training and testing

```
g.imputed_full$Parasite_challenge <-  
  as.factor(g.imputed_full$Parasite_challenge)  
  
#select the relevant columns:  
g.imputed_parasite <- g.imputed_full %>%  
  dplyr::select(c(Parasite_challenge, all_of(Genes)))  
  
# split data into training and test  
set.seed(123) # this will help us reproduce this random assignment  
# in this way we can pick the random numbers  
training.samples_parasite <- g.imputed_parasite$Parasite_challenge%>%  
  createDataPartition(p = .7, list = FALSE)  
train.data_parasite <- g.imputed_parasite[training.samples, ]  
test.data_parasite <- g.imputed_parasite[-training.samples, ]
```

Building the model_Parasite

```
#train the model  
model_Parasite <- randomForest(Parasite_challenge ~.,  
                               data = train.data_parasite, proximity = TRUE,  
                               ntree = 1500) # number of trees  
  
# save the model  
save(model_Parasite, file = "r_scripts/models/predict_infecting_parasite.RData")  
  
print(model_Parasite)
```

```
##  
## Call:  
## randomForest(formula = Parasite_challenge ~ ., data = train.data_parasite, proximity = TRUE, ntree = 1500)  
## Type of random forest: classification  
## Number of trees: 1500  
## No. of variables tried at each split: 4  
##  
## OOB estimate of error rate: 38.1%  
## Confusion matrix:  
## E_falciformis E_ferrisi uninfected class.error  
## E_falciformis 3 8 6 0.8235294  
## E_ferrisi 5 20 7 0.3750000  
## uninfected 2 4 29 0.1714286
```

OOB = 46.43, this means that only 53 % of our predictions are accurate

Quality checks

Cross-validation MSE: As a brief explanation, mean squared error (MSE) is the average of the summation of the squared difference between the actual output value and the predicted output value. Our goal is to reduce the MSE as much as possible.

Variance explained: %explained variance is a measure of how well out-of-bag predictions explain the target variance of the training set.

```
model_Parasite_cv <- rf.crossValidation(x = model_Parasite, xdata =  
                                       train.data_parasite,  
                                       p = 0.10, n = 99, ntree = 501)
```

```
## running: classification cross-validation with 99 iterations
```

```
model_Parasite_cv$fit.var.exp
```

```
## NULL
```

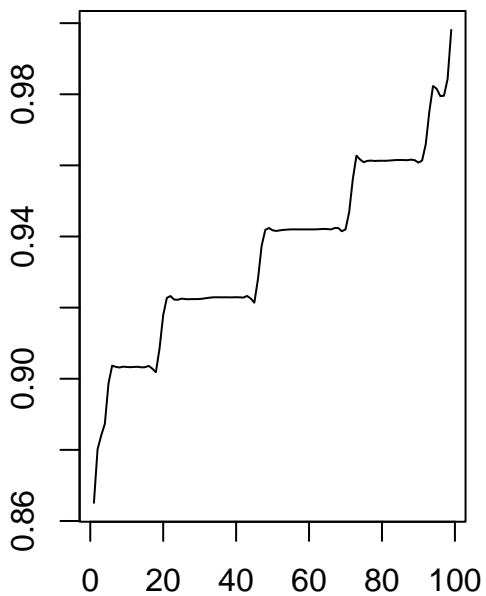
```
# Plot cross validation versus model producers accuracy
```

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

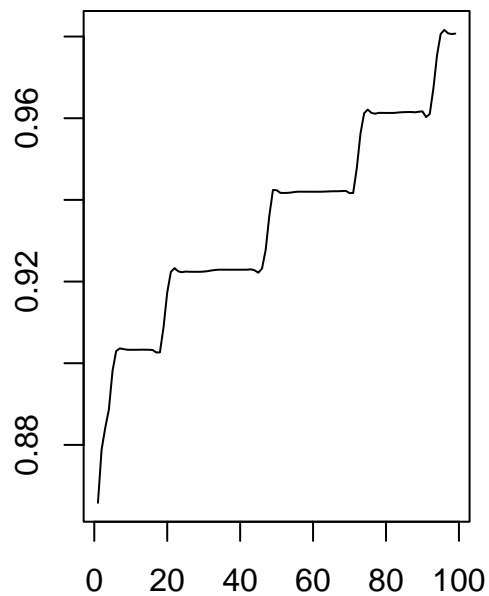
```
plot(model_Parasite_cv, type = "cv", main = "CV producers accuracy")
```

```
plot(model_Parasite_cv, type = "model", main = "Model producers accuracy")
```

CV producers accuracy



Model producers accuracy



```
# Plot cross validation versus model oob
```

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

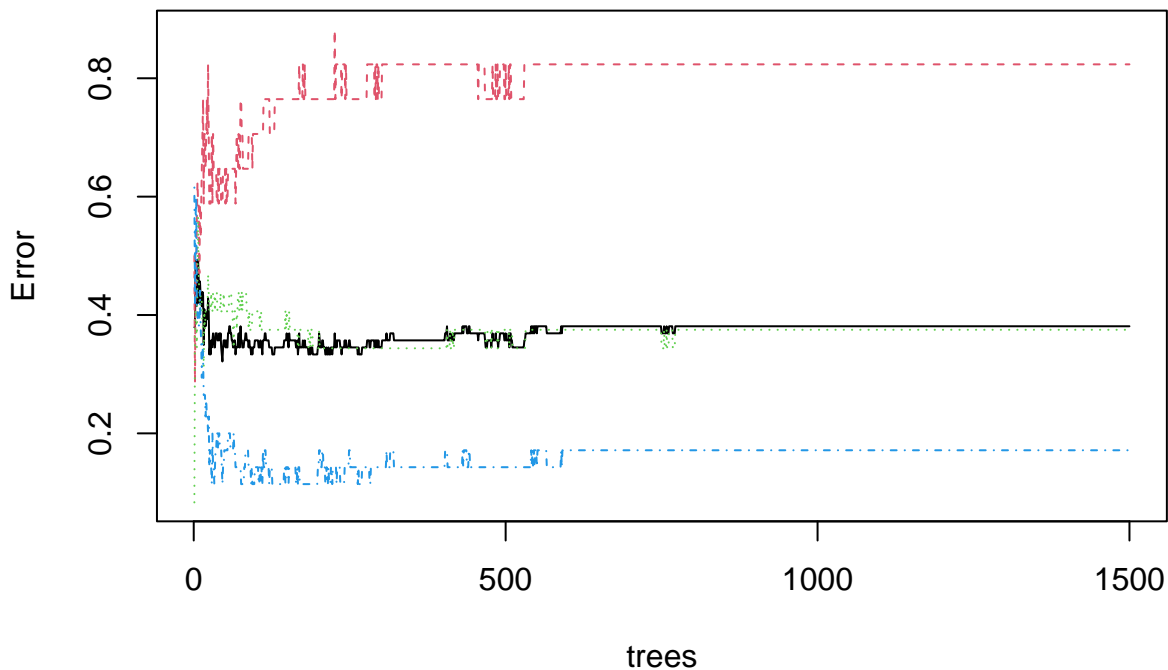
```
# plot(model_Parasite_cv, type = "cv", stat = "oob", main = "CV oob error")
```

```
#plot(model_Parasite_cv, type = "model", stat = "oob",
```

```
#   main = "Model oob error")
```

```
plot(model_Parasite)
```

model_Parasite



Testing the model: Predictions

```
#The predict() function in R is used to predict the values based on the input
# data.
predictions_parasite <- predict(model_Parasite, test.data_parasite)
# assign test.data to a new object, so that we can make changes
result_parasite <- test.data_parasite
#add the new variable of predictions to the result object
result_parasite <- cbind(result_parasite, predictions_parasite)
#add the results to a data frame containing test data and the prediction
result_parasite <- cbind(g[row.names(result_parasite), ], predictions_parasite)
```

Visualizing predictions_parasite

```
conf_matrix_parasite <-
  confusionMatrix(
    result_parasite$predictions_parasite,
    reference = result_parasite$Parasite_challenge)

print(conf_matrix_parasite)

## Confusion Matrix and Statistics
##
##               Reference
## Prediction   E_falciformis E_ferrisi uninfected
## E_falciformis          3          1          0
## E_ferrisi              2         11          1
## uninfected             0          3         11
##
## Overall Statistics
##
```

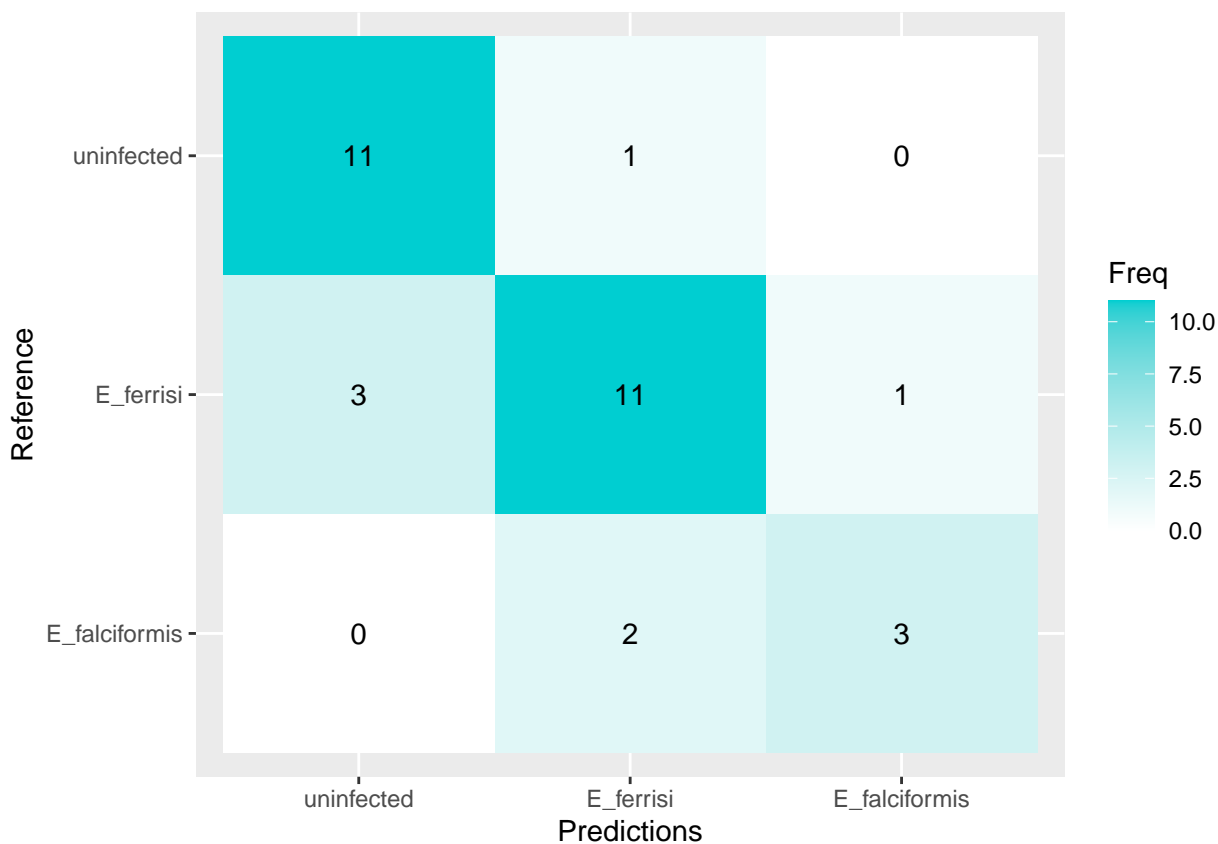
```
## Accuracy : 0.7812
## 95% CI : (0.6003, 0.9072)
## No Information Rate : 0.4688
## P-Value [Acc > NIR] : 0.0003084
##
## Kappa : 0.6422
##
## McNemar's Test P-Value : NA
##
## Statistics by Class:
##
## Class: E_falciformis Class: E_ferrisi Class: uninfected
## Sensitivity          0.60000          0.7333          0.9167
## Specificity          0.96296          0.8235          0.8500
## Pos Pred Value       0.75000          0.7857          0.7857
## Neg Pred Value       0.92857          0.7778          0.9444
## Prevalence           0.15625          0.4688          0.3750
## Detection Rate       0.09375          0.3438          0.3438
## Detection Prevalence 0.12500          0.4375          0.4375
## Balanced Accuracy    0.78148          0.7784          0.8833
```

```
conf_matrix_parasite$table
```

```
## Reference
## Prediction E_falciformis E_ferrisi uninfected
## E_falciformis      3          1          0
## E_ferrisi          2         11          1
## uninfected         0          3         11
```

```
plt <- as.data.frame(conf_matrix_parasite$table)
plt$Prediction <- factor(plt$Prediction, levels=rev(levels(plt$Prediction)))

ggplot(plt, aes(x = Prediction, y = Reference, fill= Freq)) +
  geom_tile() + geom_text(aes(label=Freq)) +
  scale_fill_gradient(low="white", high="darkturquoise") +
  labs(x = "Predictions", y = "Reference")
```



Predicting for Melting curve

Split the data again into training and testing

```
#select the relevant columns:
g.imputed_mc <- g.imputed_full %>%
  dplyr::select(c(Eim_MC, all_of(Genes)))

# split data into training and test
set.seed(123) # this will help us reproduce this random assignment
# in this way we can pick the random numbers
training.samples_mc <- g.imputed_mc$Eim_MC %>%
  createDataPartition(p = .7, list = FALSE)
train.data_mc <- g.imputed_mc[training.samples, ]
test.data_mc <- g.imputed_mc[-training.samples, ]
```

Building the model

```
#train the model
model_mc <- randomForest(Eim_MC ~., data = train.data_mc, proximity = TRUE,
  ntree = 1500) # number of trees

# save the model
save(model_mc, file = "r_scripts/models/predicting_mc.RData")

print(model_mc)
```

```
##
## Call:
## randomForest(formula = Eim_MC ~ ., data = train.data_mc, proximity = TRUE,      ntree = 1500)
##              Type of random forest: classification
##              Number of trees: 1500
## No. of variables tried at each split: 4
##
## OOB estimate of error rate: 32.14%
## Confusion matrix:
##      FALSE TRUE class.error
## FALSE    21   16  0.4324324
## TRUE     11   36  0.2340426
```

Cross-validation MSE: As a brief explanation, mean squared error (MSE) is the average of the summation of the squared difference between the actual output value and the predicted output value. Our goal is to reduce the MSE as much as possible.

Variance explained: %explained variance is a measure of how well out-of-bag predictions explain the target variance of the training set.

```
model_mc_cv <- rf.crossValidation(x = model_mc, xdata = train.data_mc,
                                p = 0.10, n = 99, ntree = 501)
```

```
## running: classification cross-validation with 99 iterations
```

```
model_mc_cv$fit.var.exp
```

```
## NULL
```

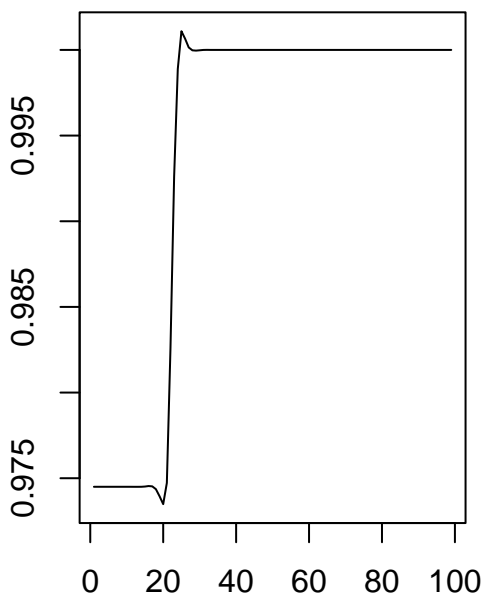
```
# Plot cross validation versus model producers accuracy
```

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

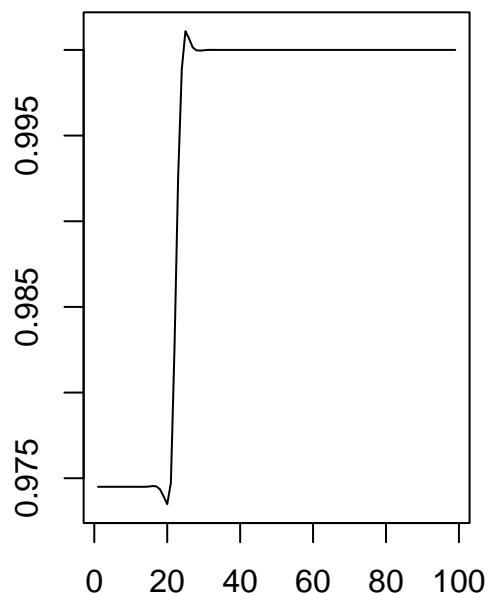
```
plot(model_mc_cv, type = "cv", main = "CV producers accuracy")
```

```
plot(model_mc_cv, type = "model", main = "Model producers accuracy")
```

CV producers accuracy



Model producers accuracy



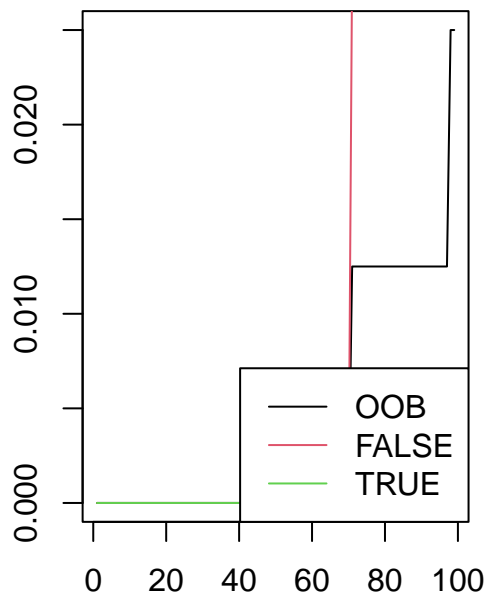
```
# Plot cross validation versus model oob
```

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

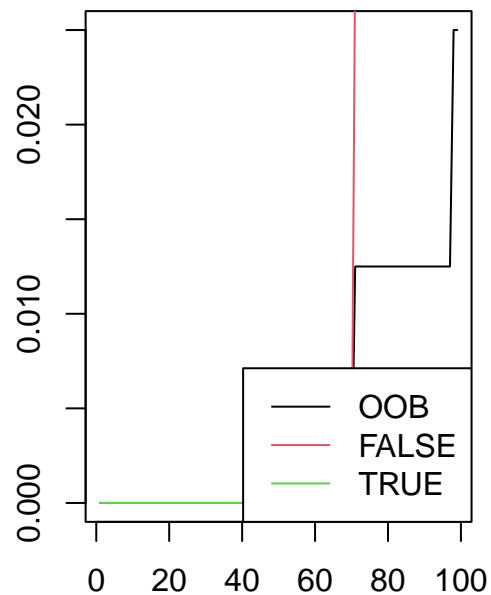
```
plot(model_mc_cv, type = "cv", stat = "oob", main = "CV oob error")
```

```
plot(model_mc_cv, type = "model", stat = "oob", main = "Model oob error")
```

CV oob error

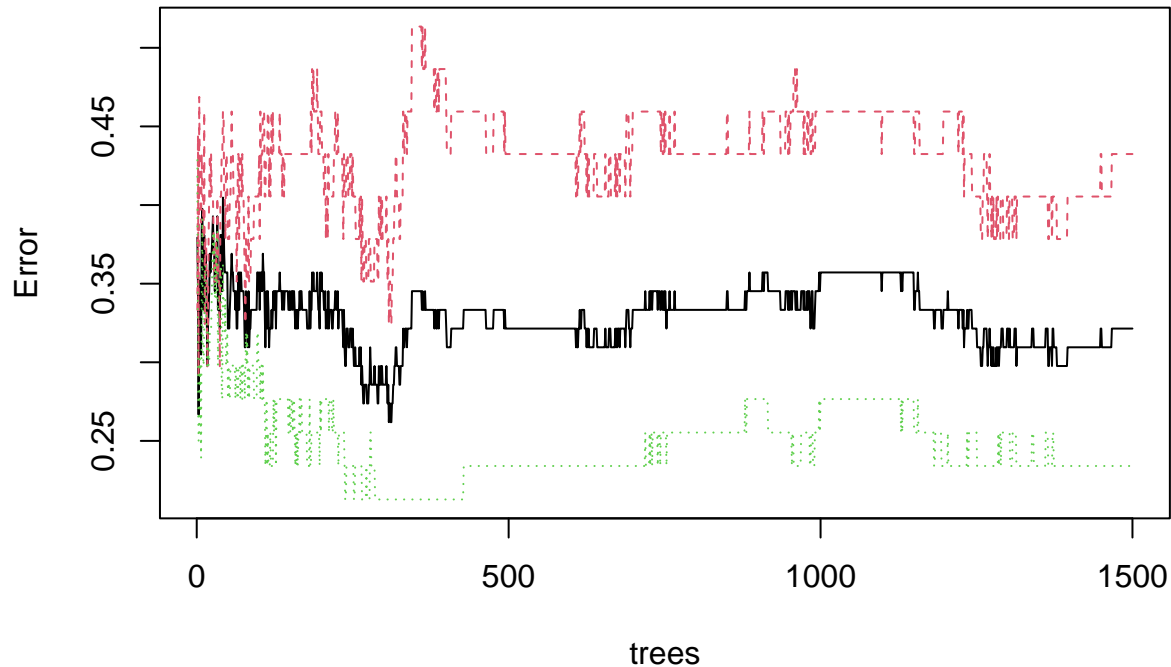


Model oob error



```
plot(model_mc)
```

model_mc



Test the model

Making predictions

#The predict() function in R is used to predict the values based on the input data.

```
predictions_mc <- predict(model_mc, test.data_mc)
```

```

# assign test.data to a new object, so that we can make changes
result_mc <- test.data_mc
#add the new variable of predictions to the result object
result_mc <- cbind(result_mc, predictions_mc)
#add the results to a data frame containing test data and the prediction
result_mc <- cbind(g[row.names(result_mc), ], predictions_mc)

```

Visualizations

```

conf_matrix_mc <-
  confusionMatrix(result_mc$predictions_mc, reference = result_mc$Eim_M)

print(conf_matrix_mc)

```

```

## Confusion Matrix and Statistics
##
##              Reference
## Prediction FALSE TRUE
##      FALSE      5    5
##      TRUE       4   18
##
##              Accuracy : 0.7188
##              95% CI : (0.5325, 0.8625)
##      No Information Rate : 0.7188
##      P-Value [Acc > NIR] : 0.5886
##
##              Kappa : 0.3271
##
##  Mcnemar's Test P-Value : 1.0000
##
##              Sensitivity : 0.5556
##              Specificity : 0.7826
##              Pos Pred Value : 0.5000
##              Neg Pred Value : 0.8182
##              Prevalence : 0.2812
##              Detection Rate : 0.1562
##      Detection Prevalence : 0.3125
##              Balanced Accuracy : 0.6691
##
##      'Positive' Class : FALSE
##

```

```
conf_matrix_mc$table
```

```

##              Reference
## Prediction FALSE TRUE
##      FALSE      5    5
##      TRUE       4   18

```

```

plt <- as.data.frame(conf_matrix_mc$table)
plt$Prediction <- factor(plt$Prediction, levels=rev(levels(plt$Prediction)))

ggplot(plt, aes(x = Prediction, y = Reference, fill= Freq)) +
  geom_tile() + geom_text(aes(label=Freq)) +
  scale_fill_gradient(low="white", high="darkturquoise") +
  labs(x = "Predictions", y = "Reference")

```