

Build Survival Model: Random Survival Forest

Mingcheng Hu

Table of contents

Load Data	2
Random Survival Forest (RSF)	4
Variable Selection	4
Cross Validation to Select the Best Number of Features	5
Model Fitting	6

```
library(tidyverse)
library(survival)
library(randomForestSRC)
library(caret)
library(survcomp)
library(parallel)
library(doParallel)
library(mcprogress) # wrap mclapply with progress bar.
library(kableExtra) # include knitr automatically

source("/work/users/y/u/youkias/BIOS-Material/BIOS992/utis/csv_utis.r")
# * Don't use setwd() for Quarto documents!
# setwd("/work/users/y/u/youkias/BIOS-Material/BIOS992/data")

adjust_type <- ifelse(exists("params"), params$adjust_type, "minimal") #
  ↳ options: "minimal", "partial", "full"
impute_type <- ifelse(exists("params"), params$impute_type, "unimputed") #
  ↳ options: "unimputed", "imputed"
include_statin <- ifelse(exists("params"), params$include_statin, "no") #
  ↳ options: "yes", "no"

n_folds <- 10
set.seed(1234)
```

```
# string of parameters
adjust_type_str <- switch(adjust_type,
  minimal = "minimal",
  partial = "partial",
  full = "full"
)
print(paste0("Model Adjustment Type: ", adjust_type_str))
```

```
[1] "Model Adjustment Type: minimal"
```

```
impute_type_str <- switch(impute_type,
  unimputed = "unimputed",
  imputed = "imputed"
)
print(paste0("Data Imputation Type: ", impute_type_str))
```

```
[1] "Data Imputation Type: unimputed"
```

Load Data

```
if (include_statin == "yes") {
  data_train <-
  ↪ read.csv(paste0("/work/users/y/u/youkias/BIOS-Material/BIOS992/data/train_data_",
  ↪ impute_type_str, "_statin.csv"),
    header = TRUE
  )
} else {
  data_train <-
  ↪ read.csv(paste0("/work/users/y/u/youkias/BIOS-Material/BIOS992/data/train_data_",
  ↪ impute_type_str, ".csv"),
    header = TRUE
  )
}

data_train <- data_train[, -1] # the first column is the index generated by
  ↪ sklearn
(dim(data_train))
```

```
[1] 28127    100
```

```
data <- select_subset(data_train, type = adjust_type)
(dim(data))
```

```
[1] 28127    48
```

```
colnames(data)
```

```
[1] "event"                "time"
[3] "HRV_SD1"              "HRV_SD2"
[5] "HRV_SD1SD2"           "HRV_S"
[7] "HRV_CSI"              "HRV_CVI"
[9] "HRV_CSI_Modified"     "HRV_PIP"
[11] "HRV_IALS"             "HRV_PSS"
[13] "HRV_PAS"              "HRV_GI"
[15] "HRV_SI"               "HRV_AI"
[17] "HRV_PI"               "HRV_C1d"
[19] "HRV_C1a"              "HRV_SD1d"
[21] "HRV_SD1a"             "HRV_C2d"
[23] "HRV_C2a"              "HRV_SD2d"
[25] "HRV_SD2a"             "HRV_Cd"
[27] "HRV_Ca"               "HRV_SDNNd"
[29] "HRV_SDNNa"            "HRV_ApEn"
[31] "HRV_ShanEn"           "HRV_FuzzyEn"
[33] "HRV_MSEn"             "HRV_CMSEn"
[35] "HRV_RCMSEn"           "HRV_CD"
[37] "HRV_HFD"              "HRV_KFD"
[39] "HRV_LZC"              "HRV_DFA_alpha1"
[41] "HRV_MFDFA_alpha1_Width" "HRV_MFDFA_alpha1_Peak"
[43] "HRV_MFDFA_alpha1_Mean" "HRV_MFDFA_alpha1_Max"
[45] "HRV_MFDFA_alpha1_Delta" "HRV_MFDFA_alpha1_Asymmetry"
[47] "HRV_MFDFA_alpha1_Fluctuation" "HRV_MFDFA_alpha1_Increment"
```

```
data <- tibble::as_tibble(data)
```

```
# * It is very hard to compare the HR as different predictors are on
  ↳ different magnitudes, so we need to normalize them.
time_col <- data$time
```

```

event_col <- data$event
data <- data %>%
  select(-c(time, event)) %>%
  mutate(across(where(is.numeric), scale)) %>%
  mutate(
    time = time_col,
    event = event_col
  )

```

Note now the interpretation of HR is different! For example, if HR=1.16 for the predictor in the univariate model fitted using scaled data, it means that each standard deviation increase is associated with 16% higher risk of event.

```
# For RSF model, we don't need to exclude the missing values
```

Random Survival Forest (RSF)

Variable Selection

The `method` argument can be set to `vh` instead for variable hunting, which should be used for problems where the number of variables is substantially larger than the sample size.

```

n_cores <- parallel::detectCores() - 1
cl <- makeCluster(n_cores)
registerDoParallel(cl)

rsf_var_select <- var.select.rfsrc(Surv(time, event) ~ .,
  data = data,
  method = "md",
  seed = 1234,
  ntree = 200,
  parallel = TRUE
) # minimal depth variable selection

stopCluster(cl)

```

```
vars_ranked <- rsf_var_select$topvars
```

Cross Validation to Select the Best Number of Features

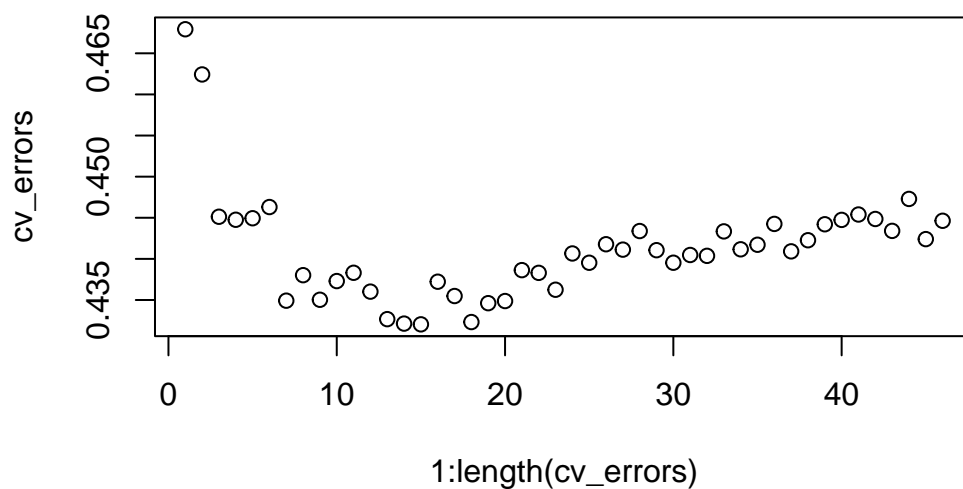
We will use 10-fold cross validation to select the best number of features used in the model.

```
set.seed(1234)
folds <- createFolds(data$event, k = n_folds) # return indices of folds

cv_errors <- pmclapply(seq(1, length(vars_ranked), by = 1),
  ↪ function(num_vars) {
    selected_vars <- vars_ranked[1:num_vars]

    fold_errors <- sapply(folds, function(fold_idx) {
      train_data_fold <- data[-fold_idx, c("time", "event", selected_vars)]
      val_data_fold <- data[fold_idx, c("time", "event", selected_vars)]
      model <- rfsrc.fast(Surv(time, event) ~ .,
        data = train_data_fold,
        ntree = 200,
        forest = TRUE
      )
      pred <- predict(model,
        newdata = val_data_fold,
        na.action = "na.impute" # * There may be missing values in the
        ↪ dataset
      )$predicted # define pred has attributes survival(sample_size*time)
      ↪ and predicted(sample_size) for risk
      # Use C-index to measure the performance of the model
      1 - concordance.index(
        pred, # pass risk prediction for first argument
        val_data_fold$time,
        val_data_fold$event
      )$c.index
    })
    mean(fold_errors)
  }, title = "Cross Validation to Select the Best Number of Features")
```

```
cv_errors <- as.numeric(cv_errors)
plot(1:length(cv_errors), cv_errors)
```



```
best_num_vars <- which.min(cv_errors)
vars_selected <- vars_ranked[1:best_num_vars]
```

```
print(paste0("The best number of features to retain is ", best_num_vars))
```

```
[1] "The best number of features to retain is 15"
```

Model Fitting

```
data_selected <- data[, c("time", "event", vars_selected)]

# Before formally fitting the model, we can tune the hyperparameters to find:
# 1. optimal mtry (possible split at each node)
# 2. optimal nodesize (minimum size of terminal nodes)
rsf_tuned <- tune.rfsrc(
  Surv(time, event) ~ .,
  data = data_selected,
)
```

```
rsf_model <- rfsrc(Surv(time, event) ~ .,  
  data = data_selected,  
  ntree = 500,  
  mtry = rsf_tuned$best.mtry,  
  nodesize = rsf_tuned$best.nodesize  
)
```

```
# We also fit the full model  
rsf_tuned_full <- tune.rfsrc(  
  Surv(time, event) ~ .,  
  data = data,  
)  
  
rsf_model_full <- rfsrc(Surv(time, event) ~ .,  
  data = data,  
  ntree = 1000,  
  mtry = rsf_tuned_full$best.mtry,  
  nodesize = rsf_tuned_full$best.nodesize  
)
```

```
# plot.rfsrc?  
# plot.variable.rfsrc?  
# print.rfsrc?
```