

Evaluate testing data (regression)

Andrew Chang

2018-11-07

Contents

0. Load Data	1
1. Scores	1
1. Scores	2
1.1 uncalibrated average	2
3. Important Features	4
4. Hyper-parameters	8

Labels: os_time

0. Load Data

```
library(tidyverse)
library(ggrepel)
library(caret)
library(pROC)
source("~/ml-pipeline/ML_performanceCheck/PerformanceUtils.R")

home <- "~/ml-pipeline/tests"
project_name <- "test_regression_1"
training_data <- "test_regr_surv_tcga_brca.csv"
label_name <- "os_time"

## output path for analyzed data
outfile <- paste0(home, "/results/", project_name)
```

300 of samples were used

101 of full features

10 runs, each run contains 3 CVs.

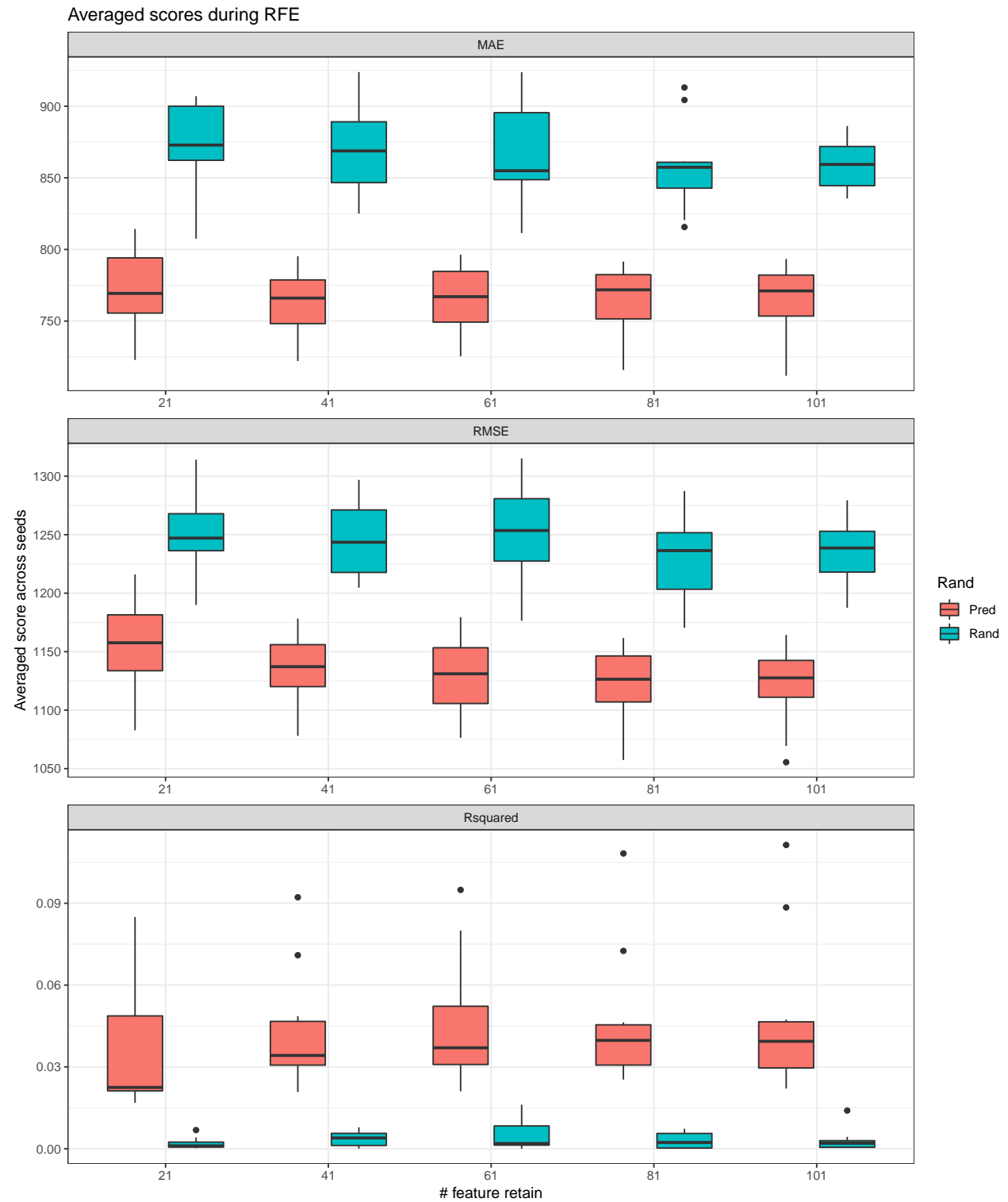
1. Scores

average = T: scores are based on pre-validation (combine predictions from all CVs per job and then calculate a single score per job).

average = F: report all the scores from each CV across entire jobs.

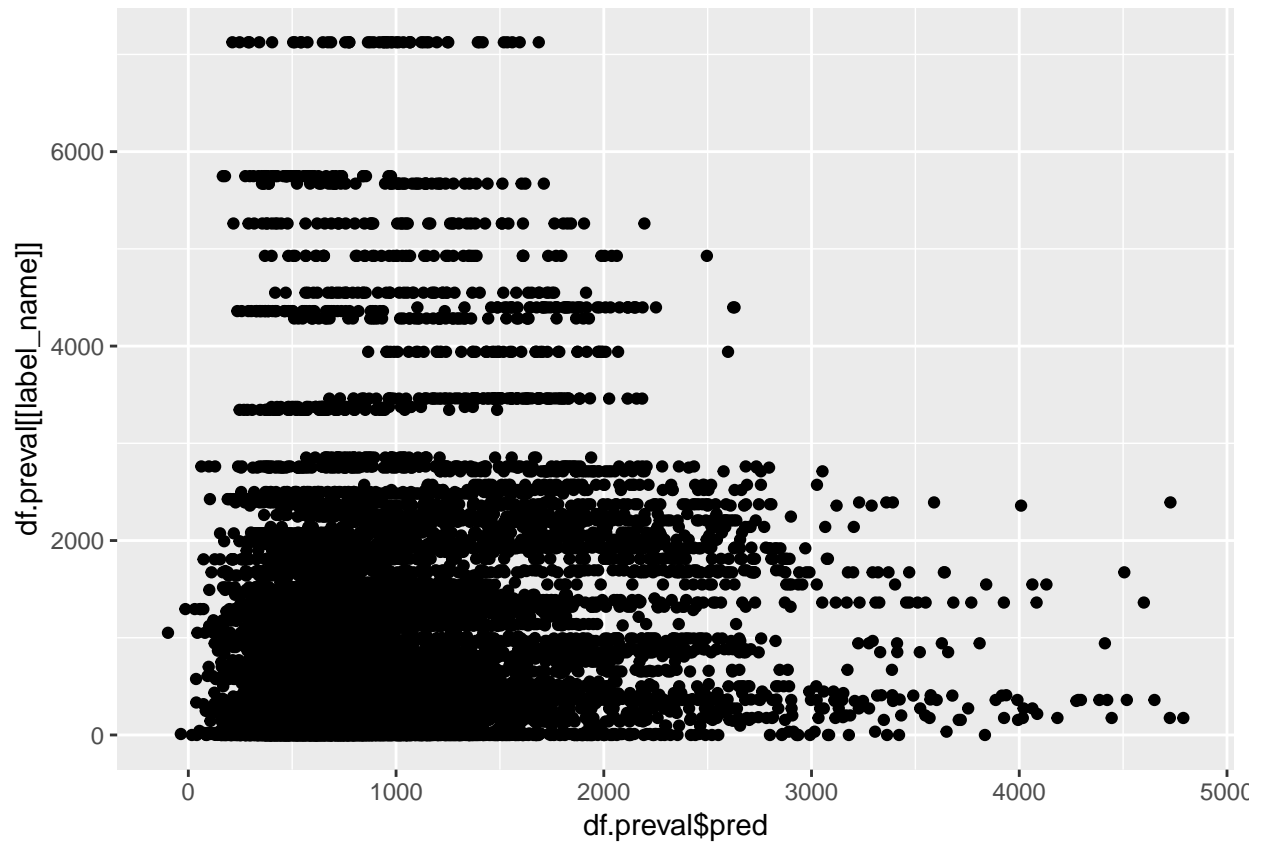
1. Scores

1.1 uncalibrated average

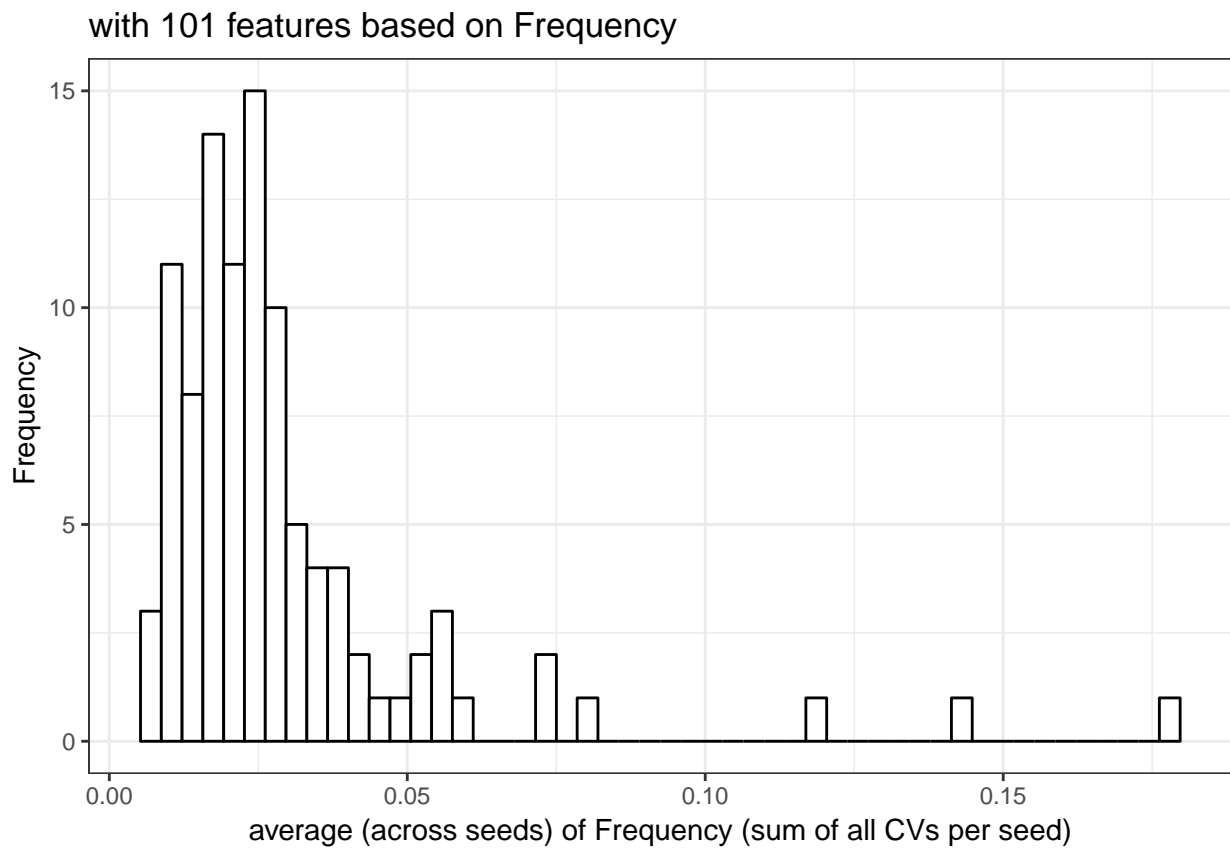


```
ggplot(data = df.preval, aes(x=df.preval$pred,  
y = df.preval[[label_name]])) +
```

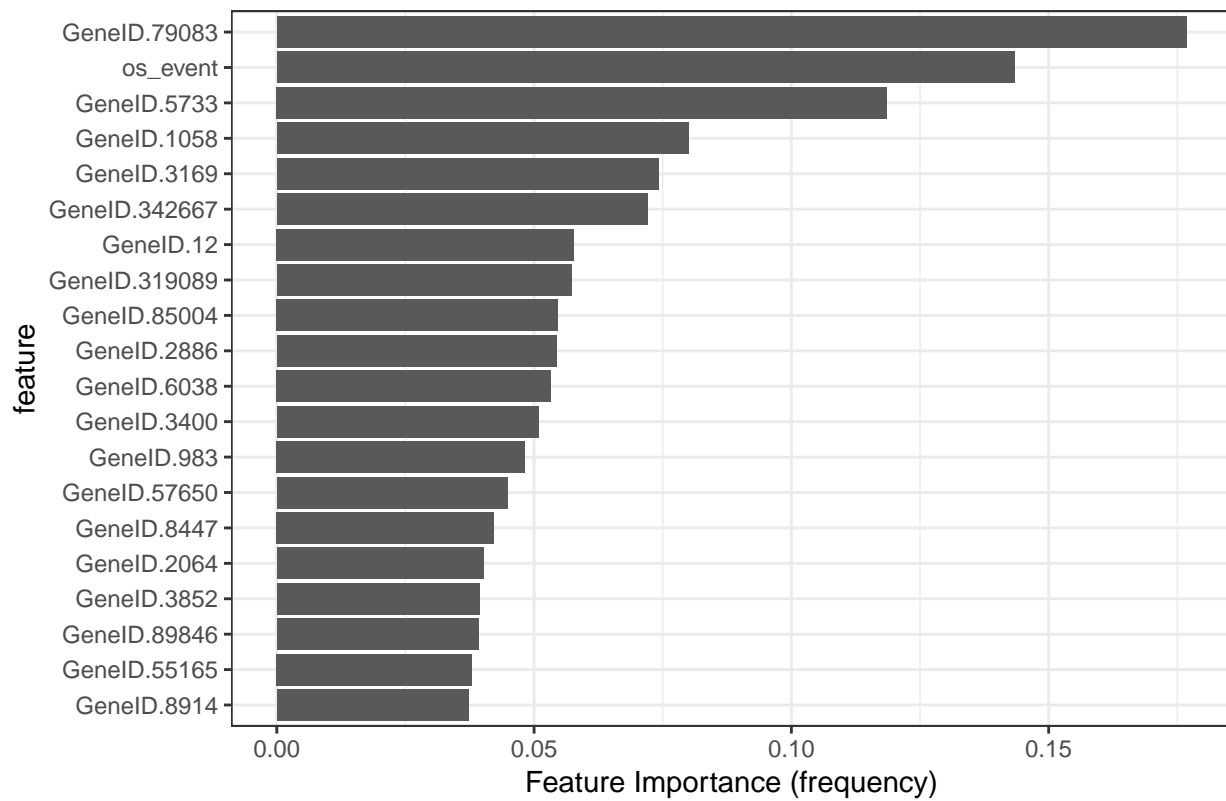
```
geom_point()
```



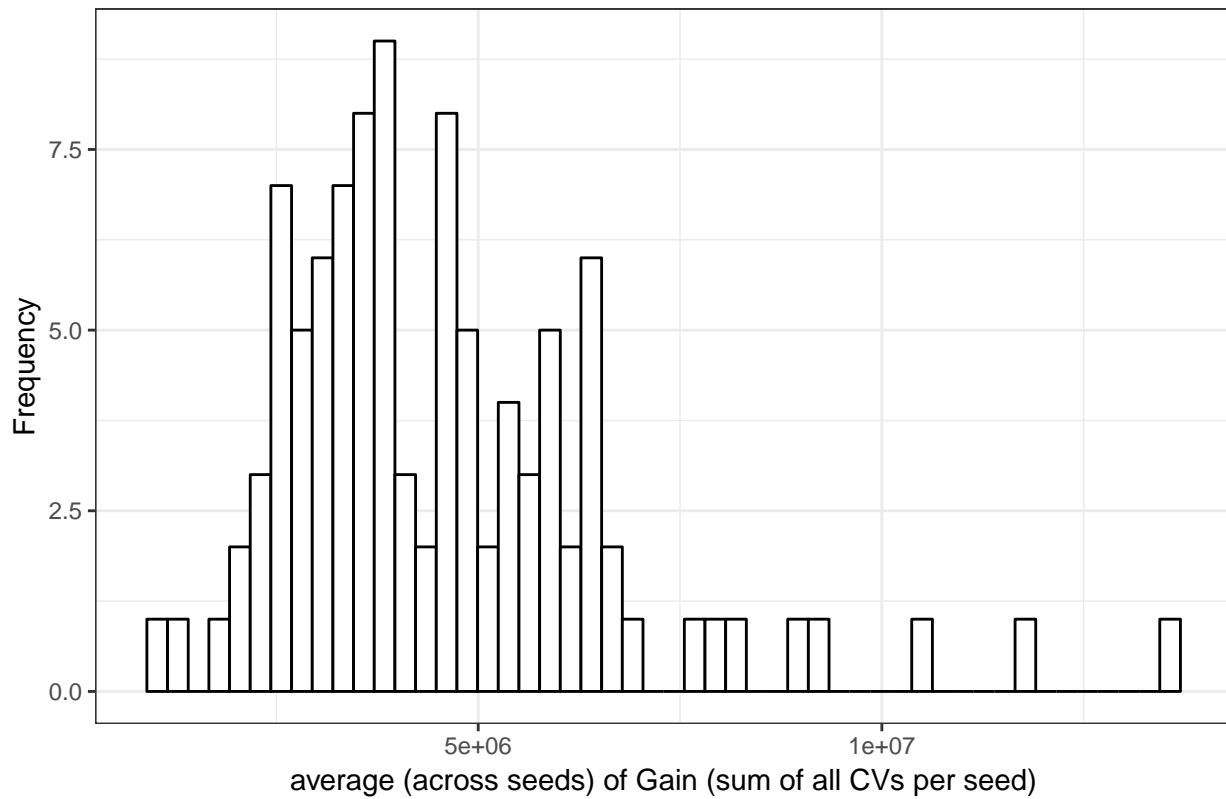
3. Important Features



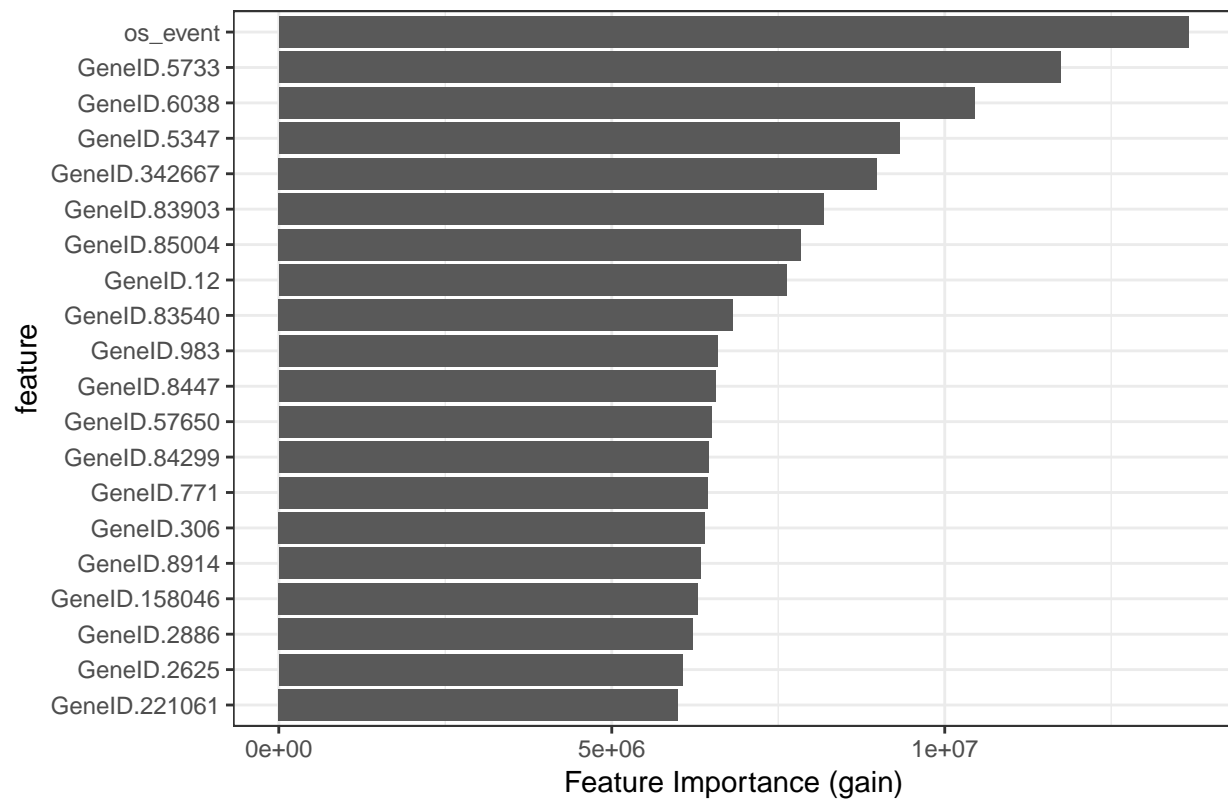
Top 20 features at 101 feature set based on Frequency

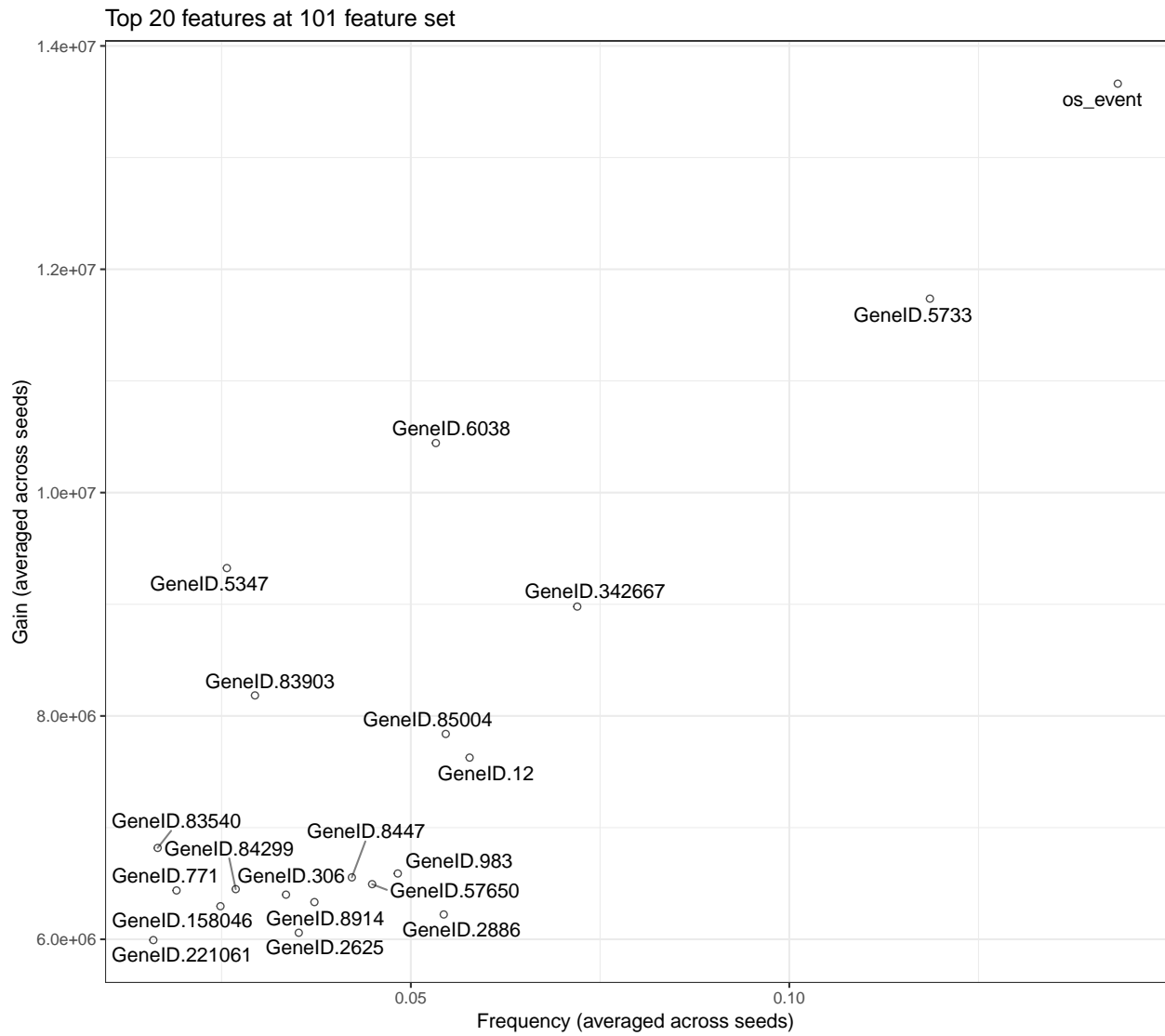


with 101 features based on Gain



Top 20 features at 101 feature set based on Gain





4. Hyper-parameters

