# Evaluate testing data (binary-class) - XGBoost

*EVE W.*

*2019-04-05*

## Contents

```
## user input
project_home <- "~/EVE/examples"
project_name <- "xgboostR_multi_1"
```

## 0. Load Data

```
## Warning: `data_frame()` is deprecated, use `tibble()`.
## This warning is displayed once per session.
```
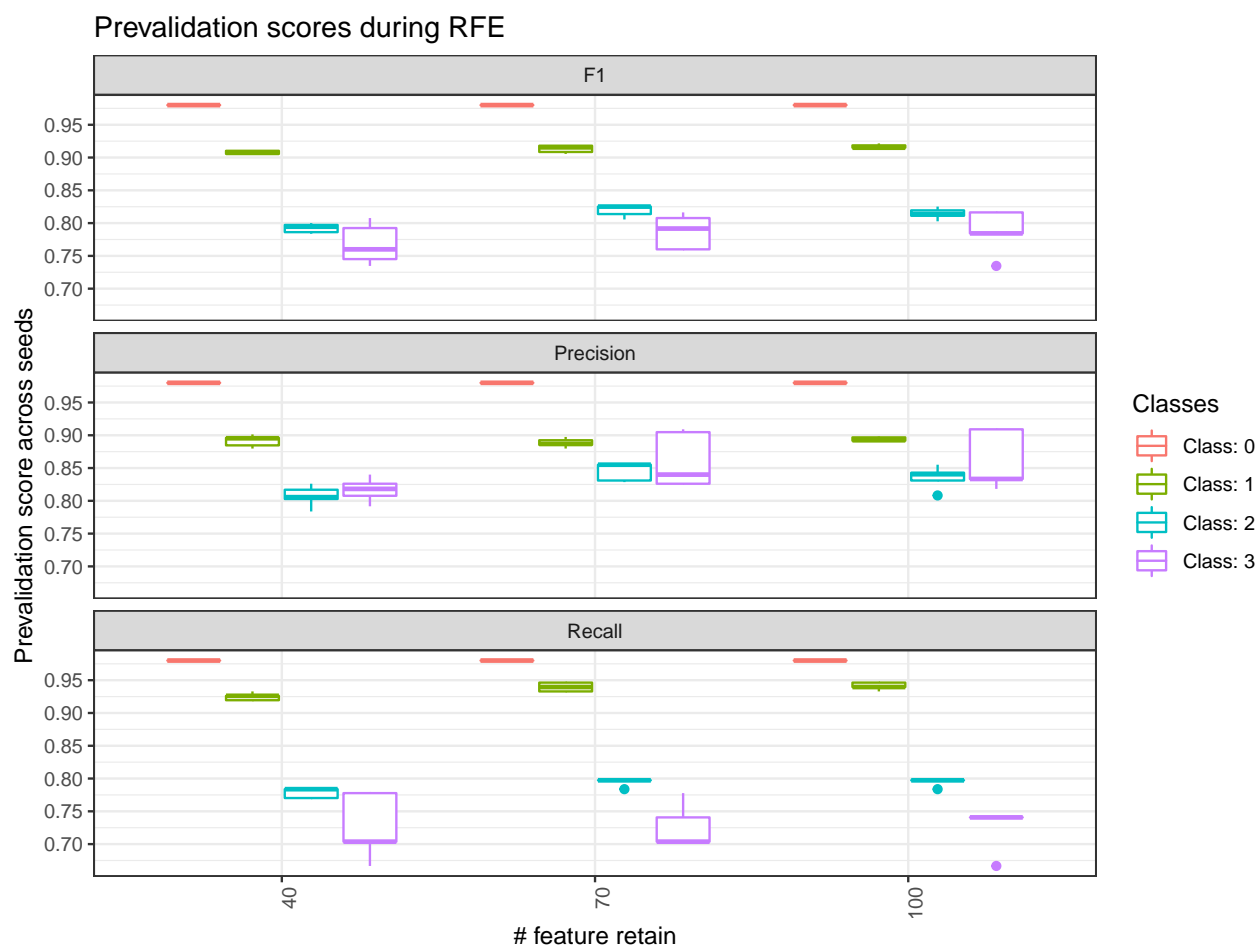
```
## Parsed with column specification:
## cols(
##   .default = col_double(),
##   Patient_ID = col_character()
## )
```

```
## See spec(...) for full column specifications.
```

```
## 300 of samples were used
```

```
## 100 of full features
```

```
## 5 runs, each run contains 300 CVs.
```

```
## Labels:
```

run with XGBoost.r evaluation metric: NA.

# 1. Scores

## 1.1 Scores per Class

### Prevalidation scores during RFE



Confusion Matrix

```
## confusion matrix at feature size = 100
## sum across 5 seeds

##            Reference
## Prediction   0    1    2    3
##          0 245    5    0    0
##          1   0  701   68   15
##          2   0   36  294   22
##          3   5    3    8   98
```
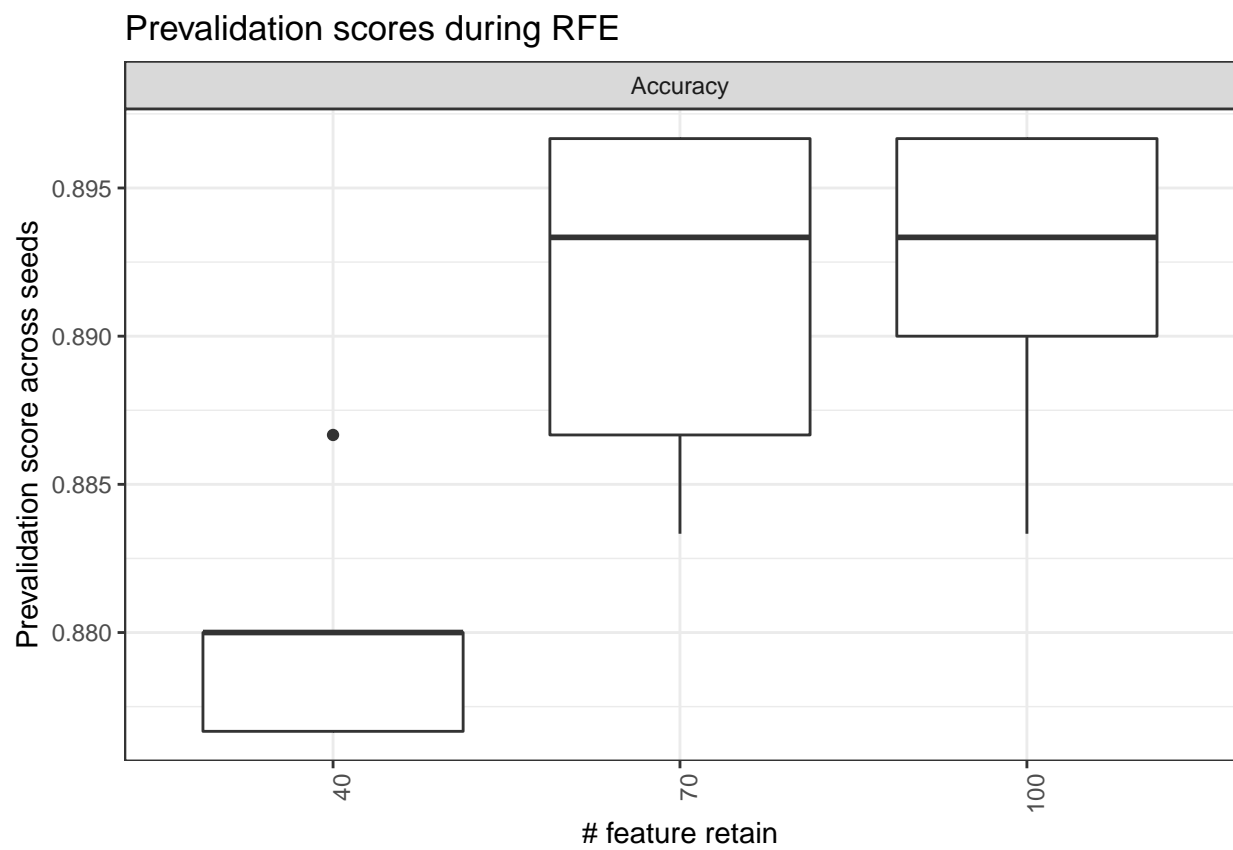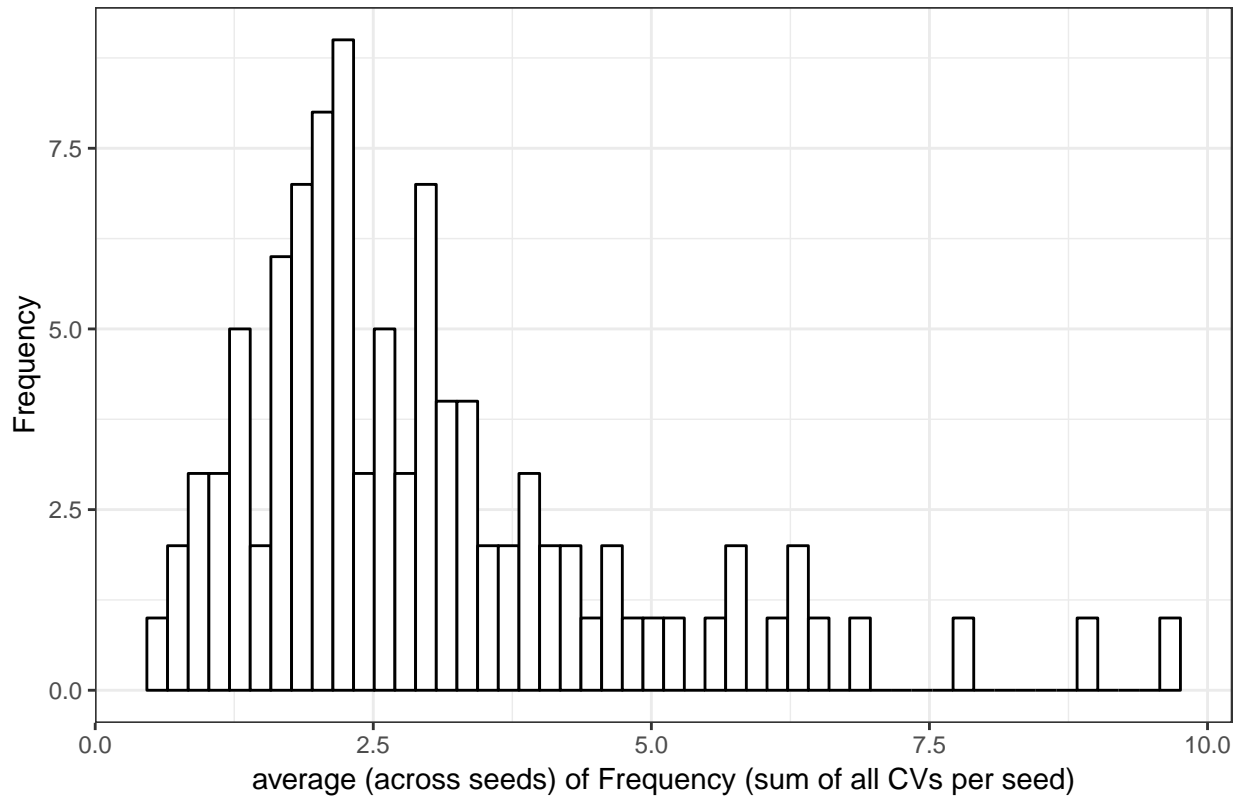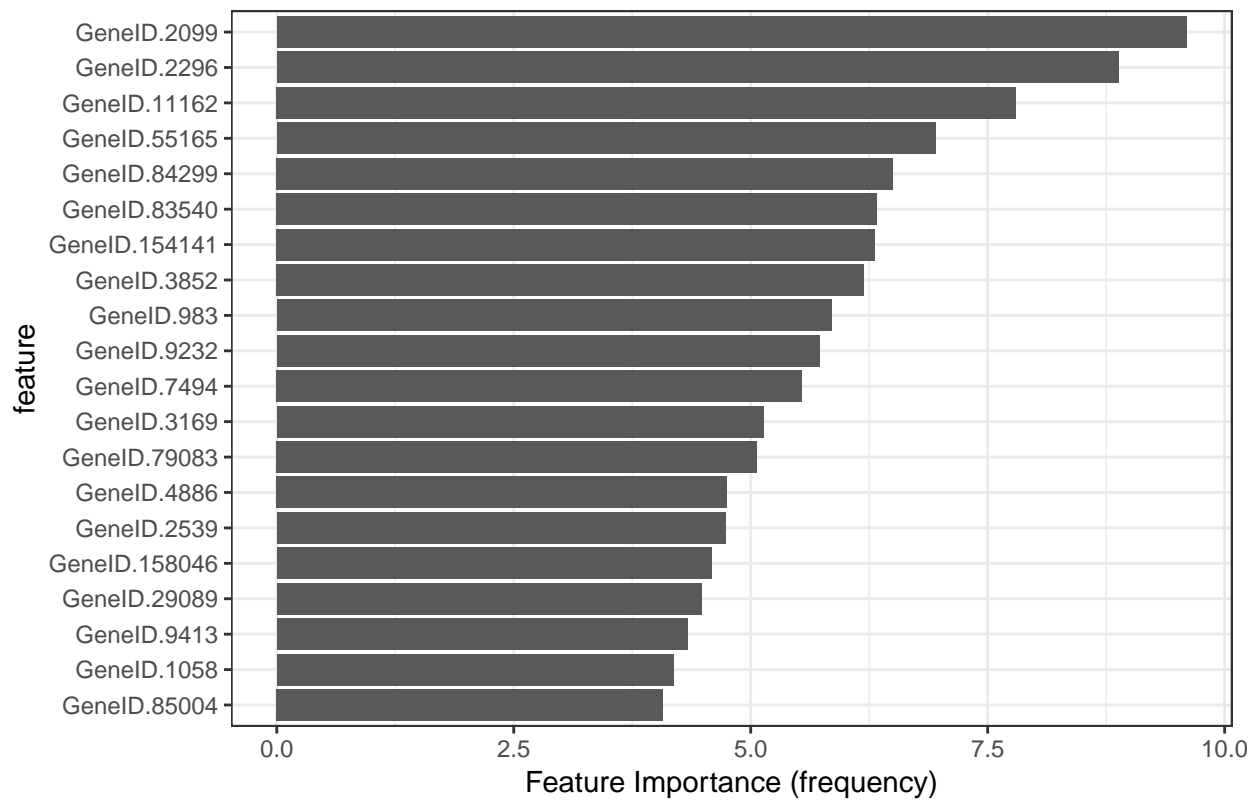
**1.2 Average score**

## Prevalidation scores during RFE



Table 1: best scores

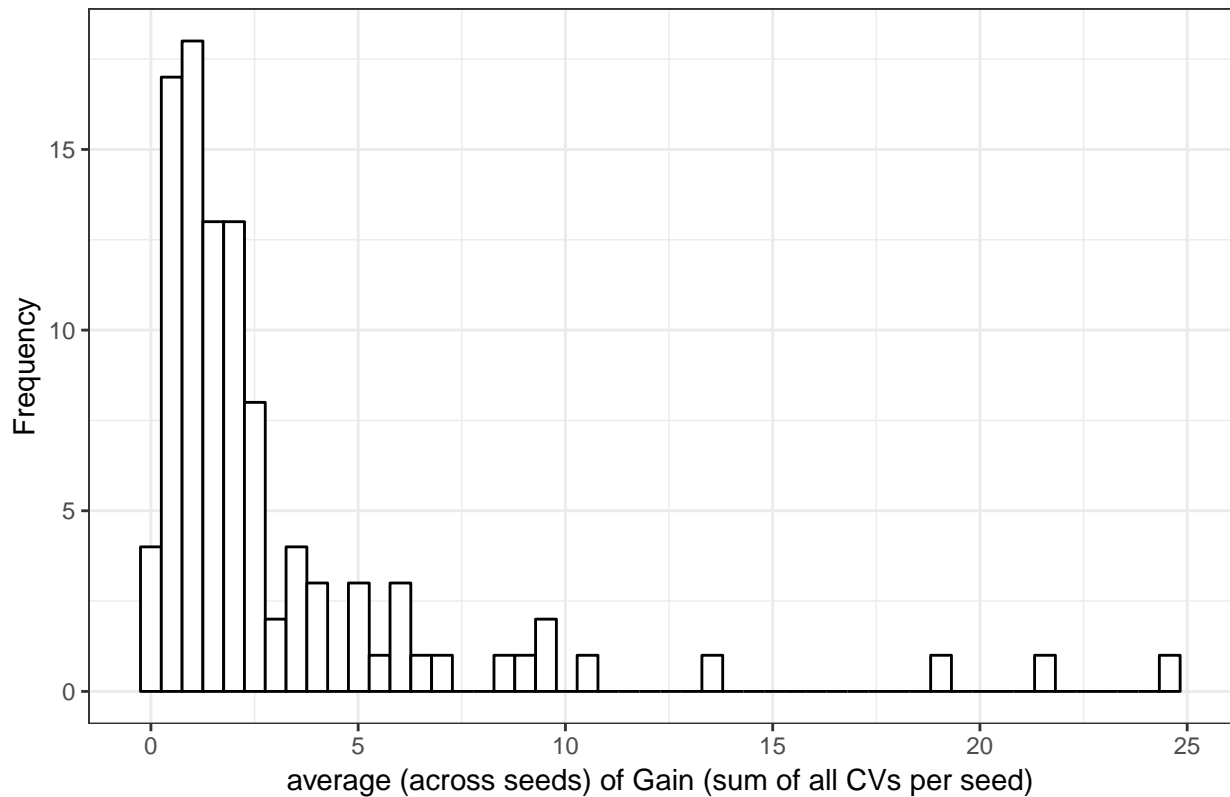| metrics | size.max | median.max | size.min | median.min |
|---|---|---|---|---|
| Accuracy | 70 | 0.893 | 40 | 0.880 |
| F1 | 100 | 0.878 | 40 | 0.861 |
| Precision | 70 | 0.893 | 40 | 0.872 |
| Recall | 100 | 0.864 | 40 | 0.847 |

## 2. Important Features

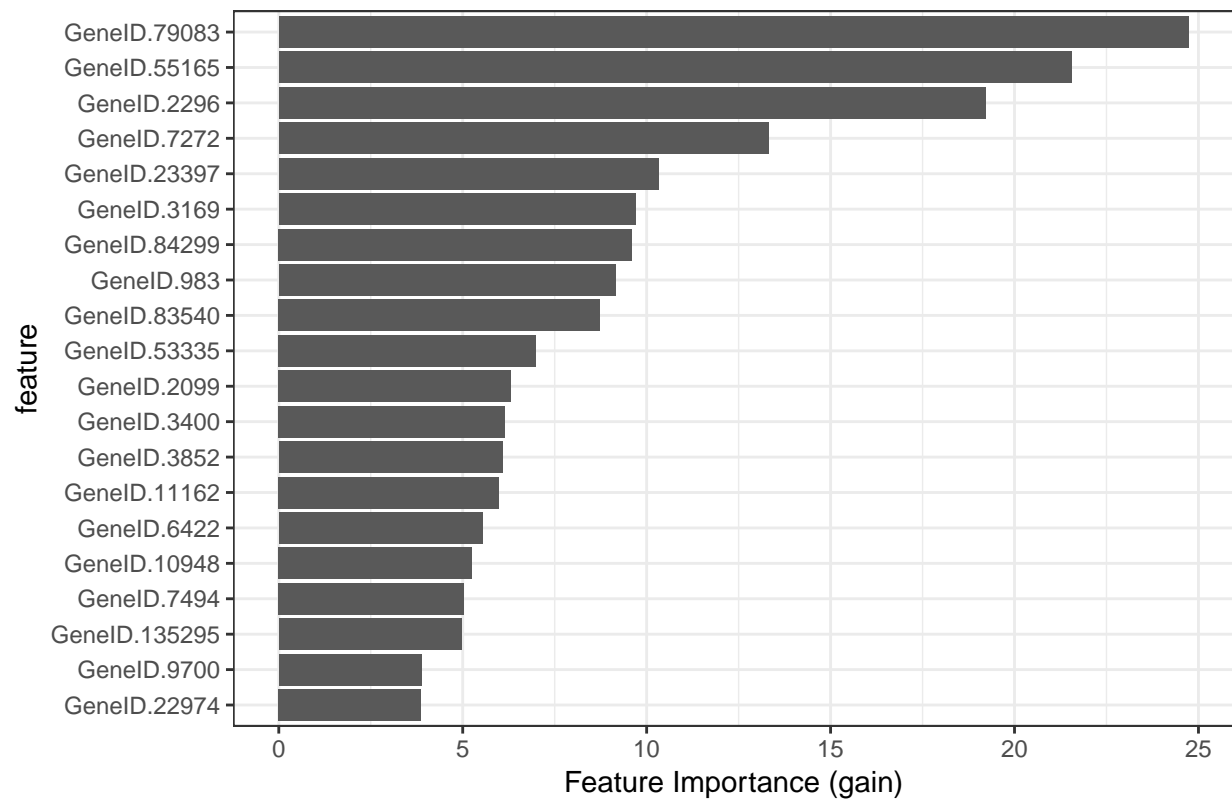with 100 features based on Frequency

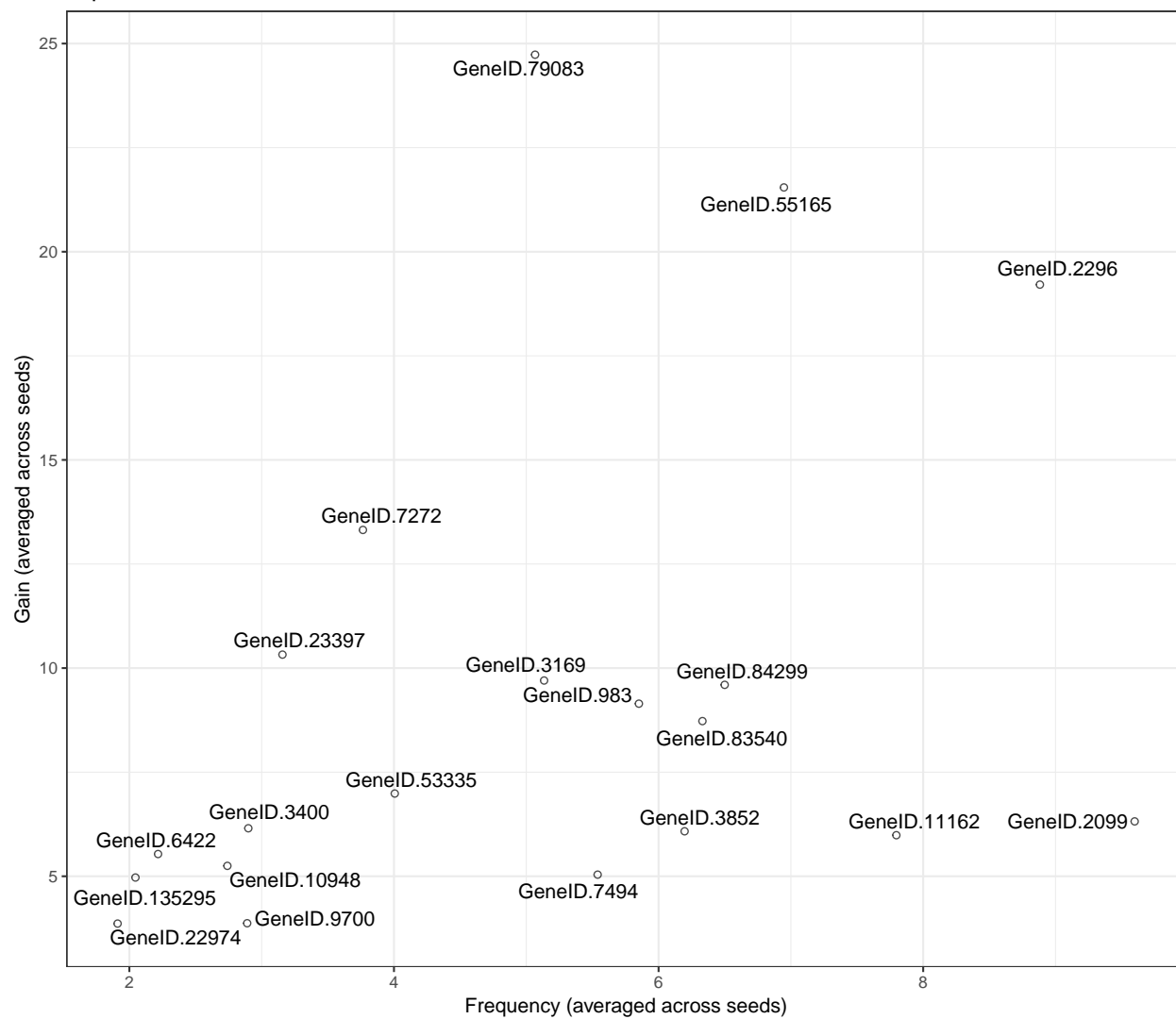Top 20 features at 100 feature set based on Frequency

with 100 features based on Gain
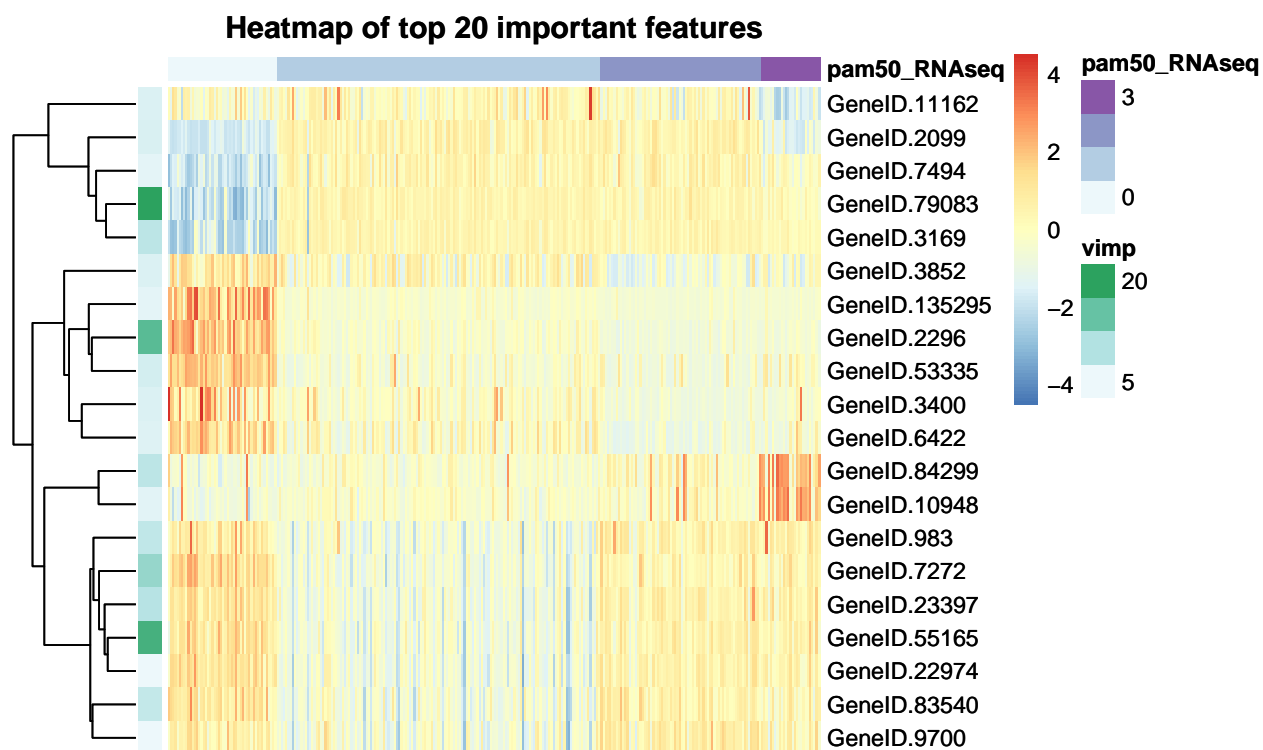
Top 20 features at 100 feature set based on Gain
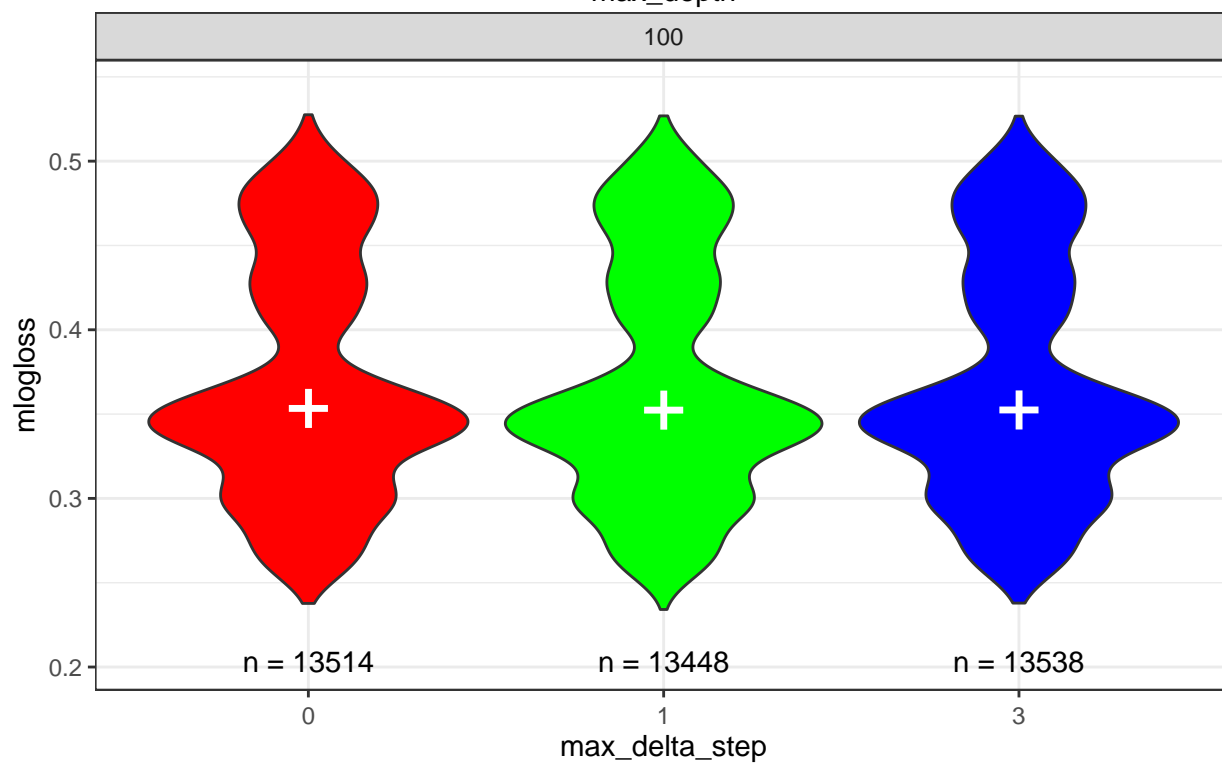
Top 20 features at 100 feature set

Heatmap of top 20 important features
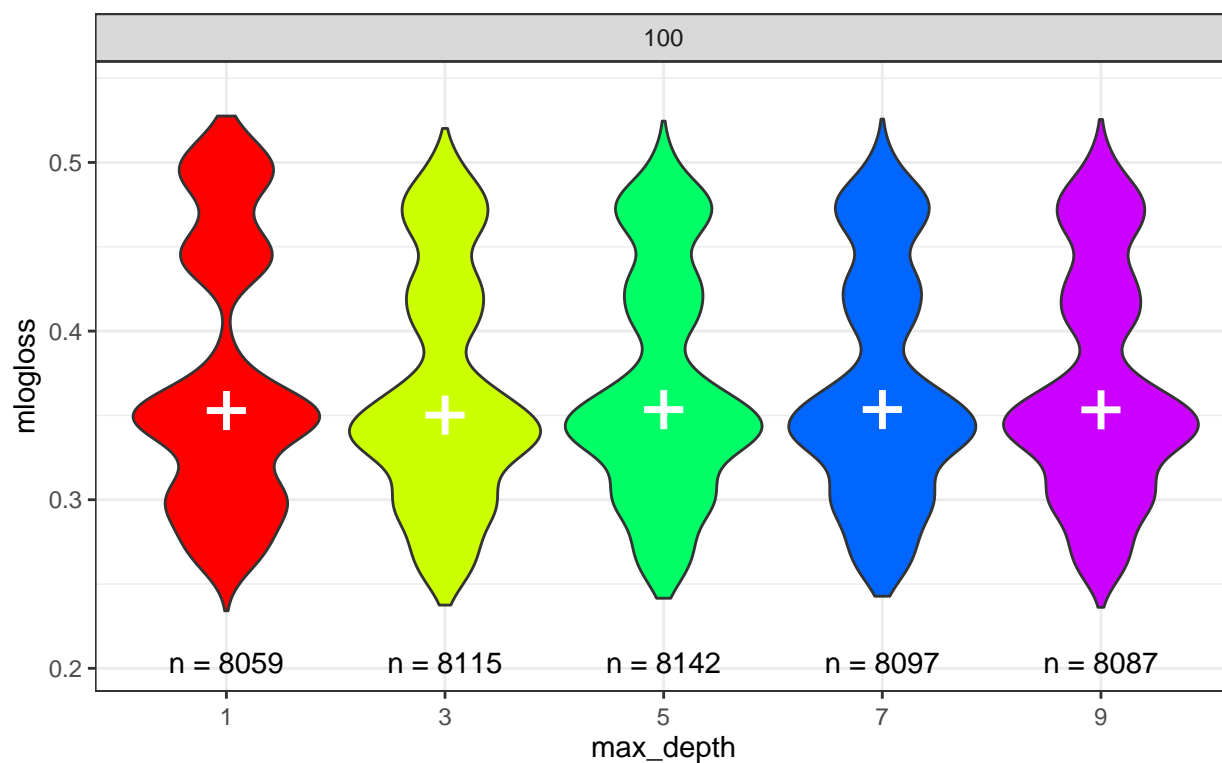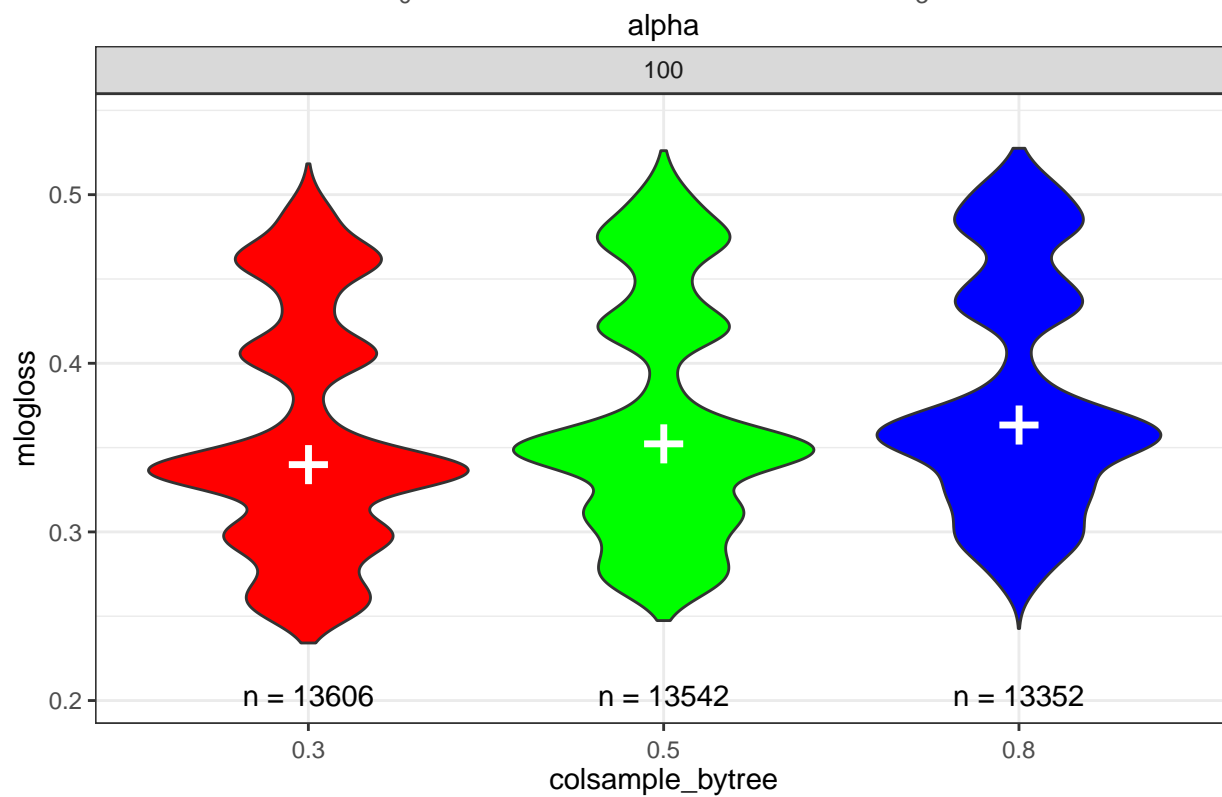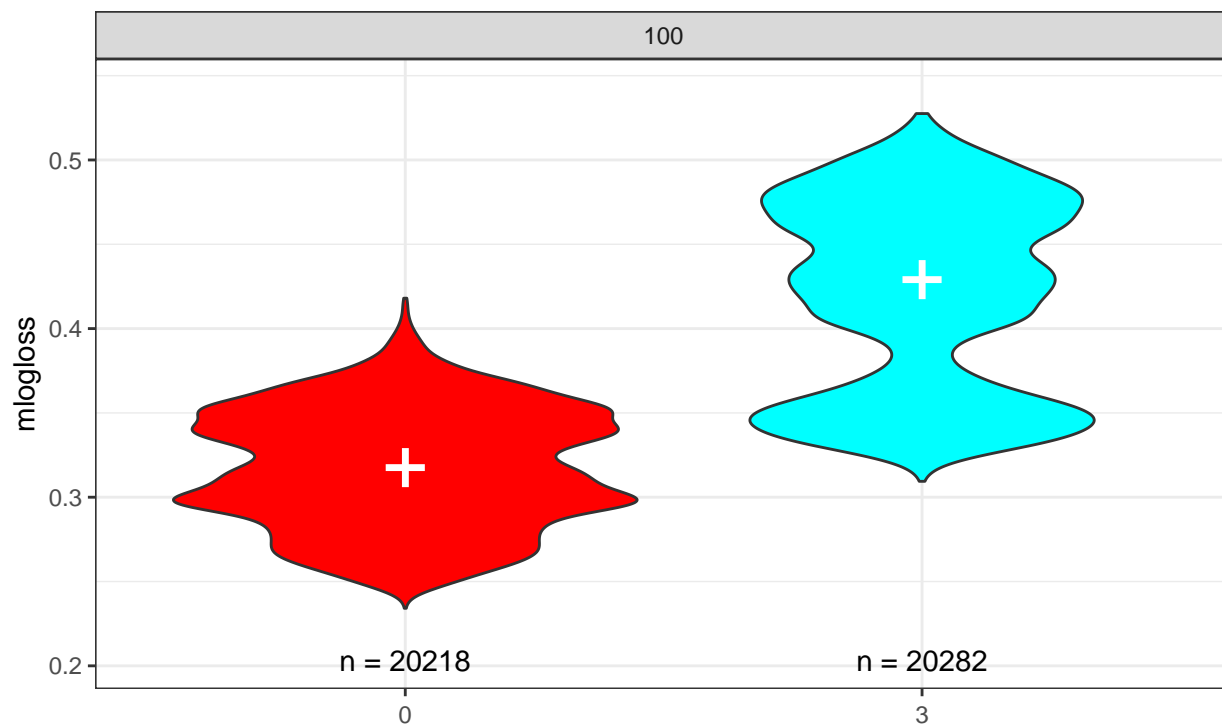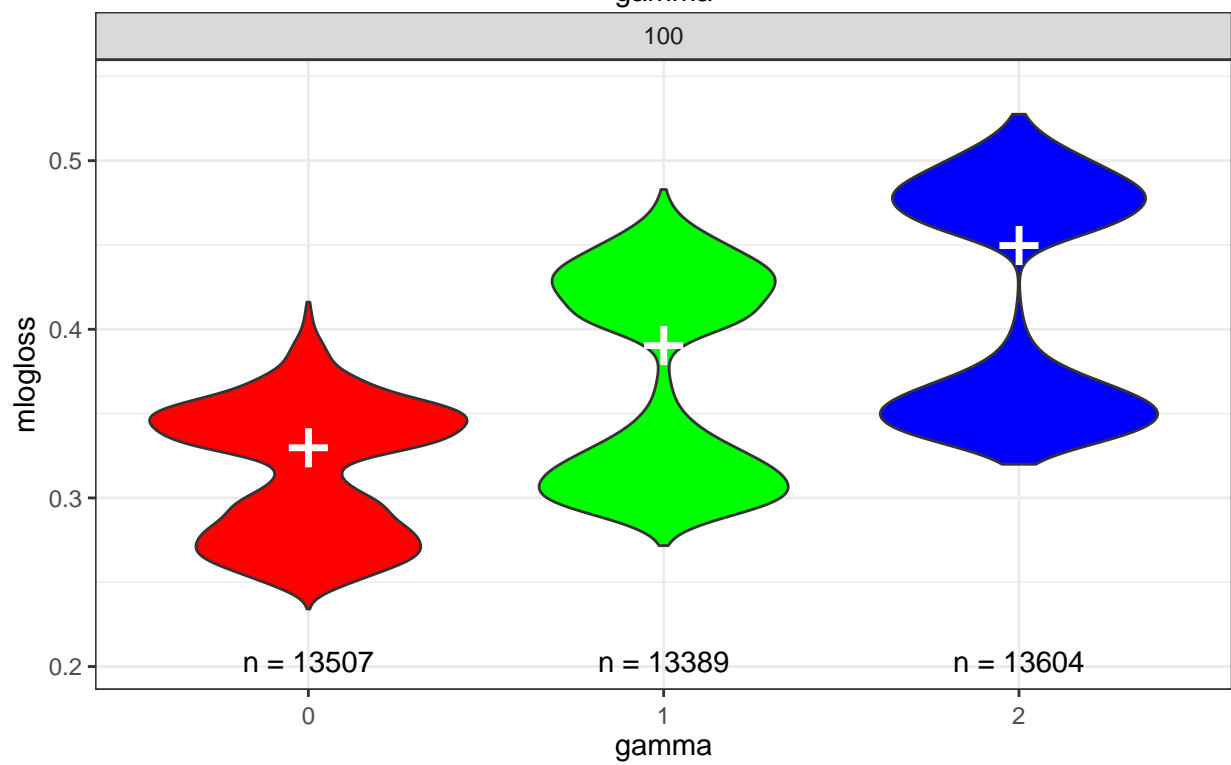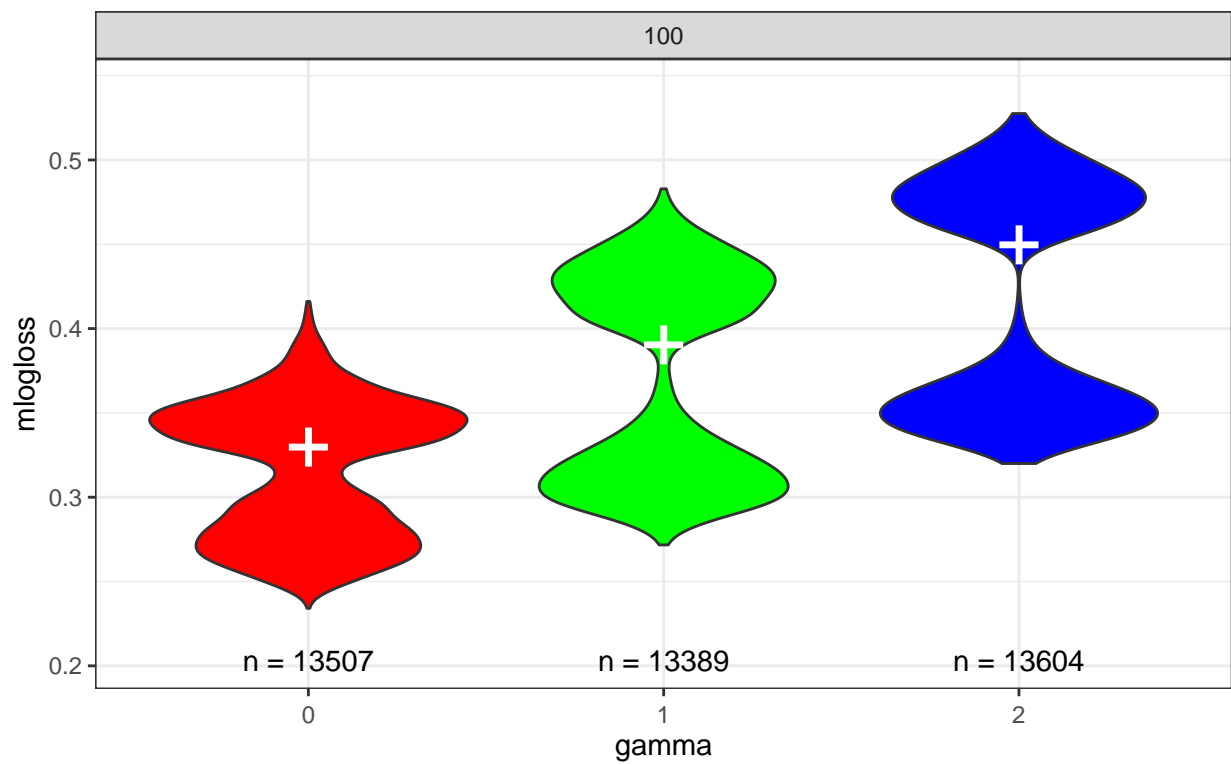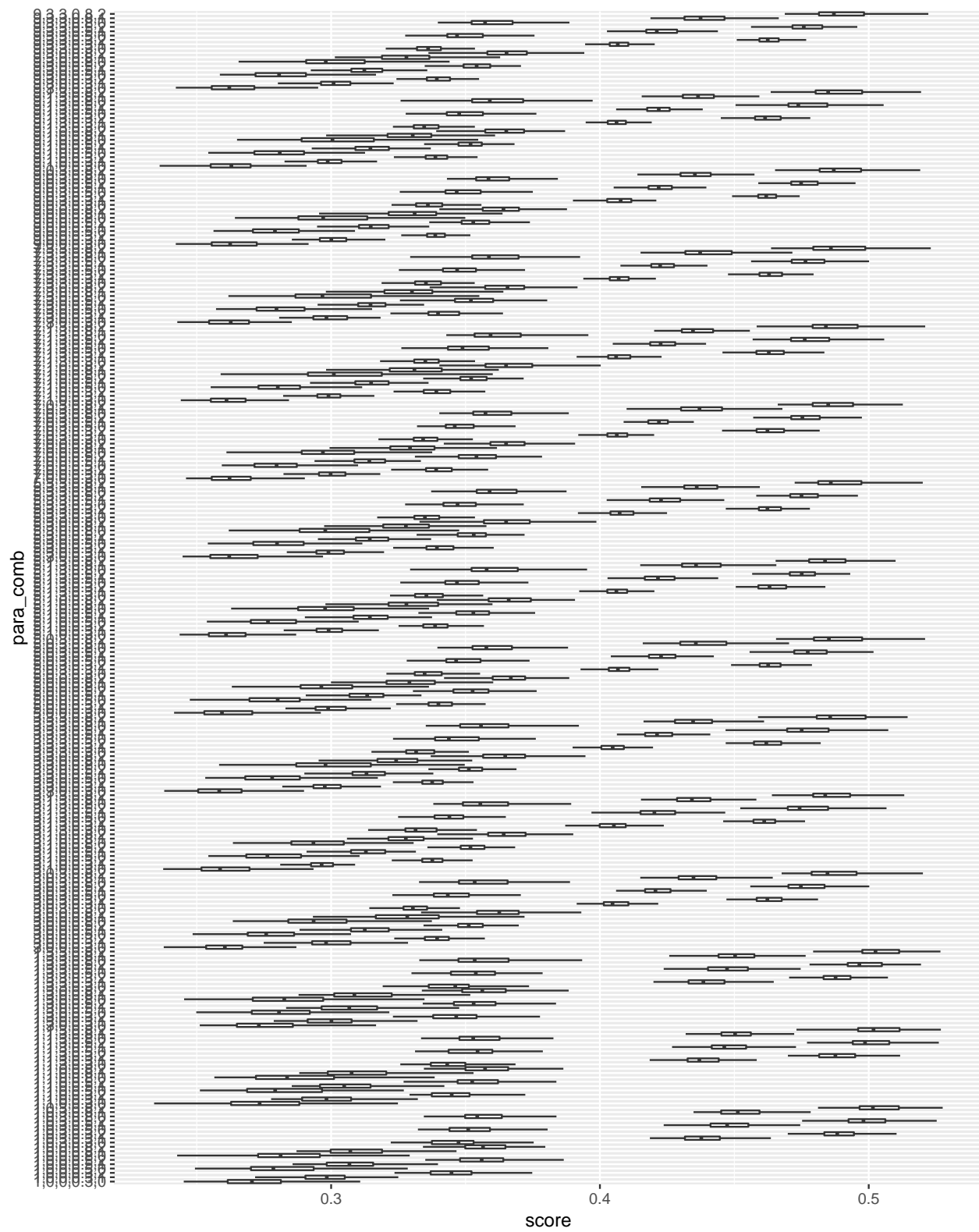
## 3. Hyper-parameters

parameter optimization file (40500 records) includes 5 seeds. Each seed generates 300 cv splits. Within each cv split, there is a 1 step RFE (at 100). So 40500 / 5 / 300 / 1 = 27 parameter combinations tried in each cv split.
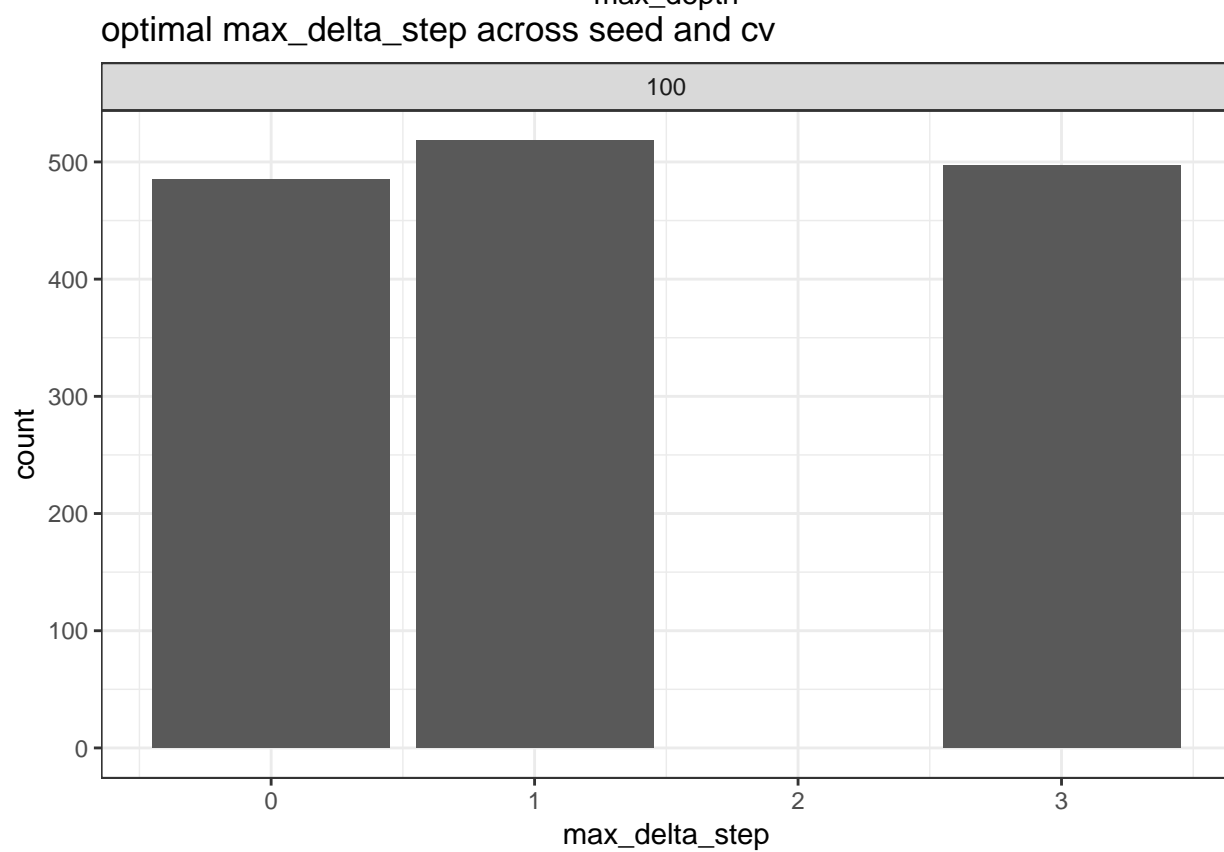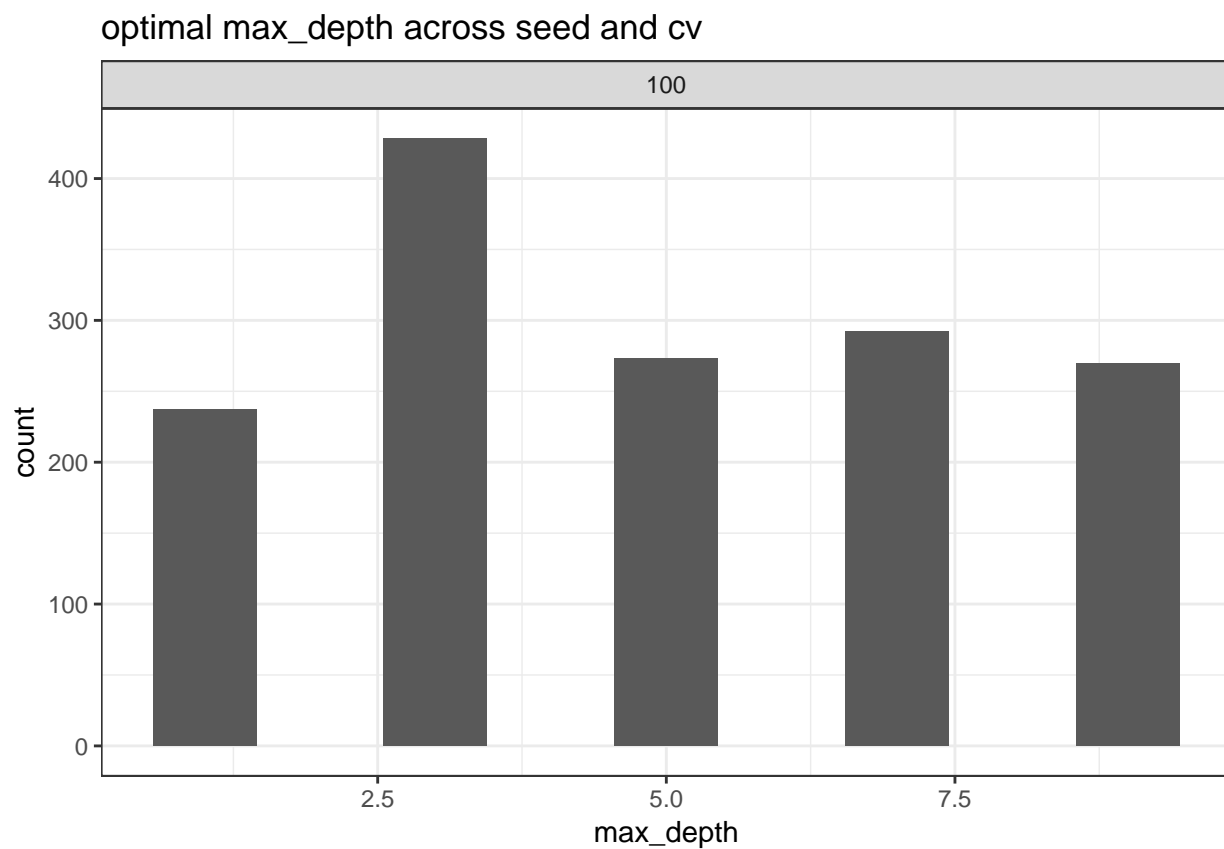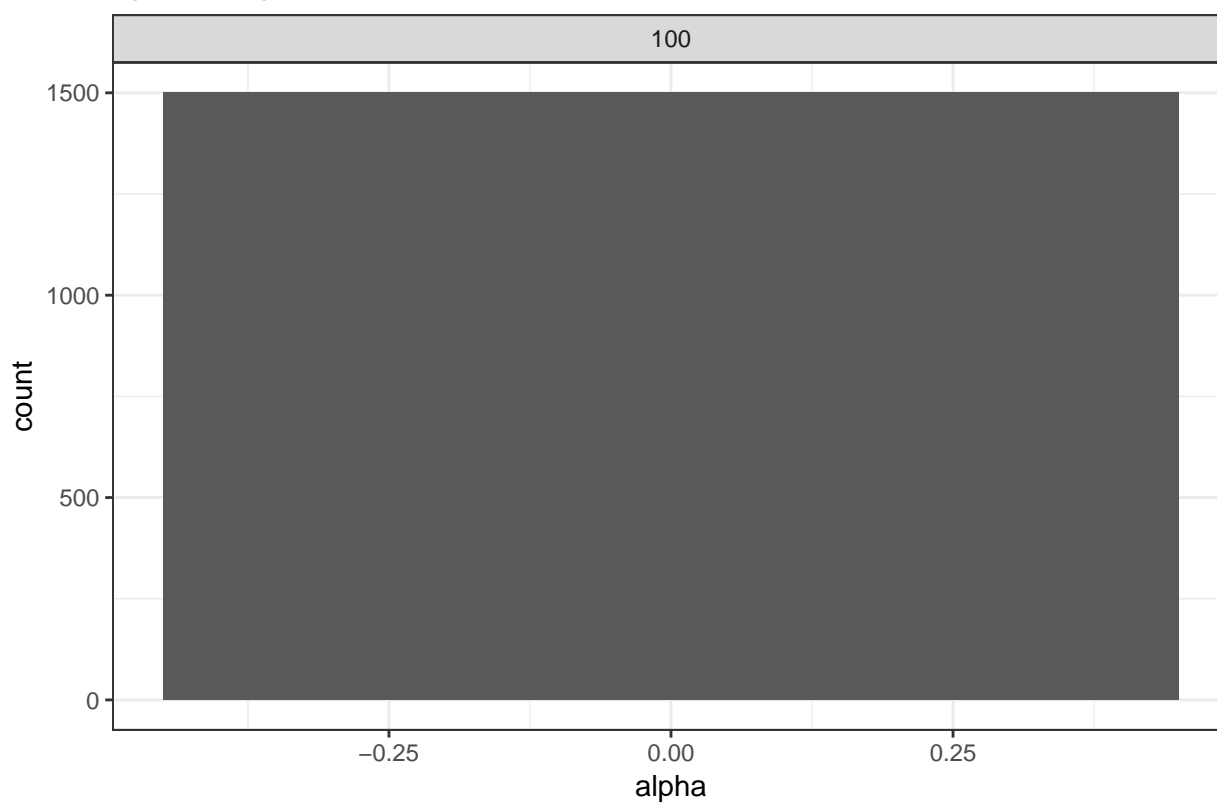
**all grid search results**

**over best parameter combo per cv**

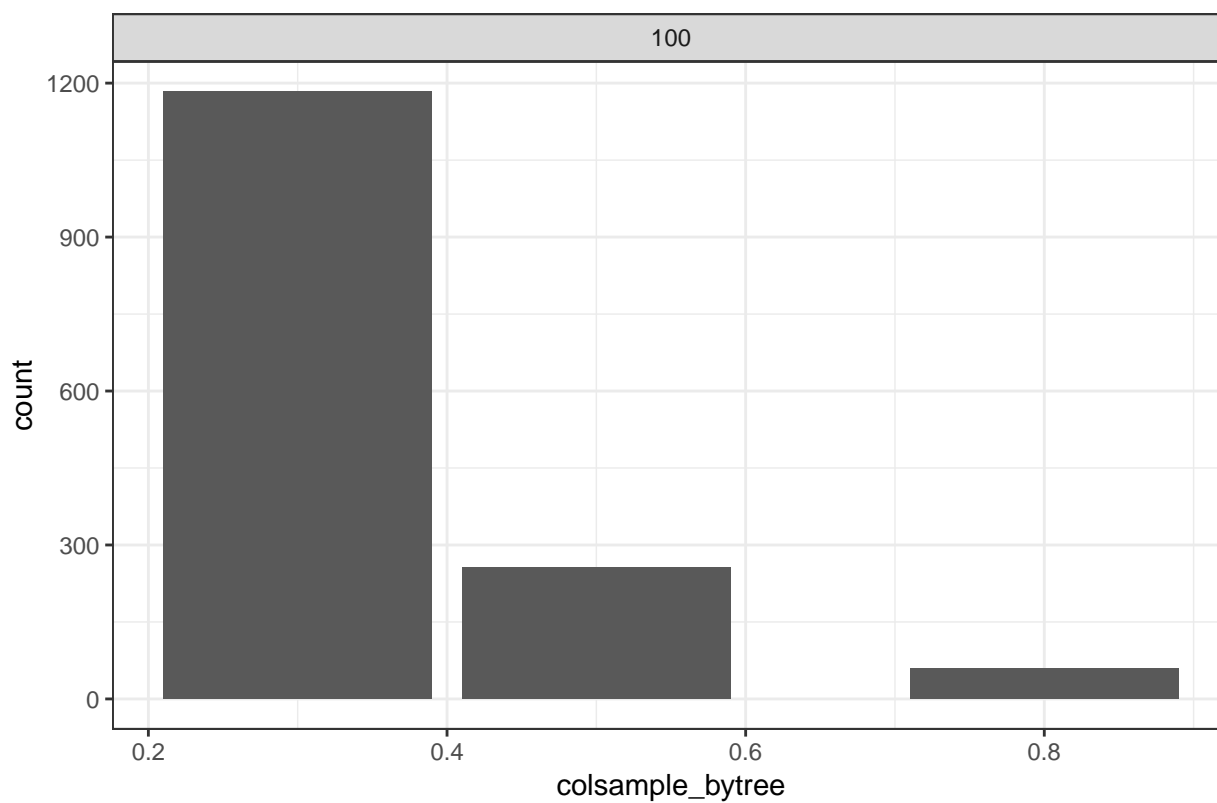Note the 2nd /3rd best parameter combinations might not be too bad either.

## optimal max_depth across seed and cv



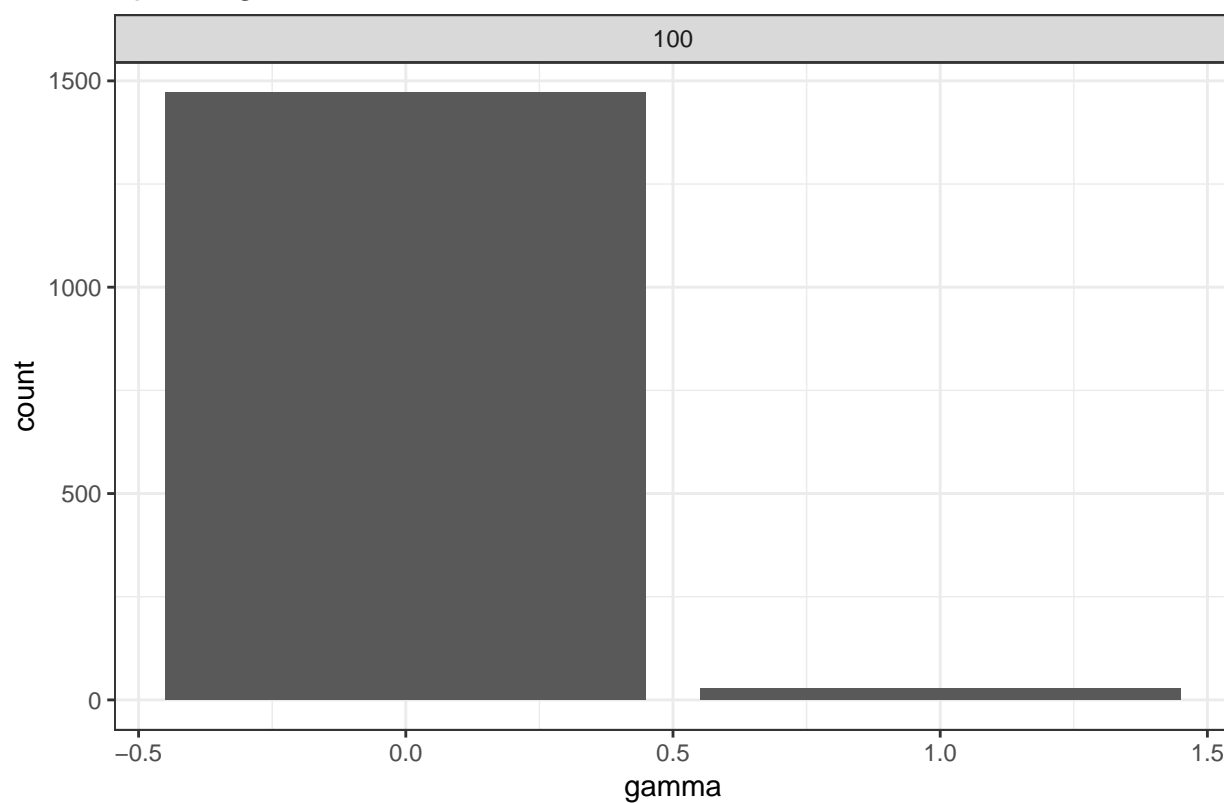## optimal max_delta_step across seed and cv

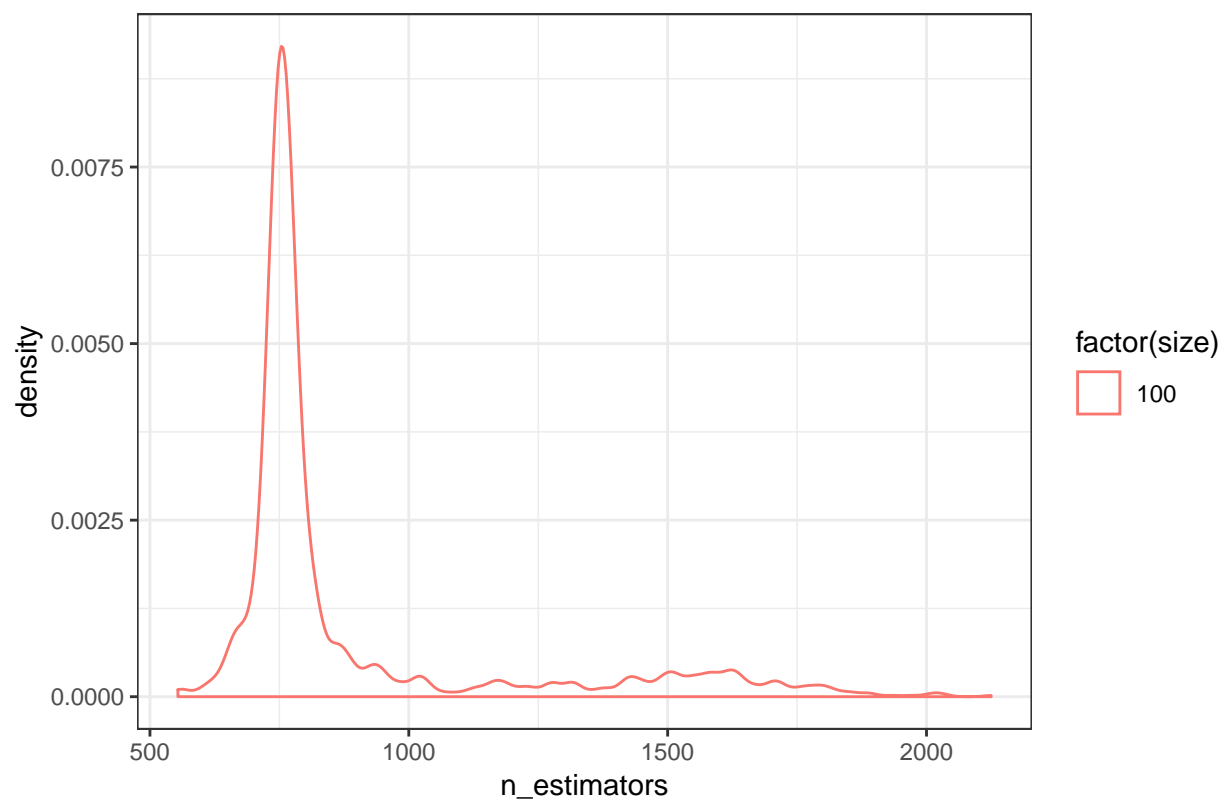## optimal alpha across seed and cv



## optimal colsample_bytree across seed and cv

optimal gamma across seed and cv

optimal n_estimator within seed and cv

**more about the best parameter combination selection**

```
select_ft_step <- 100

df1 <- subset(grid_best, size==select_ft_step & max_depth==1 &  max_delta_step == 0 )
print( paste('summary of n estimator at',select_ft_step, 'feature step'))
```

```
## [1] "summary of n estimator at 100 feature step"
```

```
print(summary(df1$n_estimators))
```

```
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##    1025    1353    1538    1510    1630    2125
```

```
df2 <- subset(df.grid, size==select_ft_step & max_depth==1 &  max_delta_step == 0 )

with(df2, plot(x = n_estimators, y=score, ylab=score_label))
```