

Evaluate testing data (multi-class) - Lasso

EVE W.

2020-04-19

Contents

0. Load Data	1
1. Scores	2
1.1 Scores per Class	2
1.2 Average score	3
2. Important Features	3

Note: The two differences between Lasso and Tree-based methods are:

1. Lasso has its own inherent feature selection process.
2. Lasso vimp will be based on how many times the feature exist in all runs.

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

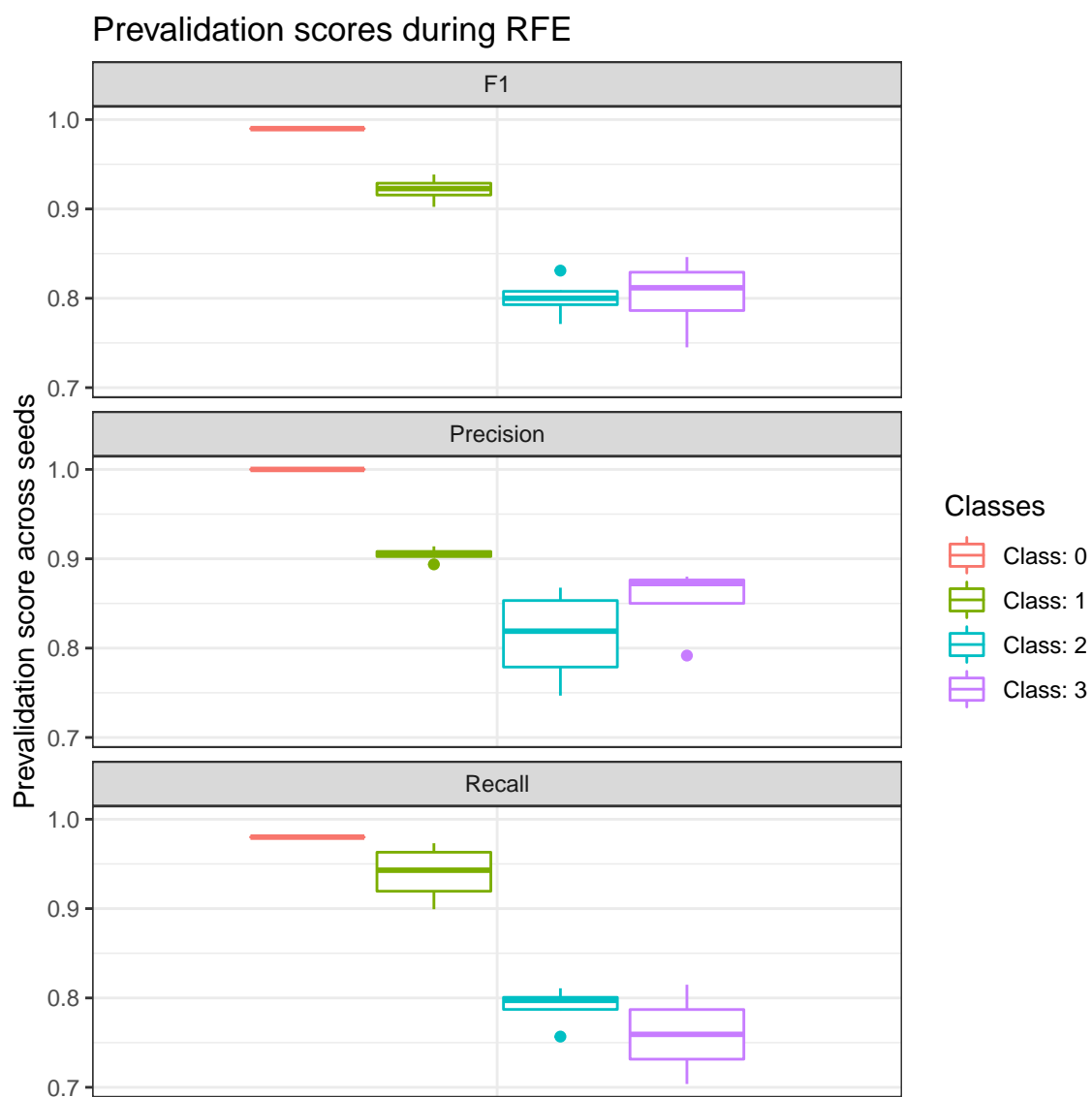
0. Load Data

```
## Error : $ operator is invalid for atomic vectors
## 300 of samples were used
## 100 of full features
## 4 runs, each run contains 3 CVs.
## Labels:
##
##    0    1    2    3
## 50 149  74  27
```

run with lasso.r with $\alpha = 0.5$.

1. Scores

1.1 Scores per Class

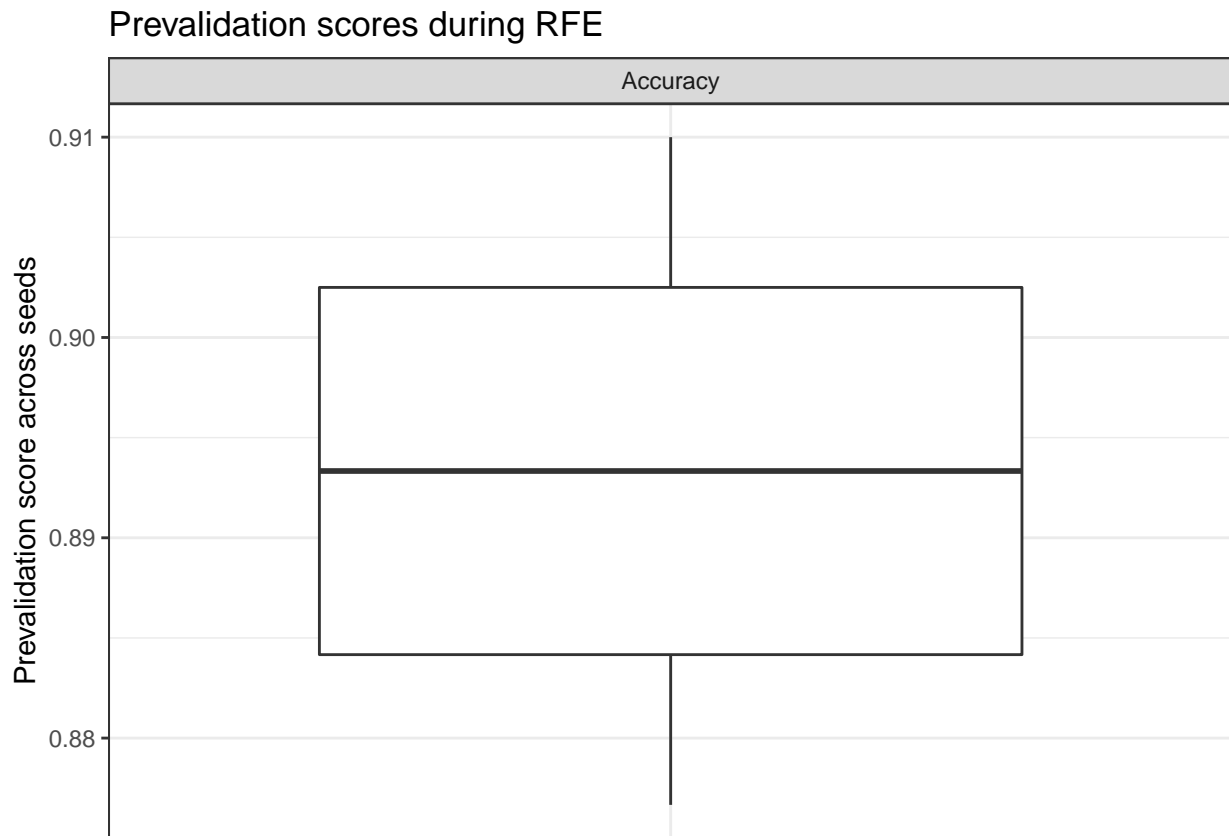


Confusion Matrix

```
## confusion matrix at feature size = 100
## sum across 4 seeds
```

```
##           Reference
## Prediction  0    1    2    3
##           0 196    0    0    0
##           1   0 560   52    7
##           2   0  36 234   19
##           3   4   0  10   82
```

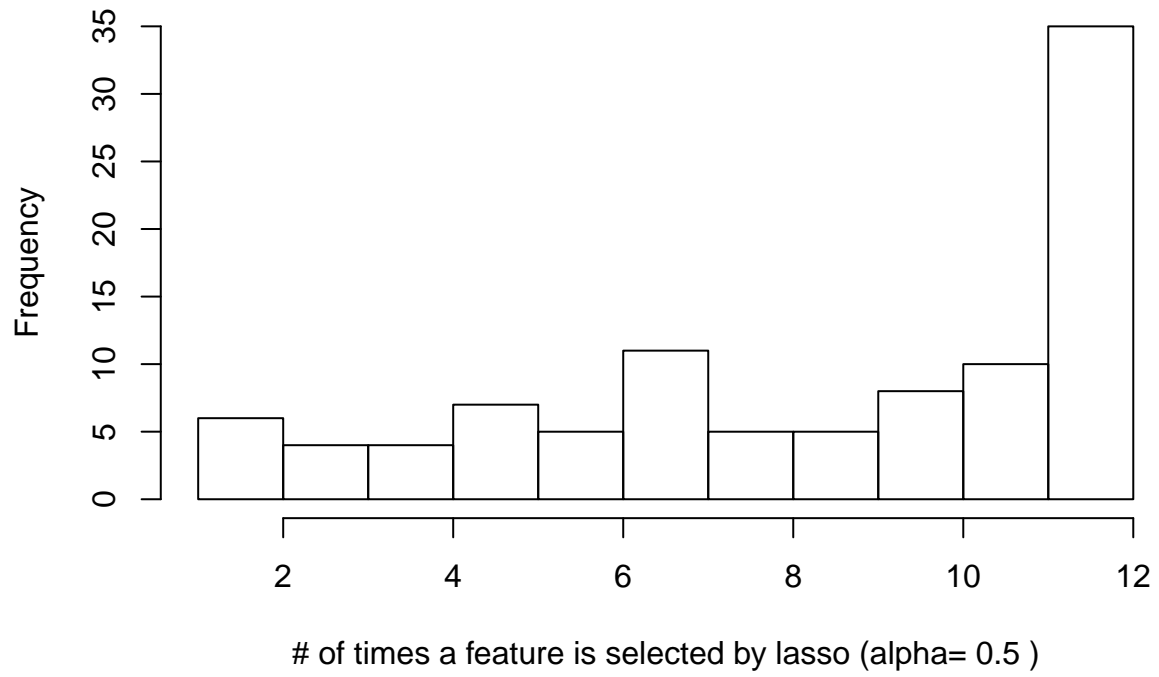
1.2 Average score



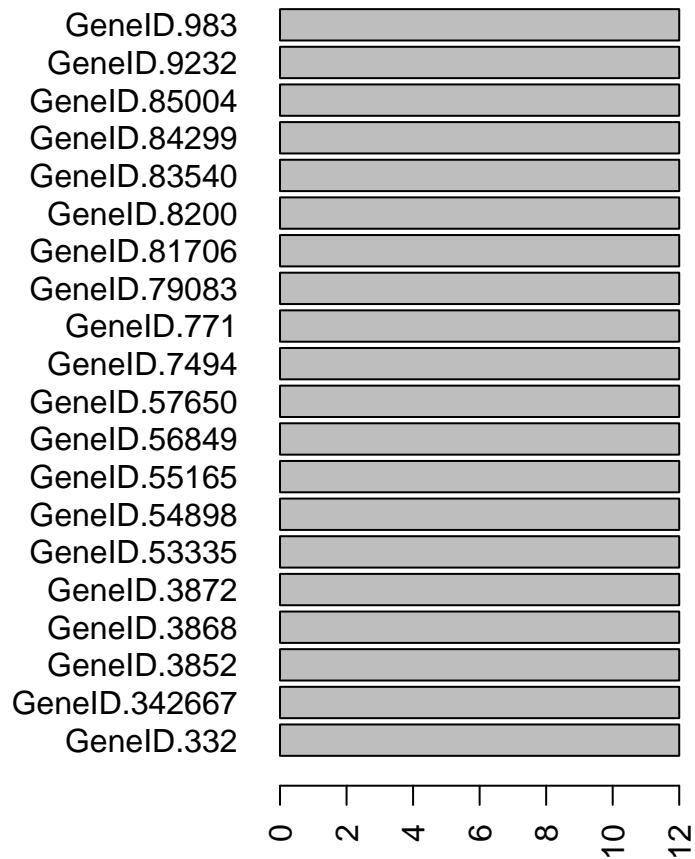
2. Important Features

For Lasso, we calculate how many times a given features is being used in all the runs.

distribution across 4 seed x 3 CV



Number of times a feature is used



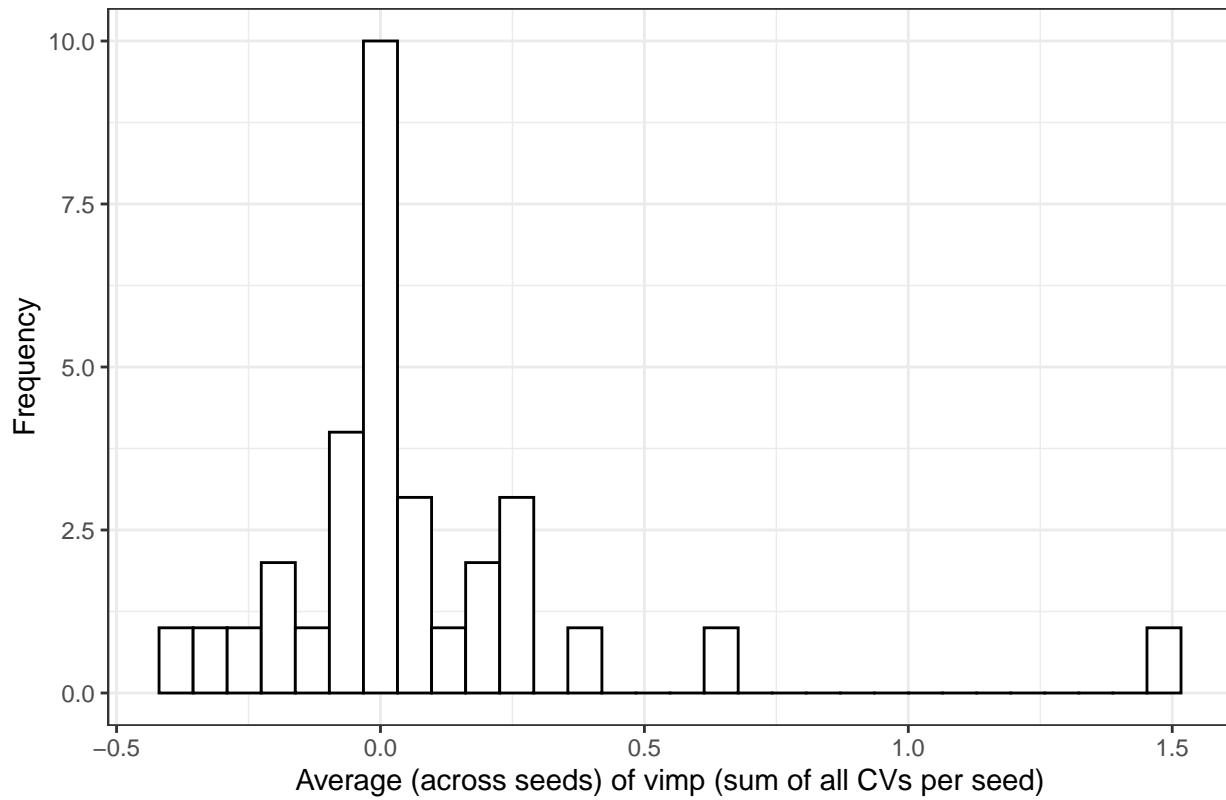
```
## (currently only Lasso has this graph)[1] "there are 100 unique features used from the 100 feature set"
## [1] "summary of number of features used in each run under 4 seeds and 3 CVs"
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  68.00  70.75   72.50   73.00   76.00   77.00
```

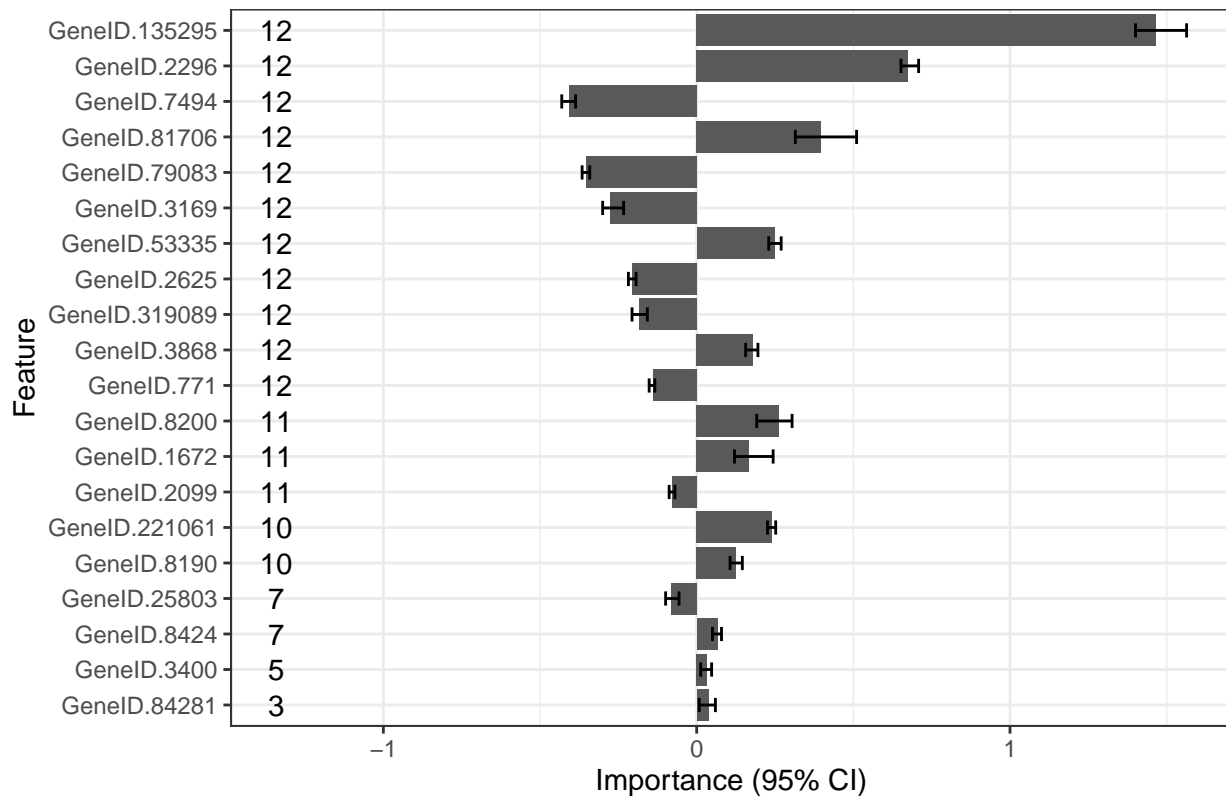
For regularized multinomial regression, glmnet does not use a reference level from the outcome variable and provides coefficients for each level of multinomial distribution. Please check out section **Regularized multinomial regression** section from here. The following barplots show the coefficients for each level.

```
## [1] "removing 644 records as their vimp is less than 1e-06"
```

Distribution of Feature Coefficient for Class 0

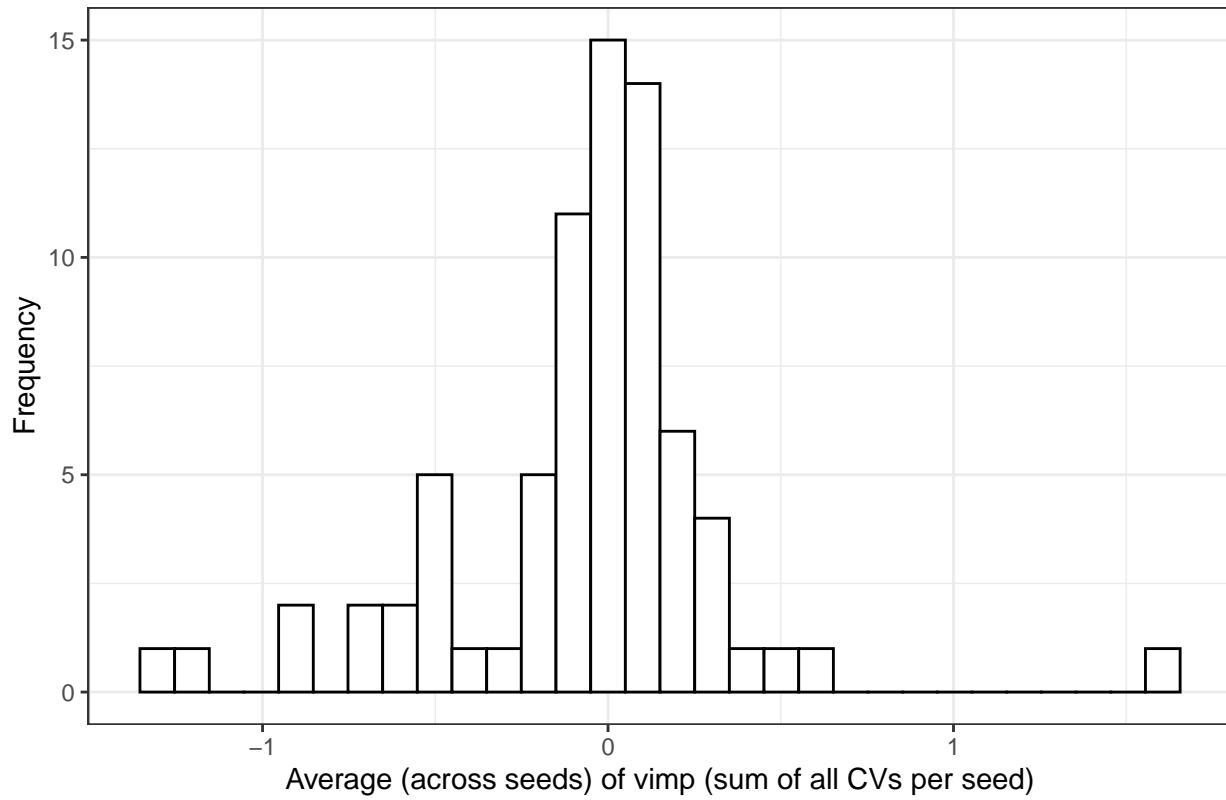


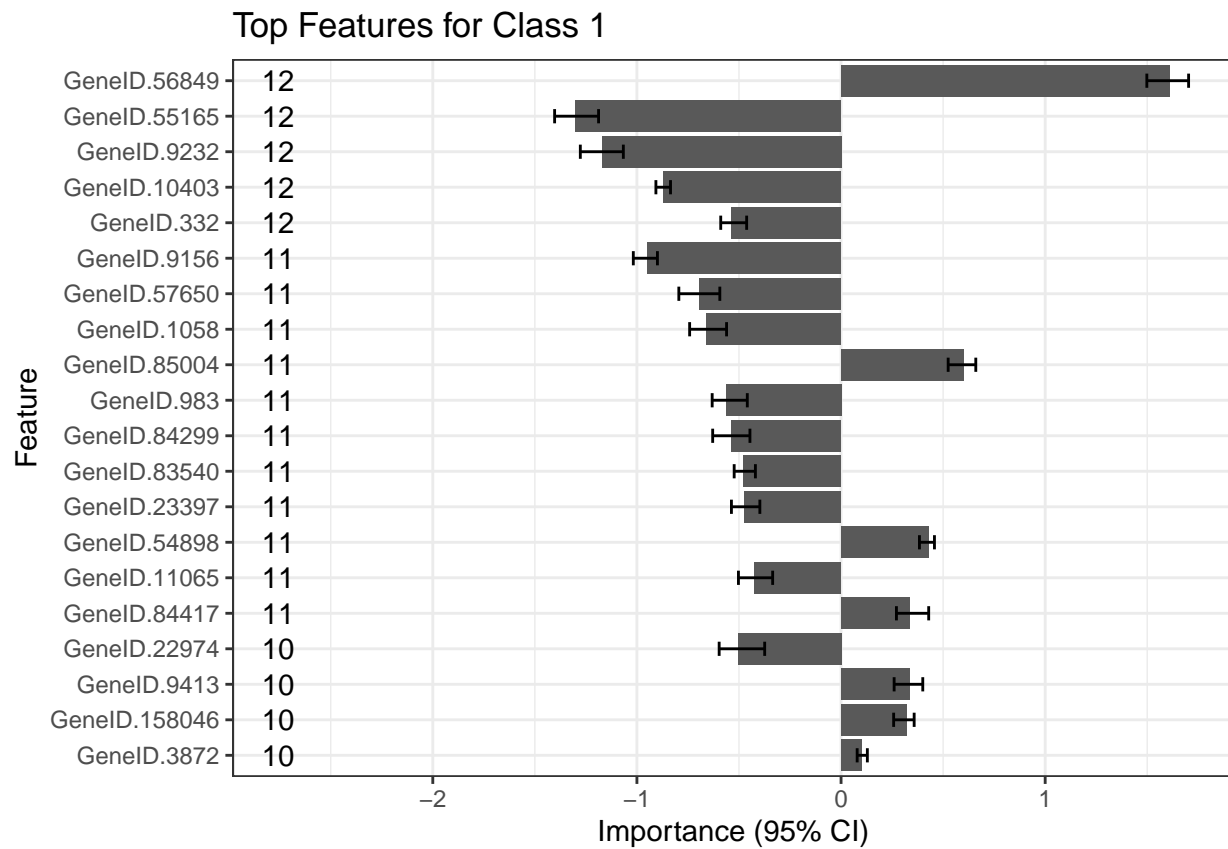
Top Features for Class 0



```
## [1] "removing 449 records as their vimp is less than 1e-06"
```

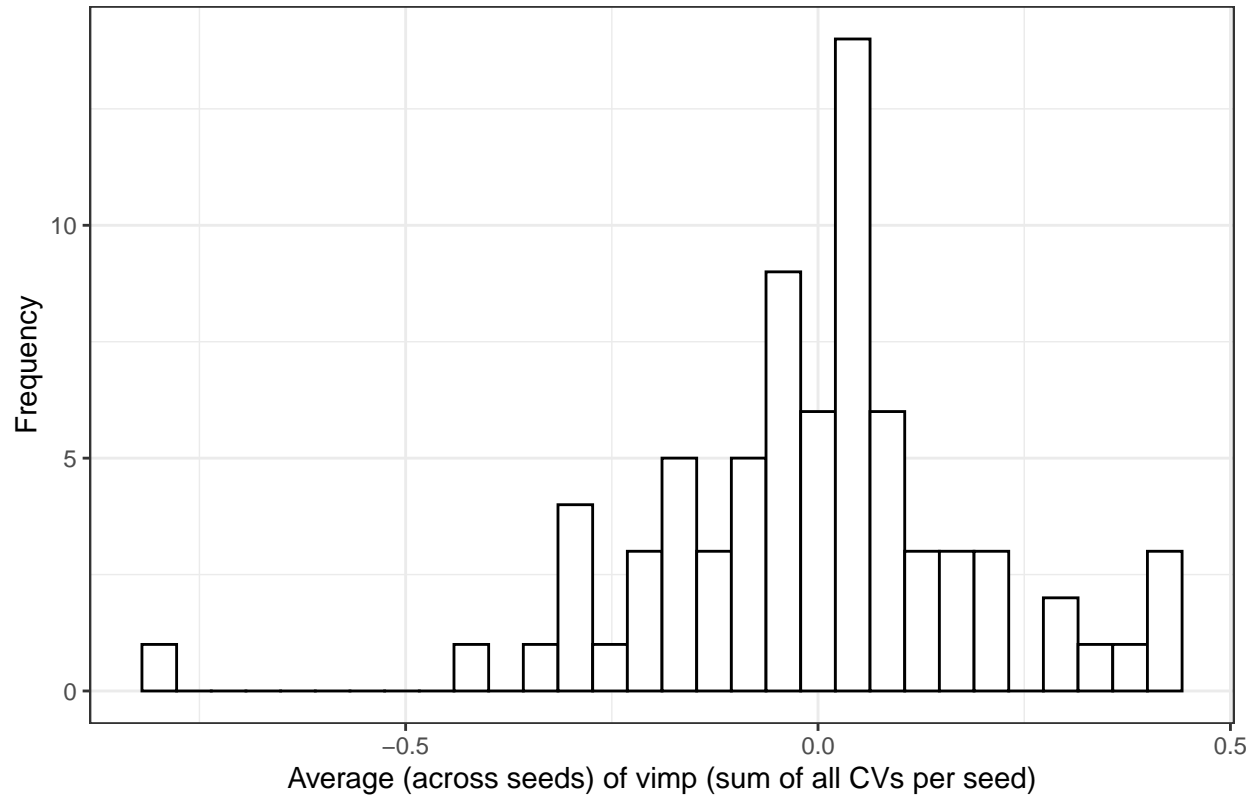
Distribution of Feature Coefficient for Class 1



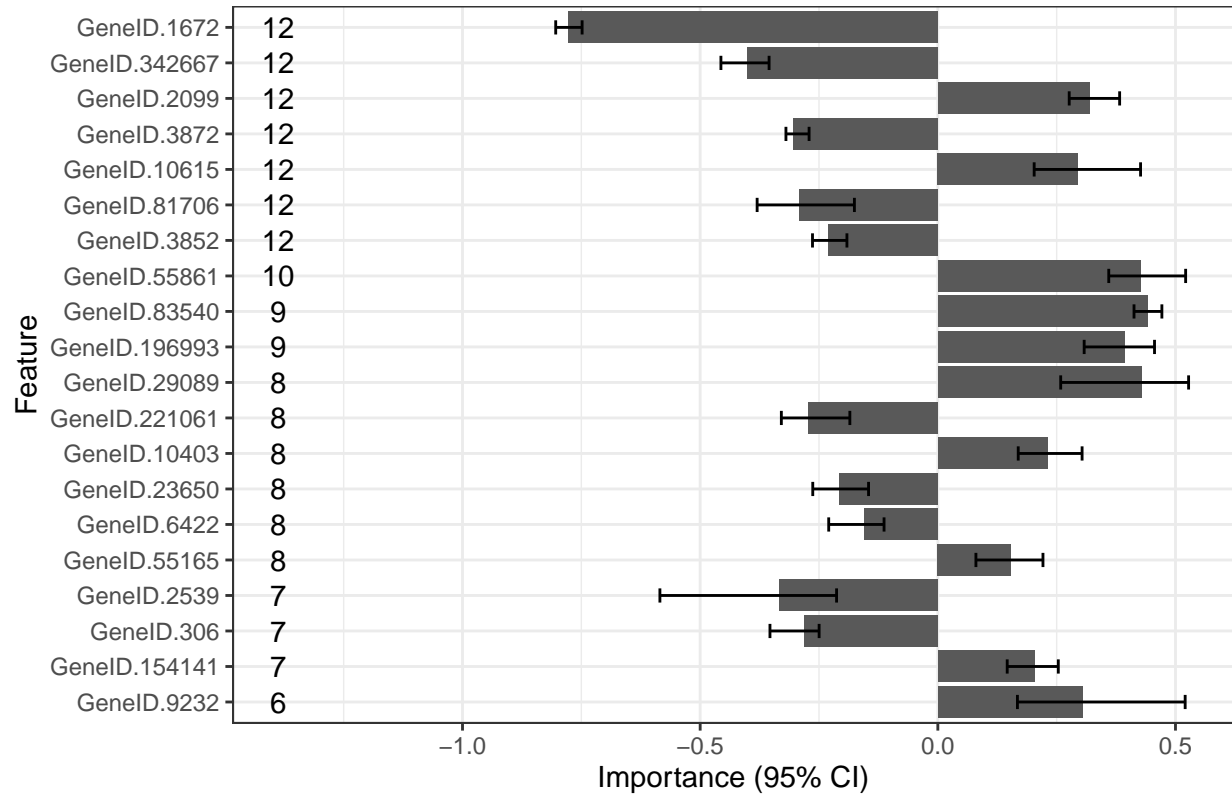


```
## [1] "removing 539 records as their vimp is less than 1e-06"
```


Distribution of Feature Coefficient for Class 2



Top Features for Class 2



```
## [1] "removing 684 records as their vimp is less than 1e-06"
```

Distribution of Feature Coefficient for Class 3

