# STA9690 Final Project

Advanced Data Mining

Christopher Lang

December 12, 2017

# The Spambase Dataset

A dataset of emails and several key features that describe the content written in the email body

| | |
|---|---|
| **Number of Emails (n)** | 2695 |
| **Number of Predictors (p)** | 57 − 3 = **54** (3 predictors are all zero) |
| **Response** | 1512 (~56%) spam<br>1183 (~44%) non-spam |

Sourced Spambase from:
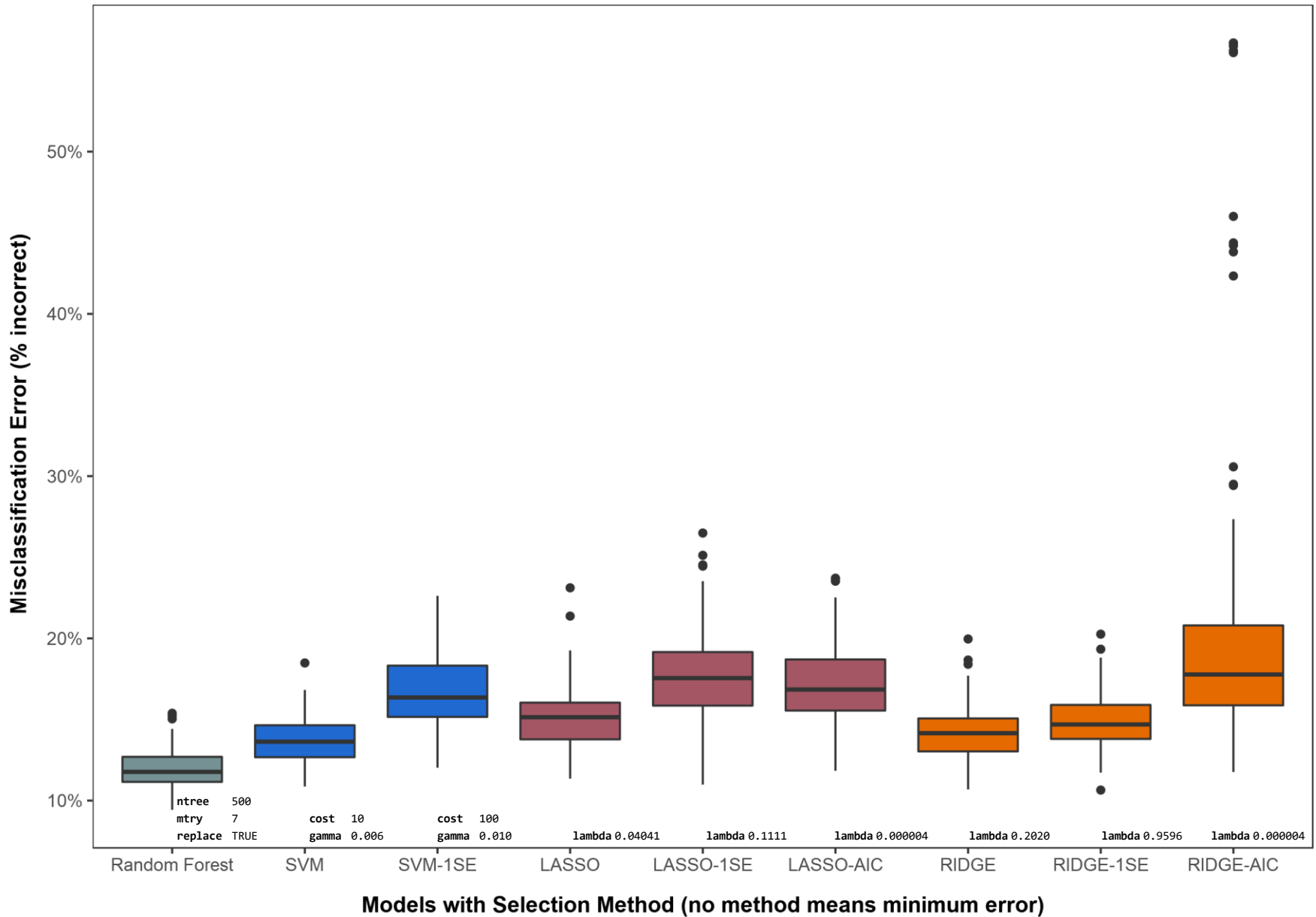
### UCI Machine Learning Repository

https://archive.ics.uci.edu/ml/datasets/spambase

## <u>Predictor Description</u>

Predictors are broken into three types of measures

| **Word Frequency**<br>% of total words in an email | | **Character Frequency**<br>% of total characters in an email | | **Capital Run Length**<br>The length of a sequence of capital letters | |
|---|---|---|---|---|---|
| | `word_freq_make` | | `char_freq_semicolon` | | `capital_run_length_longest` |
| | `word_freq_internet` | | `char_freq_exclamation` | | `capital_run_length_average` |
| | `word_freq_hp` | | `char_freq_dollarsign` | | `capital_run_length_total` |
| | …<br>… | | …<br>… | | |

45 predictors (83%)          6 predictors (11%)          3 predictors (6%)

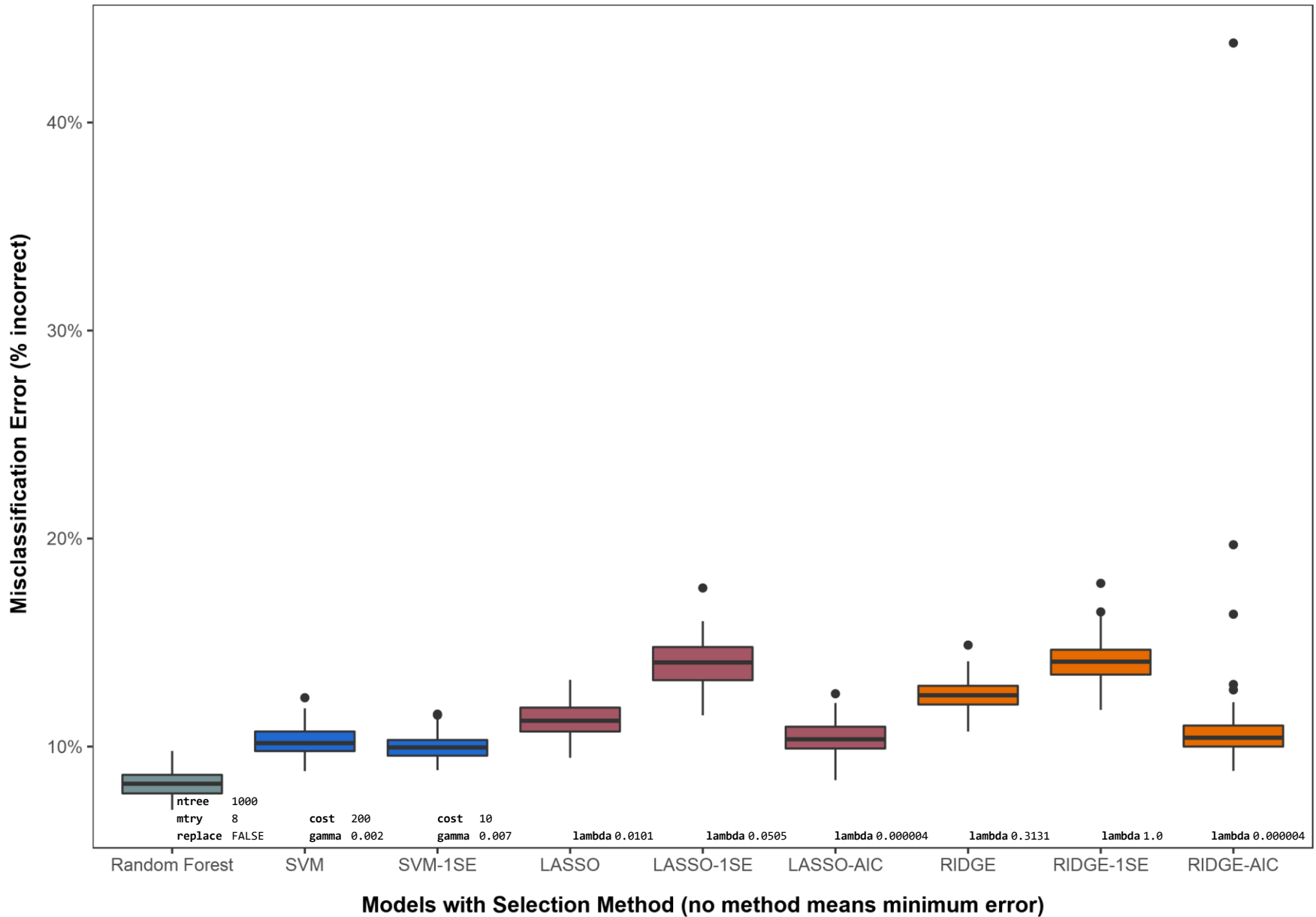**Misclassification Error for Learning Size 2p, nlearn=108**

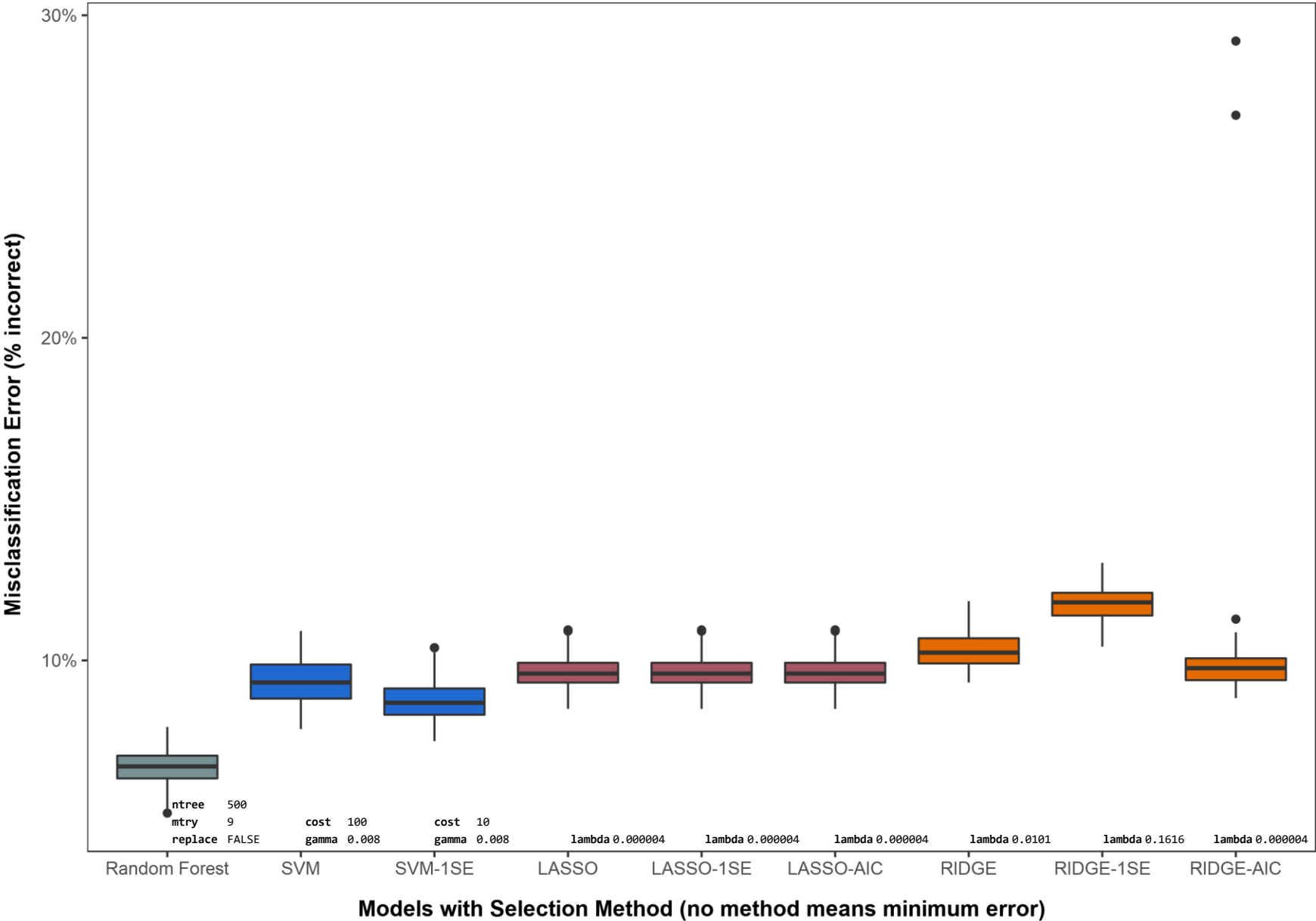Parameters tuned through a single 10-fold Cross Validation run

**Misclassification Error for Learning Size 10p, nlearn=540**

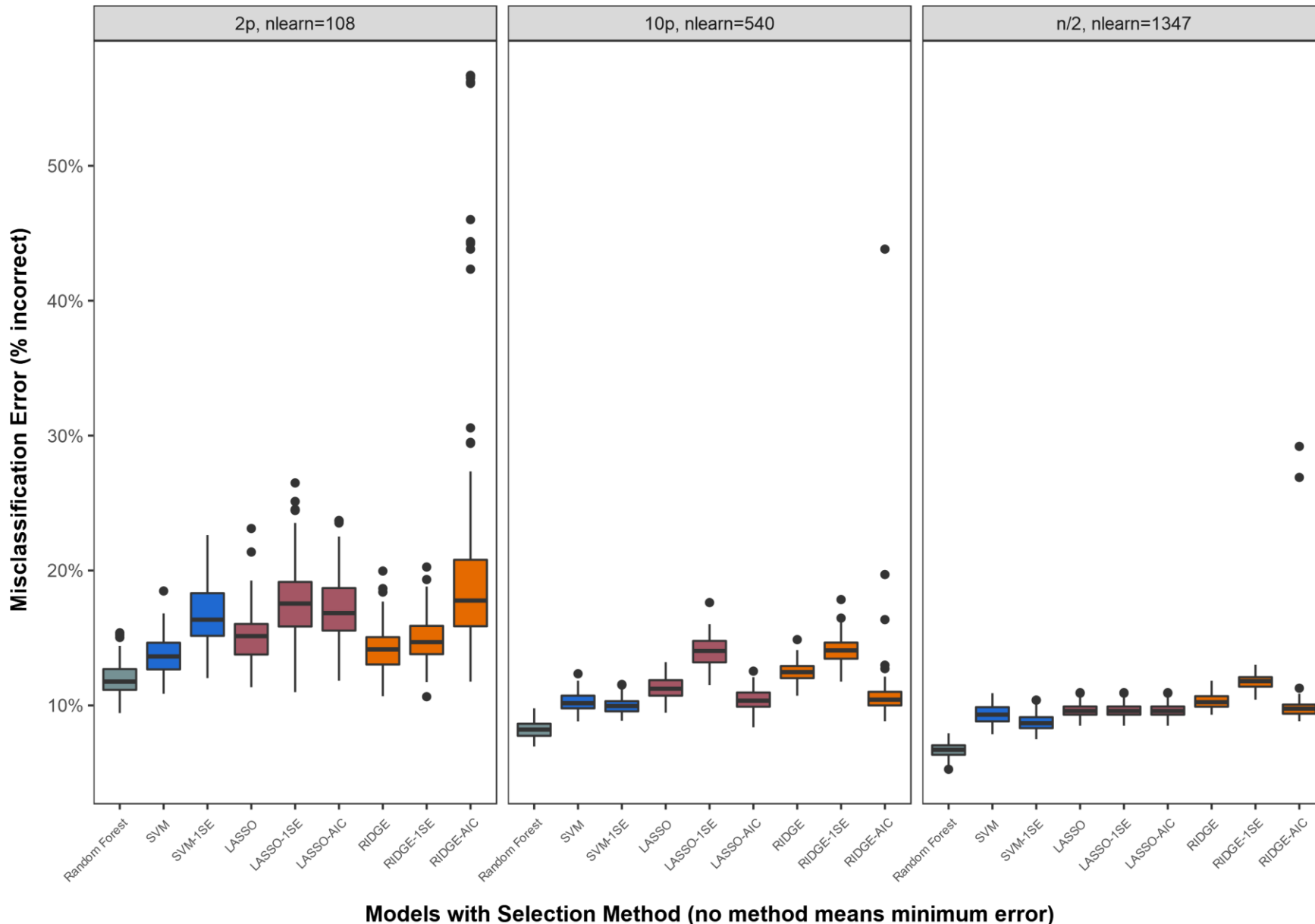Parameters tuned through a single 10-fold Cross Validation run

# Misclassification Error for Learning Size n/2, nlearn=1347

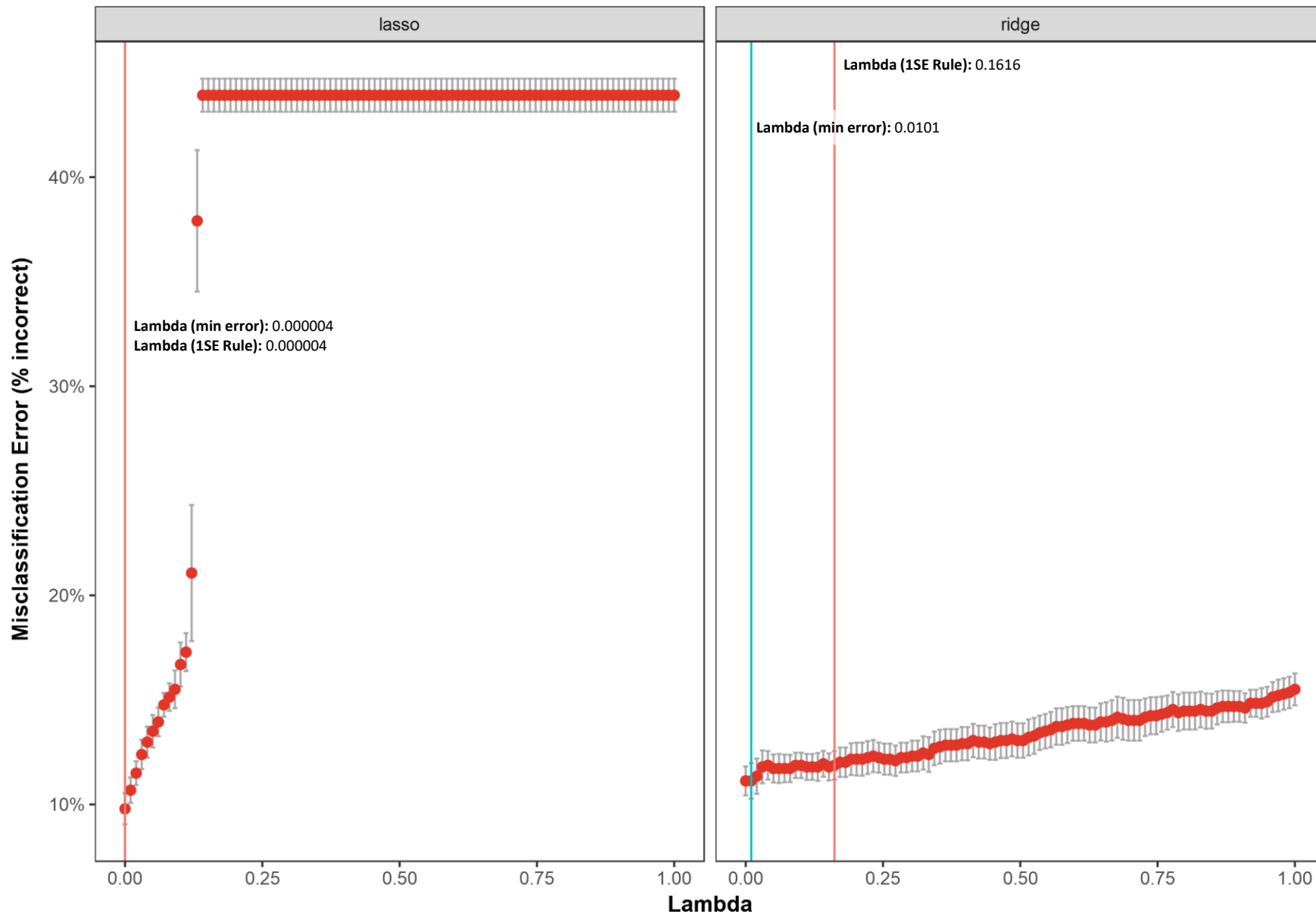Parameters tuned through a single 10-fold Cross Validation run

**Misclassification Error for All Learning Data Sizes**

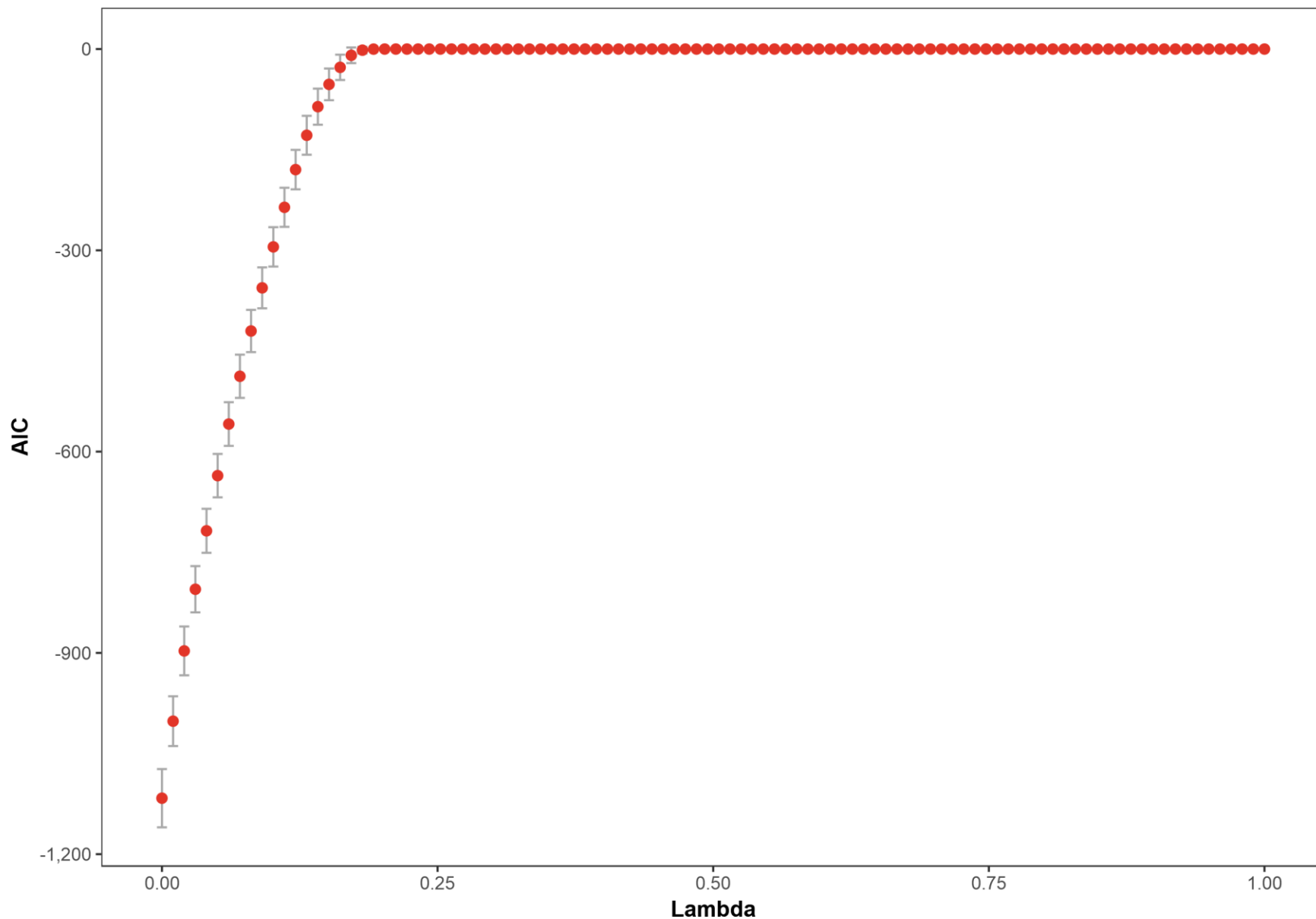Parameters tuned through a single 10-fold Cross Validation run

**Ridge Regression 10-fold Cross Validation Curve**

Curve for n/2, nlearn=1347 learning set
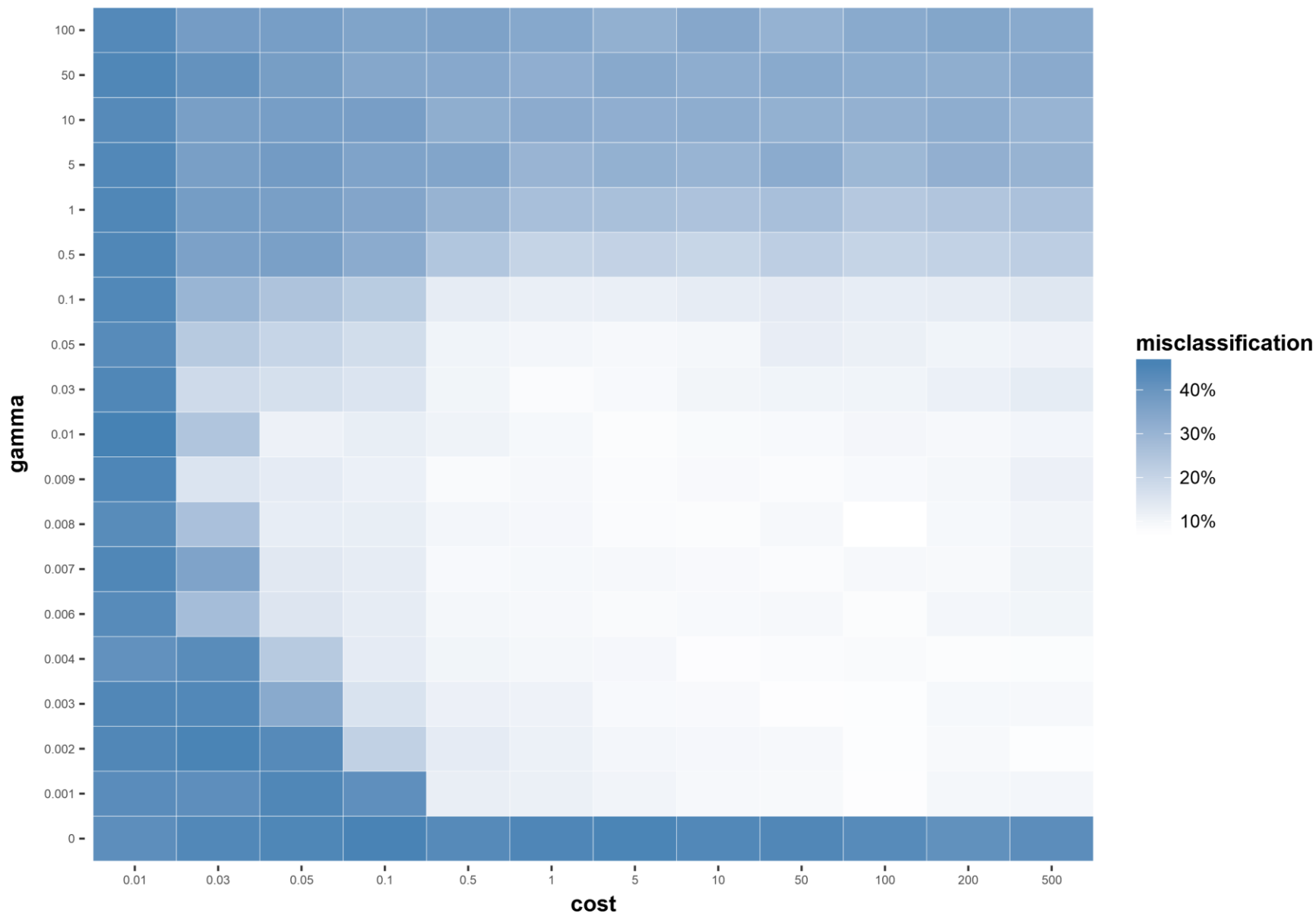
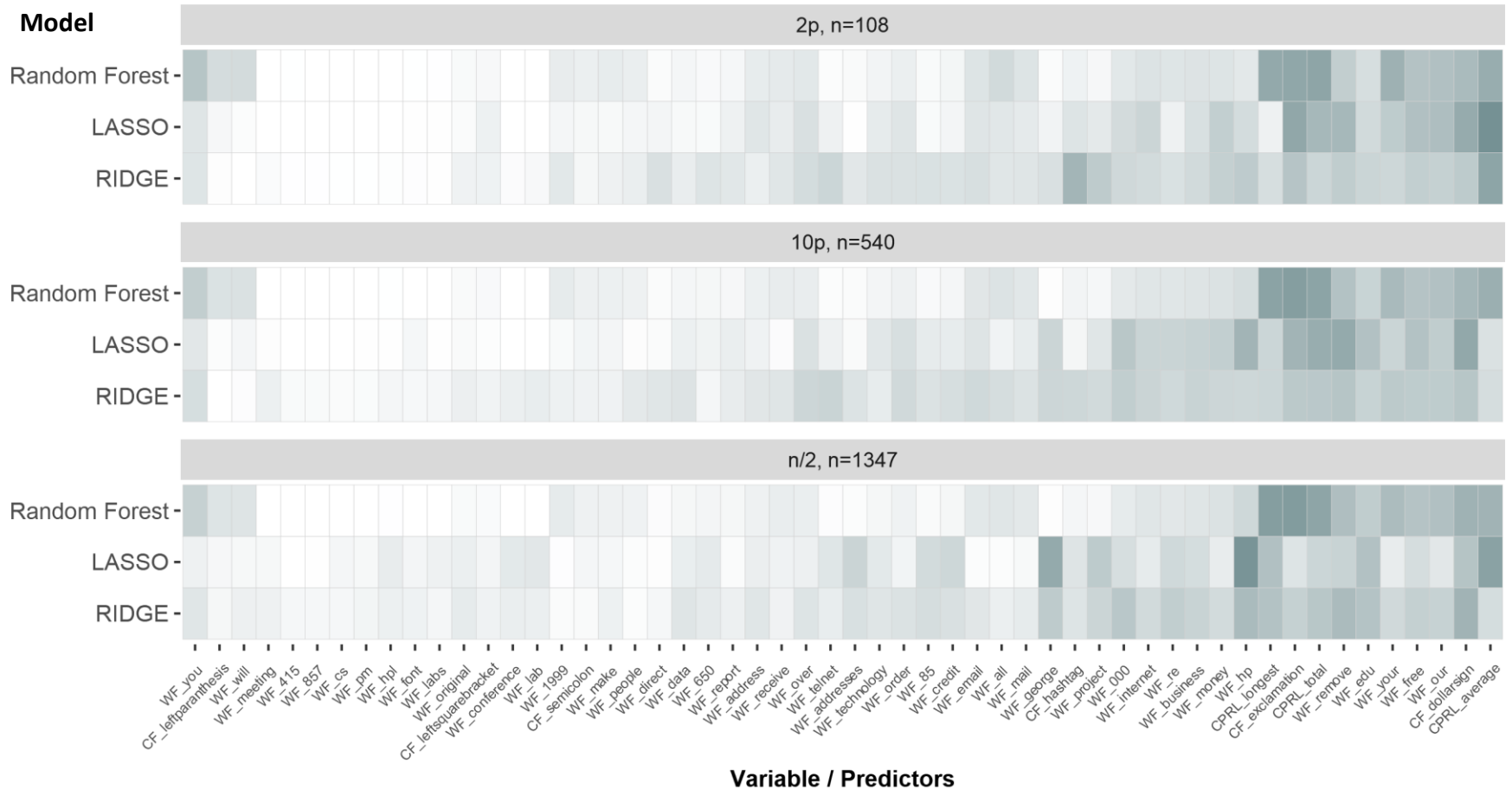**Lasso Regression AIC for 100 Repeated Sampling**

Curve for n/2, nlearn=1347 learning set

**Support Vector Machine Parameter Performance**

SVM heatmap for n/2, nlearn=1347 learning set

# Variable Importance Agreement

Averaged Coefficients and Mean Decrease Gini

**Variable Naming Scheme**

<variable type>_<word/character/measure>

**WF** – **W**ord **F**requency (% of total words in an email)
**CF** – **C**haracter **F**requency (% of total characters in an email)
**CPRL** – **Cap**ital **R**un **L**ength (The length of a sequence of capital letters)