

Statistical Modeling Course

Variable Selection Assignment

In this exercise, we will predict the number of applications, `Apps`, received using the other variables in the `College` data set. Use 5-fold cross-validation to calculate the test error using the following methods. You should fit each model five times leaving out one fifth of the data each time and calculating the test error (RMSE) on the data you left out. The total test error will be the average RMSE from the five folds. You should also calculate the standard error of the test errors, `sd(test_error)/sqrt(5)`. For best subset selection, and forward stepwise selection chose the model using BIC. For lasso, ridge, and pcr you will conduct a second level of cross validation for each of the 5 training sets. You can do this with `cv.glmnet`.

Compare your results with the null model. Code for the null model and splits is given below, make sure to compare all of your models on the same splits of the data.

```
set.seed(2020)
# Vector of data split
folds <- sample(1:5, nrow(College), replace = TRUE)

# Variable to store error from each fold
error_null <- c()

# Loop through folds and fit models with no predictors
for (i in 1:5){
  test <- College[folds == i, ]
  train <- College[folds != i, ]
  error_null[i] <- sqrt(mean((test$Apps - mean(train$Apps))^2))
}

# Print mean RMSE
mean(error_null)

## [1] 3812.055

# Print standard error of RMSE
sd(error_null)/sqrt(5)

## [1] 363.8367
```

Problem 1

- Least squares linear model
- Best subset selection (chosen using BIC)
- Forward stepwise subset selection (chosen using BIC)

Problem 2

- Ridge-regression with lambda chosen by cross-validation, report the five λ 's chosen
- Lasso with lambda chosen by cross-validation, report the five λ 's chosen

Problem 3

- Fit a PCR model, with M chosen by cross-validation. Report the test error (MSE) obtained, along with the value of M selected by cross-validation.

Problem 4

Comment on the results obtained. How accurately can we predict the number of college applications received? Is there much difference among the test errors resulting from the approaches?

Problem 5

Generate a data set with $p = 100$ features, $n = 300$ observations, and an associated quantitative response vector generated according to the model

$$p_i = \frac{e^{x_i^T \beta}}{1 + e^{x_i^T \beta}}$$
$$y_i \sim \text{Bin}(p_i, n = 1)$$

Simulate some of the features as categorical and some as numeric. Set most of the values of $\beta_p = 0$ for most but not all p .

- Using your simulated dataset split your dataset into a training set and a test set containing using an 80/20 split.
- Perform lasso and ridge regression on the training set.
- Which model has a lower test error (MSE)? Comment on your results.