

SDSC 3006 Fundamentals of Machine Learning I
Assignment 3

Deadline: November 16, Wednesday @ 10:00 PM

1. We perform best subset, forward stepwise, and backward stepwise selection on a single data set. For each approach, we obtain $p + 1$ models, containing $0, 1, 2, \dots, p$ predictors. Answer the following questions:

- (a) Which of the three models with k predictors has the smallest **training RSS**?
- (b) Which of the three models with k predictors has the smallest **test MSE**?
- (c) True or False for each statement below.
 - i. The predictors in the k - variable model identified by forward stepwise are a subset of the predictors in the $(k + 1)$ -variable model identified by forward stepwise selection.
 - ii. The predictors in the k - variable model identified by backward stepwise are a subset of the predictors in the $(k + 1)$ -variable model identified by backward stepwise selection.
 - iii. The predictors in the k - variable model identified by backward stepwise are a subset of the predictors in the $(k + 1)$ -variable model identified by forward stepwise selection.
 - iv. The predictors in the k - variable model identified by forward stepwise are a subset of the predictors in the $(k + 1)$ -variable model identified by backward stepwise selection.
 - v. The predictors in the k - variable model identified by best subset are a subset of the predictors in the $(k + 1)$ -variable model identified by best subset selection.

2. Choose the correct answer for each question below.

- (a) The **lasso**, relative to least squares, is:
 - i. More flexible and hence will give improved prediction accuracy when its increase in bias is less than its decrease in variance.
 - ii. More flexible and hence will give improved prediction accuracy when its increase in variance is less than its decrease in bias.
 - iii. Less flexible and hence will give improved prediction accuracy when its increase in bias is less than its decrease in variance.
 - iv. Less flexible and hence will give improved prediction accuracy when its increase in variance is less than its decrease in bias.
- (b) **Ridge regression**, relative to least squares, is:
 - i. More flexible and hence will give improved prediction accuracy when its increase in bias is less than its decrease in variance.
 - ii. More flexible and hence will give improved prediction accuracy when its increase in variance is less than its decrease in bias.
 - iii. Less flexible and hence will give improved prediction accuracy when its increase in bias is less than its decrease in variance.

iv. Less flexible and hence will give improved prediction accuracy when its increase in variance is less than its decrease in bias.

3. Suppose we estimate the regression coefficients in a linear regression model by minimizing

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \quad \text{subject to} \quad \sum_{j=1}^p |\beta_j| \leq s$$

for a particular value of s . Choose the correct answer for each question below.

(a) As we increase s from 0, the **training RSS** will:

- Increase initially, and then eventually start decreasing in an inverted U shape.
- Decrease initially, and then eventually start increasing in a U shape.
- Steadily increase.
- Steadily decrease.
- Remain constant.

(b) As we increase s from 0, the **test MSE** will:

- Increase initially, and then eventually start decreasing in an inverted U shape.
- Decrease initially, and then eventually start increasing in a U shape.
- Steadily increase.
- Steadily decrease.
- Remain constant.

(c) As we increase s from 0, the **variance** will:

- Increase initially, and then eventually start decreasing in an inverted U shape.
- Decrease initially, and then eventually start increasing in a U shape.
- Steadily increase.
- Steadily decrease.
- Remain constant.

(d) As we increase s from 0, the **squared (bias)** will:

- Increase initially, and then eventually start decreasing in an inverted U shape.
- Decrease initially, and then eventually start increasing in a U shape.
- Steadily increase.
- Steadily decrease.
- Remain constant.

4. Suppose we estimate the regression coefficients in a linear regression model by minimizing

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

for a particular value of λ . Choose the correct answer for each question below.

(a) As we increase λ from 0, the **training RSS** will:

- Increase initially, and then eventually start decreasing in an inverted U shape.
- Decrease initially, and then eventually start increasing in a U shape.

- iii. Steadily increase.
- iv. Steadily decrease.
- v. Remain constant.

(b) As we increase λ from 0, the **test MSE** will:

- i. Increase initially, and then eventually start decreasing in an inverted U shape.
- ii. Decrease initially, and then eventually start increasing in a U shape.
- iii. Steadily increase.
- iv. Steadily decrease.
- v. Remain constant.

(c) As we increase λ from 0, the **variance** will:

- i. Increase initially, and then eventually start decreasing in an inverted U shape.
- ii. Decrease initially, and then eventually start increasing in a U shape.
- iii. Steadily increase.
- iv. Steadily decrease.
- v. Remain constant.

(d) As we increase λ from 0, the **(squared) bias** will:

- i. Increase initially, and then eventually start decreasing in an inverted U shape.
- ii. Decrease initially, and then eventually start increasing in a U shape.
- iii. Steadily increase.
- iv. Steadily decrease.
- v. Remain constant.

5. We will predict the number of applications received in the **College** data set.

- (a) Split the data set into a training set and a test set.
- (b) Fit a linear model using least squares on the training set, and report the test error obtained.
- (c) Fit a ridge regression model on the training set, with λ chosen by cross validation. Report the test error obtained.
- (d) Fit a lasso model on the training set, with λ chosen by cross validation. Report the test error obtained, along with the number of non-zero coefficient estimates.
- (e) Fit a PCR model on the training set, with M (the number of principal components) chosen by cross validation. Report the test error obtained, along with the number of PCs selected by cross validation.

6. Suppose we produce 10 bootstrapped samples from a data set containing red and green classes. We then apply a classification tree to each bootstrapped sample and, for a specific value of X , produce 10 estimates of $P(\text{Class is Red} \mid X)$:

0.1, 0.15, 0.2, 0.2, 0.55, 0.6, 0.6, 0.65, 0.7, 0.75

There are two common ways to combine these results together into a single class prediction. One is the majority vote approach, and the other one is to classify based on the average probability.

- (a) What is the final classification under the majority vote approach?
- (b) What is the final classification under the average probability approach?

7. In the lab of tree models, a classification tree was applied to the **Carseats** data set after converting **Sales** into a qualitative response variable. Now we will seek to predict **Sales** using regression trees and related approaches, treating the response as a quantitative variable.

- (a) Split the data set into a training set and a test set.
- (b) Fit a regression tree to the training set. Plot the tree, and interpret the results. What test error rate do you obtain?
- (c) Perform tree pruning and use cross validation to determine the optimal level of tree complexity. Does pruning the tree improve the test error rate?
- (d) Use the bagging approach to analyze this dataset. What test error rate do you obtain? Find which variables are most important.
- (e) Use random forests to analyze this dataset. What test error rate do you obtain? Find which variables are most important. Describe the effect of m , the number of variables considered at each split, on the error rate obtained.

8. This problem involves the **OJ** data set which is part of the **ISLR2** package.

- (a) Create a training set containing a random sample of 800 observations, and a test set containing the remaining observations.
- (b) Fit a support vector classifier to the training data using **cost = 0.01**, with **Purchase** as the response and the other variables as predictors. Use the **summary()** function to produce summary statistics, and describe the results obtained.
- (c) What are the training and test error rates?
- (d) Use the **tune()** function to select an optimal **cost**. Consider values in the range 0.01 to 10.
- (e) Compute the training and test error rates using this new value for **cost**.
- (f) Repeat parts (b) through (e) using a support vector machine with a radial kernel. Use the default value for **gamma**.
- (g) Repeat parts (b) through (e) using a support vector machine with a polynomial kernel. Set **degree = 2**.
- (h) Overall, which approach seems to give the best results on this data?