

STAT S4240 002, Homework 4

Brian Weinstein (bmw2148)

August 13, 2015

Problem 1: James 6.8.1

- (a) The model obtained through *best subset selection* has the smallest training RSS for $k \leq p$ predictors. This is because best subset selection tests all possible subsets of size k and chooses the one that minimizes the training RSS.

There is a chance that by using *forward stepwise selection* or *backward stepwise selection* we'd end up with the same optimal subset of k predictors (i.e., the one we obtain via best subset selection), but since the k -sized subset is dependent on the covariates included in the prior steps in each of these methods, there's no guarantee that this will be the case.

- (b) Without either valuating each model on a testing set or estimating test error, none among the three models is guaranteed to have the lowest test RSS. These methods only choose the subset that results in the lowest training RSS, which in no way indicates if the model will also result in a low test RSS.
- (c)
- True. In a model identified by forward stepwise selection, the subset of predictors in a given step (with k predictors) is the base upon which the next predictor is added in the subsequent step (with $k + 1$ predictors).
 - True. In a model identified by backward stepwise selection, the subset of predictors in a given step (with k predictors) is itself a subset of the predictors used in the previous step (with $k + 1$ predictors).
 - False. The forward stepwise and backward stepwise models using k predictors are not guaranteed to use the same subset of predictors. As such, once an additional variable is added to the forward stepwise model (or equivalently, once a variable is removed from the backward stepwise model), there's no guarantee that the variables in the k -variable backward stepwise model are a subset of the $k + 1$ -variable forward stepwise model.
 - False, with similar reasoning to Part (iii). The k predictors in a forward stepwise model are the predictors that, at each incremental step, kept the error as low as possible. Similarly, the $k + 1$ predictors in a backward stepwise model are the predictors at each incremental step that kept the error as low as possible. Since each subset of predictors depends on those used in the previous step, there's no guarantee (although it may happen to be the case) that a set of predictors obtained via forward selection is a subset of those obtained via backward selection.
 - False. The set of k predictors chosen in best subset selection are those that keep the error in the k -variable model as low as possible. Independently, the set of $k + 1$ predictors chosen in best subset are those that keep the error in the $k + 1$ -variable model as low as possible. There's no guarantee that k of the variables included in the $(k + 1)$ -variable model are those that minimize the error in the k -variable model.

Problem 2: James 6.8.3

Estimating the regression coefficients in a linear regression model via

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

for a given value of s .

As we increase s from 0:

- The training RSS will (iv.) steadily decrease. As we make the constraint on $\sum_{j=1}^p |\beta_j|$ less strict, the limitation on the size of $\hat{\beta}$ disappears and we eventually get the standard “least residuals” solution. At $s = 0$ we start with only the intercept β_0 non-zero. As we increase s we allow for additional non-zero coefficients, eventually including all coefficients in the model. With all coefficients included, we minimize the training error (while potentially overfitting).
- The testing RSS will (ii.) decrease initially, and then eventually start increasing in a U shape. At $s = 0$ we start with a highly-biased model with only the intercept β_0 non-zero. As we increase s we allow for *some* of the $\beta_{j \neq 0}$ to take non-zero values (since the constraint is the same as that used in the lasso), which reduces the bias and increases the variance, for an overall drop in RSS. Once we increase s enough, *all* of the $\beta_{j \neq 0}$ will take non-zero values and we’re left with a model with very low bias, but very high variance, increasing the RSS.
- The variance will (iii.) steadily increase. The reasoning here is similar to that used in Part (ii): At $s = 0$ we start with a very inflexible model, and as we increase s to allow more non-zero coefficients, the model flexibility, and thus our variance, increases.
- The (squared) bias will (iv.) steadily decrease. The reasoning here is similar to that used in Part (ii): At $s = 0$ we start with a highly-biased, very inflexible model with only one coefficient, β_0 . As we increase s to allow more non-zero coefficients, the model flexibility, and thus our squared bias, decreases.
- The irreducible error will (v.) remain constant. The irreducible error is a function of the data we’re given, not the model we choose, so our choice of s has no effect on the size of the irreducible error.

Problem 3: James 6.8.5

(a)

Problem
3

Problem 4: James 8.4.5

10 bootstrapped samples produces estimates of $P(\text{Class is Red}|X)$ that yield 0.1, 0.15, 0.2, 0.2, 0.55, 0.6, 0.6, 0.65, 0.7, and 0.75.

Based on the majority vote approach, we’d classify this X as Red, since 6 of the 10 estimates have probabilities greater than 0.5.

Based on the average probability approach, we’d classify this X as Green, since the average probability from our 10 samples is 0.45 (< 0.5).

Problem 5: Federalist Papers Classification with Trees

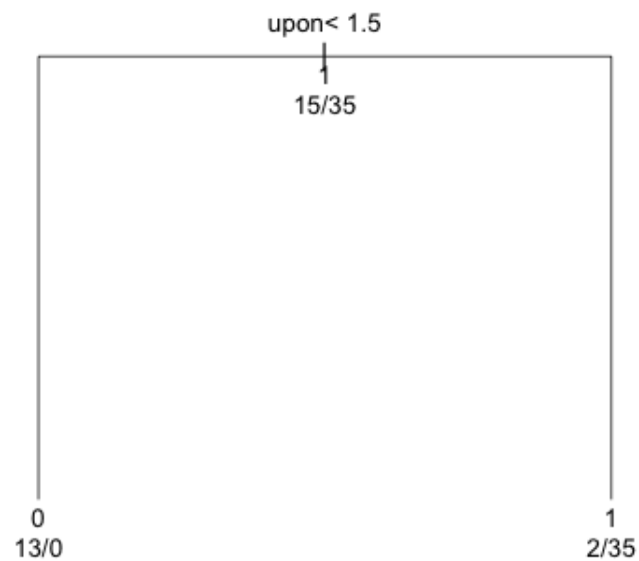
(a) Gini Impurity Coefficient Splits

i. Model Assessment

	accuracy	falseNegRate	falsePosRate
1	0.962963	0	0.09090909

- Accuracy: % of test papers that are classified correctly
- False Negative Rate: % of Hamilton papers incorrectly classified as Madison
- False Positive Rate: % of Madison papers incorrectly classified as Hamilton

ii. Tree Plot

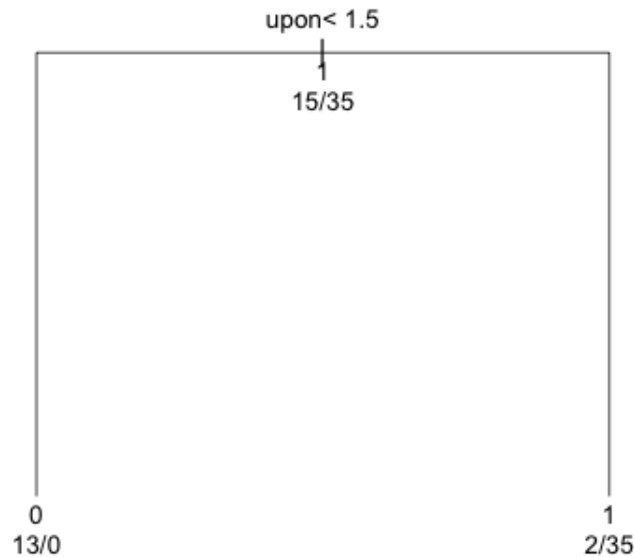


(b) Information Gain Splits

i. Model Assessment

	accuracy	falseNegRate	falsePosRate
1	0.962963	0	0.09090909

ii. Tree Plot



- iii. The trees generated from Gini Impurity Coefficient splits and Information Gain splits ended up being identical here.

Problem 6: Federalist Papers Classification with Regularized Logistic Regression

- (a) We should always center and scale our data whenever we use a regularization method. Without scaling our data, the variables with larger values would have smaller coefficients (to indicate some level of importance), which would decrease the penalty for these predictors. Similarly, variables with smaller values would have much larger coefficients (to indicate the same level of importance), which would increase the penalty for these predictors.

In our dataset, since each of the word counts are on different scales (i.e., some words are much more common than others), performing unregularized logistic regression wouldn't penalize each of the variables in a comparable way.

- (b) Ridge Regression
- Model Assessment

	accuracy	falseNegRate	falsePosRate
1	0.5925926	0	1

- Important Words

The 10 words with the largest coefficients (in absolute value) are shown below.

	word	coef	absCoef
1	februari	-0.01667849	0.01667849
2	upon	0.01572523	0.01572523
3	whilst	-0.01428038	0.01428038
4	within	-0.01367440	0.01367440
5	sever	-0.01344406	0.01344406
6	1783	-0.01302052	0.01302052
7	form	-0.01193143	0.01193143
8	member	-0.01168190	0.01168190
9	5	-0.01153821	0.01153821
10	although	-0.01127392	0.01127392

(c) Lasso

i. Model Assessment

	accuracy	falseNegRate	falsePosRate
1	0.8888889	0	0.2727273

ii. Important Words

The 10 words with the largest coefficients (in absolute value) are shown below.

	word	coef	absCoef
1	whilst	-1.3557246	1.3557246
2	februari	-1.1650795	1.1650795
3	upon	1.1493109	1.1493109
4	form	-0.5714197	0.5714197
5	although	-0.4741114	0.4741114
6	within	-0.4227536	0.4227536
7	sever	-0.3698089	0.3698089
8	lesser	-0.3483888	0.3483888
9	anim	-0.1694737	0.1694737
10	indirect	-0.1641048	0.1641048

Both ridge and lasso include “februari”, “upon”, “whilst”, “within”, “sever”, “form”, and “although” in the their 10 most important words.

“1783”, “member”, and “5” are included in the top 10 ridge words, but not in the top 10 lasso; and “lesser”, “anim”, and “indirect” are included in lasso, but not in ridge.

In general the ridge weights are smaller (in absolute value) than those from lasso.

Problem 7: Federalist Papers Classification with Feature Selection

(a)

Problem
7

Problem 8: James 8.4.10

Using boosting to predict Salary in the Hitters data set.

(a) Removing the observations with no salary information, and log-transforming the salaries.

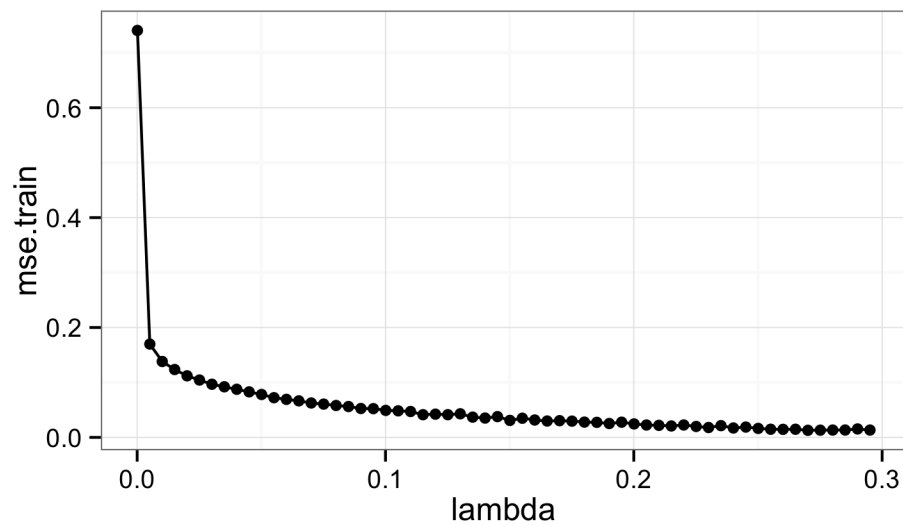
```
# remove observations with no salary info
Hitters <- Hitters[!is.na(Hitters$Salary), ]

# log transform the salaries
Hitters$Salary <- log(Hitters$Salary) #
```

(b) Creating training and testing sets.

```
Hitters.train <- Hitters[1:200, ]
Hitters.test <- Hitters[-(1:200), ]
```

(c) Boosting on the training set with 1,000 trees for a range of values of the shrinkage parameter $\lambda \in \{0.0001, 0.0051, 0.0101, 0.0151, \dots, 0.2951\}$. Training MSE vs λ shown below.



Problem 9: James 10.7.1

(a)

Problem
9

Todo list

Problem 3	2
Problem 7	5
Problem 9	6