

Homework 5

Statistics W4240: Data Mining

Columbia University

Due Tuesday, November 25 (All Sections)

For your .R submission, submit a single file labeled **hw05.R**. The write up should be saved as a .pdf of size less than 4MB. **DO NOT** submit .rar, .tar, .zip, .docx, or other file types.

Problem 1. (10 Points) James 6.8.1

Problem 2. (10 Points) James 6.8.2 (a)-(b)

Problem 3. (10 Points) James 6.8.4

Problem 4. (10 Points) James 8.4.3

Problem 5. (10 Points) James 8.4.5

Problems 6 and 7 use classification trees and logistic regression to classify the Federalist Papers.

Question 6. (20 Points) Use your code from Homework 4 to read in the Federalist Papers and create document term matrices `dtm.hamilton.train`, `dtm.hamilton.test`, `dtm.madison.train`, and `dtm.madison.test`. Create a set of labels for each document term matrix, with Madison documents given values 0 and Hamilton documents given values 1. Combine the document term matrices and labels to create two data frames: one that includes all training data and one that includes all testing data. Make sure that the column labels for the covariates are the dictionary words (hint: use `as.vector(dictionary$words)` to get a vector of words) and the column label for the response is `y`.

- (10 Points) Use tree classification to predict the author using the training data. Apply the model to the testing data. Specifically, in R use `rpart` classification with Gini impurity coefficient splits. Then compute the proportion classified correctly, the proportion of false negatives, and the proportion of false positives. Plot the tree with labeled splits.
- (10 Points) Now use tree classification again, but this time with information gain splits, to predict the author. Apply the model to the testing data. Then compute the proportion classified correctly, the proportion of false negatives, and the proportion of false positives. Plot the tree with labeled splits. Are there any differences between the two plots? If so, what are they and why do you think they arose?

Question 7. (20 Points) Create centered and scaled versions of your document term matrices. (Do not center and scale the labels.) We will use these for regularized logistic regression with `glmnet`.

- a. (6 Points) Could we use an unregularized logistic regression model with this data set? Why or why not?
- b. (7 Points) Use `glmnet` to fit a ridge regression model on the training data. Apply the model to the testing data. Then compute the proportion classified correctly, the proportion of false negatives, and the proportion of false positives. Find the 10 most important words according to the model along with their coefficients.
- c. (7 Points) Use `glmnet` to fit a lasso regression model on the training data. Apply the model to the testing data. Then compute the proportion classified correctly, the proportion of false negatives, and the proportion of false positives. Find the 10 most important words according to the model along with their coefficients. Compare the “important” words selected by ridge and lasso. Are the words different? What about their relative weights?