

Homework 06

STAT 430, Fall 2017

Due: Friday, October 27, 11:59 PM

Please see the [homework instructions document](#) for detailed instructions and some grading notes. Failure to follow instructions will result in point reductions.

For this homework we will use data found in [wisc-trn.csv](#) and [wisc-tst.csv](#) which contain train and test data respectively. [wisc.csv](#) is provided but not used. This is a modification of the Breast Cancer Wisconsin (Diagnostic) dataset from the UCI Machine Learning Repository. Only the first 10 feature variables have been provided. (And these are all you should use.)

- [UCI Page](#)
- [Data Detail](#)

You should consider coercing the response to be a factor variable.

You should use the `caret` package and training pipeline to complete this homework. Any time you use the `train()` function, first run `set.seed(1337)`.

Exercise 1 (Tuning KNN with caret)

[6 points] Train a KNN model using all available predictors, **no data preprocessing**, 5-fold cross-validation, and a well chosen value of the tuning parameter. Consider $k = 1, 3, 5, 7, \dots, 101$. Store the tuned model fit to the training data for later use. Plot the cross-validated accuracies as a function of the tuning parameter.

Exercise 2 (More Tuning KNN with caret)

[6 points] Train a KNN model using all available predictors, predictors scaled to have mean 0 and variance 1, 5-fold cross-validation, and a well chosen value of the tuning parameter. Consider $k = 1, 3, 5, 7, \dots, 101$. Store the tuned model fit to the training data for later use. Plot the cross-validated accuracies as a function of the tuning parameter.

Exercise 3 (Random Forest?)

[6 points] Now that we've introduced `caret`, it becomes extremely easy to try different statistical learning methods. Train a random forest using all available predictors, **no data preprocessing**, 5-fold cross-validation, and well a chosen value of the tuning parameter. Using `caret` to perform the tuning, there is only a single tuning parameter, `mtry`. Consider `mtry` values between 1 and 10. Store the tuned model fit to the training data for later use. Report the cross-validated accuracies as a function of the tuning parameter using a well formatted table.

Exercise 4 (Concept Checks)

[1 point each] Answer the following questions based on your results from the three exercises. Format your answer to this exercise as a table with one column indicating the part, and the other column for your answer. See the `rmarkdown` source for a template of this table.

- (a) What value of k is chosen for KNN without predictor scaling?
- (b) What is the cross-validated accuracy for KNN without predictor scaling?
- (c) What is the test accuracy for KNN without predictor scaling?
- (d) What value of k is chosen for KNN **with** predictor scaling?
- (e) What is the cross-validated accuracy for KNN **with** predictor scaling?
- (f) What is the test accuracy for KNN **with** predictor scaling?
- (g) Do you think that KNN is performing better with or without predictor scaling?
- (h) What value of `mtry` is chosen for the random forest?
- (i) Using the random forest, what is the (estimated) probability that the 10th observation of the test data is a cancerous tumor?
- (j) Using the random forest, what is the (test) sensitivity?
- (k) Using the random forest, what is the (test) specificity?
- (l) Based on these results, is the random forest or KNN model performing better?