

Homework 05

STAT 430, Fall 2017

Due: Friday, October 13, 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.

Exercise 1 (Detecting Cancer with KNN)

[7 points] For this exercise 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. Use KNN with all available predictors. For simplicity, do not scale the data. (In practice, scaling would slightly increase performance on this dataset.) Consider $k = 1, 3, 5, 7, \dots, 51$. Plot train and test error vs k on a single plot.

Use the seed value provided below for this exercise.

```
set.seed(314)
```

Exercise 2 (Logistic Regression Decision Boundary)

[5 points] Continue with the cancer data from Exercise 1. Now consider an additive logistic regression that considers only two predictors, `radius` and `symmetry`. Plot the test data with `radius` as the x axis, and `symmetry` as the y axis, with the points colored according to their tumor status. Add a line which represents the decision boundary for a classifier using 0.5 as a cutoff for predicted probability.

Exercise 3 (Sensitivity and Specificity of Cancer Detection)

[5 points] Continue with the cancer data from Exercise 1. Again consider an additive logistic regression that considers only two predictors, `radius` and `symmetry`. Report test sensitivity, test specificity, and test accuracy for three classifiers, each using a different cutoff for predicted probability:

- $c = 0.1$
- $c = 0.5$
- $c = 0.9$

Consider `M` to be the “positive” class when calculating sensitivity and specificity. Summarize these results using a single well-formatted table.

Exercise 4 (Comparing Classifiers)

[7 points] Use the data found in `hw05-trn.csv` and `hw05-tst.csv` which contain train and test data respectively. Use `y` as the response. Coerce `y` to be a factor after importing the data if it is not already.

Create pairs plot with ellipses for the training data, then train the following models using both available predictors:

- Additive Logistic Regression
- LDA (with Priors estimated from data)
- LDA with Flat Prior
- QDA (with Priors estimated from data)
- QDA with Flat Prior
- Naive Bayes (with Priors estimated from data)

Calculate test and train error rates for each model. Summarize these results using a single well-formatted table.

Exercise 5 (Concept Checks)

[1 point each] Answer the following questions based on your results from the three exercises.

- (a) Which k performs best in Exercise 1?
- (b) In Exercise 4, which model performs best?
- (c) In Exercise 4, why does Naive Bayes perform poorly?
- (d) In Exercise 4, which performs better, LDA or QDA? Why?
- (e) In Exercise 4, which prior performs better? Estimating from data, or using a flat prior? Why?
- (f) In Exercise 4, of the four classes, which is the easiest to classify?
- (g) [Not Graded] In Exercise 3, which classifier would be the best to use in practice?