

## COEN 140 Machine Learning and Data Mining

### Bonus Homework 4 (100 points)

**Due: 12:10pm, Thursday, Mar 1, 2018**

#### Spam classification using logistic regression

Consider the email spam data set. This consists of 4601 email messages, from which 57 features have been extracted. These are as follows:

- 48 features, in  $[0, 100]$ , giving the percentage of words in a given message which match a given word on the list. The list contains words such as “business”, “free”, “george”, etc. (The data was collected by George Forman, so his name occurs quite a lot.)
- 6 features, in  $[0, 100]$ , giving the percentage of characters in the email that match a given character on the list. The characters are ; ( [ ! \$ #
- Feature 55: The average length of an uninterrupted sequence of capital letters (max is 40.3, mean is 4.9)
- Feature 56: The length of the longest uninterrupted sequence of capital letters (max is 45.0, mean is 52.6)
- Feature 57: The sum of the lengths of uninterrupted sequence of capital letters (max is 25.6, mean is 282.2)

Download the data at <http://www.cse.scu.edu/~yfang/coen140/spambase.zip>. The data is split into a training set (of size 3065) and a test set (of size 1536).

One can imagine performing several kinds of preprocessing to this data. Try each of the following separately:

- a. Standardize the columns so they all have mean 0 and unit variance.
- b. Transform the features using  $\log(\mathbf{x}_{ij} + 0.1)$ .
- c. Binarize the features using  $I(\mathbf{x}_{ij} > 0)$ .

For each version of the data, fit a logistic regression model using gradient descent. Report the error rate on the training and test sets. Turn in your code and numerical results.