

## Exercise: Spam classification using logistic regression

Consider the email spam data set discussed on p300 of (?). 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 52.6, mean is 45.0)
- Feature 57: The sum of the lengths of uninterrupted sequence of capital letters (max is 25.6, mean is 282.2)

Load the data from `spamData.mat`, which contains 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:

1. Standardize the columns so they all have mean 0 and unit variance.
2. Transform the features using  $\log(x_{ij} + 0.1)$ .
3. Binarize the features using  $\mathbb{I}(x_{ij} > 0)$ .

For each version of the data, fit a logistic regression model. Use cross validation to choose the strength of the  $\ell_2$  regularizer. Report the mean error rate on the training and test sets. You should get numbers similar to this:

method	train	test
stnd	0.082	0.079
log	0.052	0.059
binary	0.065	0.072

(The precise values will depend on what regularization value you choose.) Turn in your code and numerical results. (See also Exercise “Spam classification using naive Bayes”.)