

# Machine Learning: Homework 3

## Instructions

- **Collaboration policy:** Homeworks must be done individually, except where otherwise noted in the assignments. “Individually” means each student must hand in their own answers, and each student must write and use their own code in the programming parts of the assignment. It is acceptable for students to collaborate in figuring out answers and to help each other solve the problems, though you must in the end write up your own solutions individually, and you must list the names of students you discussed this with. We will be assuming that, as participants in an undergraduate course, you will be taking the responsibility to make sure you personally understand the solution to any work arising from such collaboration.
- **Online submission:** You must submit your solutions online on [NYU Classes](#). You need to submit (1) a PDF which contains the solutions to all questions (2) `x.py` or `x.ipynb` files for the programming questions. We recommend that you use L<sup>A</sup>T<sub>E</sub>X, but we will accept scanned / pictured solutions as well.

## Problem 1: Logistic Regression

In problem1 and problem2, we cover two probabilistic linear classifiers we have covered in class so far. First, a discriminative linear classifier: logistic regression. Second, a generative linear classifier: Linear discriminant analysis (LDA). Both the algorithms find a linear decision boundary that separates the data into two classes, but make different assumptions. Our goal in this problem is to get a deeper understanding of the similarities and differences of these two algorithms.

1. **[15 Points]** We have seen in the lecture that given training data set  $(x_i, y_i)$ , the cross entropy loss for logistic regression is

$$J(w) = -\frac{1}{n} \sum_{i=1}^n [y_i \log(p_w(x_i)) + (1 - y_i) \log(1 - p_w(x_i))]$$

Recall that in lecture we were able to prove the gradient of  $J(w)$  w.r.t.  $w$  is

$$\frac{\partial J(w)}{\partial w} = -\frac{1}{n} \sum_{i=1}^n (y_i - p_w(x_i)) x_i$$

where  $w$  and  $x_i$  are all  $p \times 1$  vectors.

We have provided two data files `logistic.x.txt`, `logistic.y.txt`. These files contain the inputs ( $x_i \in R^2$ ) and outputs ( $y_i \in \{-1, 1\}$ ), respectively for a binary classification problem, with one training example per row. In `logistic_regression.ipynb`, implement gradient descent method to optimize  $J(w)$ , and apply it to fit a logistic regression model to the given data (remember to transform the target variable  $y$  to  $(y_i \in \{0, 1\})$ ). What are the coefficients  $w$  resulting from your fit? (Remember to include the intercept term.)

2. **[15 Points]** In `logistic_regression.ipynb`, plot the training data (your axes should be  $x_1$  and  $x_2$ , corresponding to the two coordinates of the inputs, and you should use a different symbol for each point plotted to indicate whether that example had label 1 or 0). Also plot on the same figure the decision boundary fit by logistic regression. (This should be a straight line showing the boundary separating the region where  $p_w(x) > 0.5$  from where  $p_w(x) \leq 0.5$ )

## Problem 2: Linear Discriminative Analysis

Recall that in linear discriminant analysis (LDA), we make prediction based on the following posterior probability

$$P(y = 1|x) = \frac{P(x|y = 1) * P(y = 1)}{P(x|y = 1) * P(y = 1) + P(x|y = 0)P(y = 0)}$$
$$P(x|y = 0) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0))$$
$$P(x|y = 1) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1))$$

Here  $\mu_0, \mu_1, \Sigma$  and the prior probability  $\phi = P(y = 1)$  are our model parameters (Note that for LDA we assume the negative class and the positive class share the same co-variance matrix  $\Sigma$ )

1. **[15 Points]** Suppose we have already learned  $\phi, \mu_0, \mu_1, \Sigma$  and now want to predict  $y$  given a new point  $x$ . To show that LDA results in a classifier that has a linear decision boundary, show the posterior distribution can be written as

$$P(y = 1|x; \mu_0, \mu_1, \Sigma) = \frac{1}{1 + \exp(-\theta_0 + \theta^T x)}$$

where  $\theta \in R^d$  and  $\theta_0 \in R$  are appropriate functions of  $\phi, \mu_0, \mu_1, \Sigma$

2. **[20 Points]** Coding problem. In `gda.ipynb`, fill in the code to calculate  $\phi, \mu_0, \mu_1, \Sigma$ . Note that the MLE for  $\mu_0, \mu_1, \Sigma$  can be obtained using

$$\hat{\mu}_0 = \frac{\sum_{X_i \in C_0} x_i}{n_0}$$
$$\hat{\mu}_1 = \frac{\sum_{X_i \in C_1} x_i}{n_1}$$
$$\hat{\Sigma} = \frac{\sum_{X_i \in C_0} (x_i - \hat{\mu}_0)(x_i - \hat{\mu}_0)^T + \sum_{X_i \in C_1} (x_i - \hat{\mu}_1)(x_i - \hat{\mu}_1)^T}{n}$$

where  $n_0$  and  $n_1$  are the number of training data points in the negative and positive class respectively. After we get the parameter estimation for  $\phi, \mu_0, \mu_1, \Sigma$ , we can make predictions based on the posterior probability. Fit the model using the breast cancer data (using the first two features only) that is available in `sklearn`. After fitting the model using training data, get the predictions for the test data, and print out the accuracy score for the test data. Plot the testing data and use a different symbol to indicate the ground truth label. Also plot on the same figure the decision boundary given by LDA. What are the coefficients of the linear decision boundary?

## Problem 3: Naive Bayes

In this problem we will implement the Naive Bayes classifier and apply it to spam email classification problem. Save your code in `naive_bayes.ipynb`.

**Data files:** We have provided you with two files: `spam_train.txt`, `spam_test.txt`. Each row of the data files corresponds to a single email. The first column gives the label (1=spam, 0=not spam).

**Pre-processing:** The dataset included for this exercise is based on a subset of the SpamAssassin Public Corpus. Figure 1 shows a sample email that contains a URL, an email address (at the end), numbers, and dollar amounts. While many emails would contain similar types of entities (e.g., numbers, other URLs, or other email addresses), the specific entities (e.g., the specific URL or specific dollar amount) will be different in almost every email. Therefore, one method often employed in processing emails is to "normalize" these

> Anyone knows how much it costs to host a web portal?  
 > Well, it depends on how many visitors youre expecting. This can be anywhere from  
 less than 10 bucks a month to a couple of \$100. You should checkout  
<http://www.rackspace.com/> or perhaps Amazon EC2 if youre running something big..  
 To unsubscribe yourself from this mailing list, send an email to: [groupname-unsubscribe@egroups.com](mailto:groupname-unsubscribe@egroups.com)

Figure 1: Sample e-mail in SpamAssassin corpus before pre-processing.

anyon know how much it cost to host a web portal well it depend on how mani visitor  
 your expect thi can be anywher from less than number buck a month to a coupl of  
 dollarnumb you should checkout httpaddr or perhap amazon ecnumb if your run  
 someth big to unsubscrib yourself from thi mail list send an email to emailaddr

Figure 2: Preprocessed version of the sample email from Figure1.

values, so that all URLs are treated the same, all numbers are treated the same, etc. For example, we could replace each URL in the email with the unique string "httpaddr" to indicate that a URL was present. This has the effect of letting the spam classifier make a classification decision based on whether any URL was present, rather than whether a specific URL was present. This typically improves the performance of a spam classifier, since spammers often randomize the URLs, and thus the odds of seeing any particular URL again in a new piece of spam is very small.

We have already implemented the following email preprocessing steps: lower-casing; removal of HTML tags; normalization of URLs, e-mail addresses, and numbers. In addition, words are reduced to their stemmed form. For example, "discount", "discounts", "discounted" and "discounting" are all replaced with "discount". Finally, we removed all non-words and punctuation. The result of these preprocessing steps is shown in Figure 2. below.

1. **[20 Points]** In order to train a Naive Bayes classifier, we need to transform the data into feature vectors. To do that we first build a vocabulary list using only the 5000 e-mail training set by finding all words that occur across the training set. Note that we assume that the data in the test sets is completely unseen when we train our model, and thus we do not use any information contained in them. **Ignore all words that appear in fewer than  $T=26$  e-mails of the 5000 e-mail in training set – this is both a means of preventing over-fitting (to be discussed in class) and of improving scalability.** To do this, write a function `create_wordlist(data)` which takes the training data as input and returns a Python list containing all the words that occur in at least 26 emails.

Then we will implement a Naive Bayes classifier only considering the words in the vocabulary. Take the first 4000 emails as training data and the last 1000 emails as validation data. Fit the Naive Bayes classifier using the first 4000 emails. Report the validation error rate on the last 1000 emails when threshold  $T=26$ .

$$\text{error rate} = \frac{\# \text{ incorrectly classified}}{\text{total number of classifications made}}$$

Recall that we make classification based on the posterior probability

$$P(y = 1|x) = \frac{P(x|y = 1)P(y = 1)}{P(x|y = 1)P(y = 1) + P(x|y = 0)P(y = 0)}$$

And under the independence assumption, we have

$$P(x|y) = P(x_1, x_2, \dots, x_d|y) = P(x_1|y)P(x_2|y)\dots P(x_d|y)$$

where  $d$  is the size of the vocabulary.

The maximum likelihood estimates for the parameters are

$$\begin{aligned} P(y = 1) &= \frac{\sum_{i=1}^n I(y_i = 1)}{n} \\ P(x_j = 1|y = 1) &= \frac{\sum_{i=1}^n I(y_i = 1 \ \& \ x_j = 1)}{\sum_{i=1}^n I(y_i = 1)} \\ P(x_j = 1|y = 0) &= \frac{\sum_{i=1}^n I(y_i = 0 \ \& \ x_j = 1)}{\sum_{i=1}^n I(y_i = 0)} \end{aligned}$$

In the situation where the testing data contain new words which the training data has never seen before, we may get 0 posterior probability. **To avoid that, we use the following equations as the estimation of the parameters:**

$$\begin{aligned} P(x_j = 1|y = 1) &= \frac{\sum_{i=1}^n I(y_i = 1 \ \& \ x_j = 1) + 1}{\sum_{i=1}^n I(y_i = 1) + 2} \\ P(x_j = 1|y = 0) &= \frac{\sum_{i=1}^n I(y_i = 0 \ \& \ x_j = 1) + 1}{\sum_{i=1}^n I(y_i = 0) + 2} \end{aligned}$$

2. **[15 Points]** To improve the accuracy, feel free to change the dimensionality of features by using a different vocabulary threshold  $T$ . For example, you can try threshold  $T$  from 1 to 34. Plot the training error rate and validation error rate vs. the threshold. Choose the threshold with the minimal validation error rate and report the corresponding test error rate.