

# ECE368: Probabilistic Reasoning

## Lab 1: Naïve Bayes Classifier

### 1 Naïve Bayes Classifier for Spam Filtering

In this lab, we use a Naïve Bayes Classifier to build a spam email filter based on whether and how many times each word in a fixed vocabulary occurs in the email. Suppose that we need to classify a set of  $N$  emails, and each email  $n$  is represented by  $\{\mathbf{x}_n, y_n\}$ ,  $n = 1, 2, \dots, N$ , where  $y_n$  is the class label which takes the value

$$y_n = \begin{cases} 1 & \text{if email } n \text{ is spam,} \\ 0 & \text{if email } n \text{ is non-spam (also called ham),} \end{cases} \quad (1)$$

and  $\mathbf{x}_n$  is a feature vector of the email  $n$ . We use a multinomial model to construct the feature vector  $\mathbf{x}_n$ . Let  $\mathcal{W} = \{w_1, w_2, \dots, w_D\}$  be the set of the words (called the vocabulary) that appear at least once in the training set. The feature vector  $\mathbf{x}_n$  is defined as a  $D$ -dimensional vector  $\mathbf{x}_n = [x_{n1}, x_{n2}, \dots, x_{nD}]$ , where each entry  $x_{nd}$ ,  $d = 1, 2, \dots, D$  is the number of occurrences of word  $w_d$  in email  $n$ . Thus the total number of words in email  $n$  can be expressed as  $l_n = x_{n1} + x_{n2} + \dots + x_{nD}$ .

We assume that each email  $n$  of length  $l_n$  is generated by a sequence of  $l_n$  independent events that randomly draw words from the vocabulary  $\mathcal{W}$ . (This is known as the naïve Bayes assumption.) For each event, let  $p(w_d | y_n = 1)$  be the probability that word  $w_d$  is picked, given that the email belongs to spam; let  $p(w_d | y_n = 0)$  be the probability that word  $w_d$  is picked, given that the email belongs to ham. Note that  $p(w_d | y_n = 1)$  and  $p(w_d | y_n = 0)$  are different, which gives us a way to classify spam vs. ham. For example, words like “dollar”, “winner” would be more likely to occur in spam than in ham. Also, note that both  $p(w_d | y_n = 1)$ ,  $d = 1, 2, \dots, D$  and  $p(w_d | y_n = 0)$ ,  $d = 1, 2, \dots, D$  should sum to one, i.e.,

$$\sum_{d=1}^D p(w_d | y_n = 1) = 1, \quad (2)$$

$$\sum_{d=1}^D p(w_d | y_n = 0) = 1. \quad (3)$$

The probabilities  $p(w_d | y_n = 1)$ ,  $p(w_d | y_n = 0)$ ,  $d = 1, \dots, D$  should be learned from the training data.

We make use of the word frequencies to model each email  $n$  probabilistically. Since each word in the email is seen as independently drawn from the vocabulary  $\mathcal{W}$ , the distribution of the feature vector  $\mathbf{x}_n$  given label  $y_n$  can be seen as a multinomial distribution as follows,

$$p(\mathbf{x}_n | y_n) = \frac{(x_{n1} + x_{n2} + \dots + x_{nD})!}{(x_{n1})!(x_{n2})! \dots (x_{nD})!} \prod_{d=1}^D p(w_d | y_n)^{x_{nd}}. \quad (4)$$

We assume that the prior class distribution  $p(y_n)$  is modeled as

$$p(y_n = 1) = \pi, \quad (5)$$

$$p(y_n = 0) = 1 - \pi, \quad (6)$$

where  $\pi$  is a fixed parameter (e.g., 0.5).

In the following, we first estimate the probabilities  $p(w_d | y_n = 1)$ ,  $p(w_d | y_n = 0)$ ,  $d = 1, \dots, D$  using the training set; we then build a classifier based on Bayes' rule and make predictions on the testing set.

Download `classifier.zip` from Quercus and unzip the file. The spam emails for training are in the subfolder `/data/spam/`. The ham emails for training are in the subfolder `/data/ham/`. The unlabeled emails for testing are in the subfolder `/data/testing/`.

Please answer the questions below and complete the routine `classifier.py`. File `util.py` contains a few functions/classes that will be helpful in writing the code for the classifier. *Please update the student numbers at the beginning of the code. This adjustment will create a unique dataset, ensuring each group gets distinct results. Implementing this change is essential for the assignment.*

## Questions

1. Training. We estimate the conditional probability distribution of the  $D$ -ary random variable as specified by  $p(w_d | y_n = 1)$  and  $p(w_d | y_n = 0)$ ,  $d = 1, \dots, D$ , from the training data using a *bag-of-words* model as follows. For notational simplicity, we define  $p_d = p(w_d | y_n = 1)$  and  $q_d = p(w_d | y_n = 0)$ .
  - (a) We put all the words from the spam emails in the training set in a bag and simply count the number of occurrences of each word  $w_d$ ,  $d = 1, \dots, D$ . We do the same for ham emails. The maximum likelihood estimates of  $p_d$  and  $q_d$  based on these counts are not the most appropriate to use when the probabilities are very close to 0 or to 1. For example, some words that occur in one class may not occur at all in the other class. In this problem, we use the technique of “Laplace smoothing” to deal with this problem. Please write down such an estimator for  $p_d$  and  $q_d$  as functions of the training data  $\{\mathbf{x}_n, y_n\}$ ,  $n = 1, 2, \dots, N$  using Laplace smoothing for the  $D$ -ary random variable.
  - (b) Complete the function `learn_distributions` in file `classifier.py`. In `learn_distributions`, you first build the vocabulary  $\{w_1, \dots, w_D\}$  by accounting for all the words that appear in the training set at least once; you then estimate  $p_d$  and  $q_d$ ,  $d = 1, 2, \dots, D$  using your expressions in part (a).
2. Testing. We classify the unlabeled emails in `/data/testing/` using the trained classifier.
  - (a) Let  $\{\mathbf{x}, y\}$  be a data point from the testing set whose class label  $y$  is unknown. Write down the maximum a posteriori (MAP) rule to decide whether  $y = 1$  or  $y = 0$  based on the feature vector  $\mathbf{x}$ . The  $d$ -th entry of  $\mathbf{x}$  is denoted by  $x_d$ . Please incorporate  $p_d$  and  $q_d$  in your expression. Please think carefully how to treat words that do not appear in either ham or spam training sets. Please assume that  $\pi = 0.5$ .
  - (b) Complete the function `classify_new_email` in file `classifier.py` to implement the MAP rule, and run it on the testing set. There are two types of errors in classifying unlabeled emails: Type 1 error is defined as the event that a spam email is misclassified as ham; Type 2 error is defined as the event that a ham email is misclassified as spam. Write down the numerical values of these two numbers of errors made by your classifier on the testing data. To avoid numerical underflow in your code, please work with the log probability  $\log p(y|\mathbf{x})$  in your code.
  - (c) In practice, Type 1 error and Type 2 error lead to different consequences (or costs). Therefore, we may wish to trade off one type of error against the other in designing the classifier. For example, we usually want to achieve a very low Type 2 error since the cost of missing a useful email can be severe, while we can tolerate a relative high Type 1 error as it merely causes inconvenience. Please provide a way to modify the decision rule in the classifier such that these two types of error can be traded off. In other words, change the decision rule in a way such that Type 2 error would decrease at a cost of Type 1 error, and vice versa. Test your method on the testing set and provide the following plot: Let the  $x$ -axis be the number of Type 1 errors and the  $y$ -axis be the number of Type 2 errors in the testing data set. Plot at least 10 points corresponding to different

pairs of Type 1 and Type 2 errors, as a result of adjusting the classification rule. The two end points of the plot should be: 1) the one with zero Type 1 error; and 2) the one with zero Type 2 error. The code should be included in file `classifier.py`.

3. Why do we need Laplace smoothing? Briefly explain what would go wrong if we do use the maximum likelihood estimators in the training process.

*This lab is adapted from the originals by Greg Wornell and Wei Yu.*