

# ECE368: Probabilistic Reasoning

## Lab 1: Classification with Multinomial and Gaussian Models

Name: Ryan Song

Student Number: 1003833658

You should hand in: 1) A scanned .pdf version of this sheet with your answers (file size should be under 2 MB); 2) one figure for Question 1.2.(c) and two figures for Question 2.1.(c) in the .pdf format; and 3) two Python files classifier.py and lda\_qda.py that contain your code. All these files should be uploaded to Quercus.

### 1 Naïve Bayes Classifier for Spam Filtering

- (a) Write down the estimators for  $p_d$  and  $q_d$  as functions of the training data  $\{x_n, y_n\}, n = 1, 2, \dots, N$  using the technique of "Laplace smoothing". (1 pt)

$$p_d = \frac{\# \text{ of occurrences in spam} + 1}{\text{Total \# words in spam} + \text{Total \# of distinct words}}$$

$$q_d = \frac{\# \text{ of occurrences in HAM} + 1}{\text{Total \# in ham} + \text{Total \# of distinct words}}$$

- (b) Complete function learn\_distributions in python file classifier.py based on the expressions. (1 pt)
- (a) Write down the MAP rule to decide whether  $y = 1$  or  $y = 0$  based on its feature vector  $\mathbf{x}$  for a new email  $\{x, y\}$ . The  $d$ -th entry of  $\mathbf{x}$  is denoted by  $x_d$ . Please incorporate  $p_d$  and  $q_d$  in your expression. Please assume that  $\pi = 0.5$ . (1 pt)

$$\sum_{i=0}^D x_i \cdot \log(p_{-d}[\text{word in question}]) \stackrel{?}{\geq} \sum_{i=0}^D x_i \cdot \log(q_{-d}[\text{word}])$$

the word that  $x_i$  is describing.

- (b) Complete function classify\_new\_email in classifier.py, and test the classifier on the testing set. The number of Type 1 errors is 2, and the number of Type 2 errors is 4. (1 pt)
- (c) Write down the modified decision rule in the classifier such that these two types of error can be traded off. Please introduce a new parameter to achieve such a trade-off. (0.5 pt)

introduce parameter  $\gamma$ :

$$\sum_{i=0}^D x_i \log(p_{-d}[w]) \stackrel{?}{\geq} \gamma \cdot \sum_{i=0}^D x_i \log(q_{-d}[w])$$

use  $\gamma$  to tradeoff the types of errors.

Write your code in file classifier.py to implement your modified decision rule. Test it on the testing set and plot a figure to show the trade-off between Type 1 error and Type 2 error. In the figure, the  $x$ -axis should be the number of Type 1 errors and the  $y$ -axis should be the number of Type 2 errors. Plot at least 10 points corresponding to different pairs of these two types of error in your figure. The two end points of the plot should be: 1) the point with zero Type 1 error; and 2) the point with zero Type 2 error. Please save the figure with name **nbc.pdf**. (1 pt)



- (d) If we do not use Laplace smoothing and simply use maximum likelihood estimation in the training phase, what will go wrong? What kind of emails such a classifier would fail to classify? (0.5 pt)

Any <sup>new</sup> email which has a word that exists in SPAM but not HAM will be classified as SPAM and vice versa. In our specific estimator, it wouldn't even work as  $\log(0)$  is undefined.

## 2 Linear/Quadratic Discriminant Analysis for Height/Weight Data

1. (a) Write down the maximum likelihood estimates of the parameters  $\mu_m$ ,  $\mu_f$ ,  $\Sigma$ ,  $\Sigma_m$ , and  $\Sigma_f$  as functions of the training data  $\{x_n, y_n\}, n = 1, 2, \dots, N$ . (1 pt)

$$\begin{aligned}\hat{\mu}_m &= \frac{1}{n_m} \sum_{x_i \in \text{all men}} x_i & \hat{\mu}_f &= \frac{1}{n_f} \sum_{x_i \in \text{all female}} x_i \\ \hat{\Sigma} &= \frac{1}{n} \sum_{x_i \in \text{all samples}} (x_i - \hat{\mu}_t)(x_i - \hat{\mu}_t)^T & \hat{\mu}_t &= \frac{1}{n} \sum_{x_i \in \text{all samples}} x_i \\ \hat{\Sigma}_m &= \frac{1}{n_m} \sum_{x_i \in \text{all men}} (x_i - \hat{\mu}_m)(x_i - \hat{\mu}_m)^T \\ \hat{\Sigma}_f &= \frac{1}{n_f} \sum_{x_i \in \text{all female}} (x_i - \hat{\mu}_f)(x_i - \hat{\mu}_f)^T\end{aligned}$$

- (b) In the case of LDA, write down the decision boundary as a linear equation of  $x$  with parameters  $\mu_m$ ,  $\mu_f$ , and  $\Sigma$ . Note that we assume  $\pi = 0.5$ . (0.5 pt)

$$\begin{aligned}\log \pi_m - \frac{1}{2} \mu_m^T \Sigma^{-1} \mu_m + (\Sigma^{-1} \mu_c)^T x &\stackrel{m}{\geq} \log \pi_f - \frac{1}{2} \mu_f^T \Sigma^{-1} \mu_f + (\Sigma^{-1} \mu_f)^T x \\ (\Sigma^{-1} \mu_c)^T x - \frac{1}{2} \mu_m^T \Sigma^{-1} \mu_m &\stackrel{m}{\geq} (\Sigma^{-1} \mu_f)^T x - \frac{1}{2} \mu_f^T \Sigma^{-1} \mu_f\end{aligned}$$

In the case of QDA, write down the decision boundary as a quadratic equation of  $x$  with parameters  $\mu_m$ ,  $\mu_f$ ,  $\Sigma_m$ , and  $\Sigma_f$ . Note that we assume  $\pi = 0.5$ . (0.5 pt)

$$\begin{aligned}-\frac{1}{2} \log |\Sigma_m| - \frac{1}{2} (x - \mu_m)^T \Sigma_m^{-1} (x - \mu_m) \\ \stackrel{m}{\geq} -\frac{1}{2} \log |\Sigma_f| - \frac{1}{2} (x - \mu_f)^T \Sigma_f^{-1} (x - \mu_f)\end{aligned}$$

- (c) Complete function `discrimAnalysis` in `lda_qda.py` to visualize LDA and QDA models and the corresponding decision boundaries. Please name the figures as `lda.pdf`, and `qda.pdf`. (1 pt)

2. The misclassification rates are 13/110 for LDA, and 13/110 for QDA. (1 pt)