ECE368: Probabilistic Reasoning
## Lab 1: Classification with Multinomial and Gaussian Models

Name: Raymond Hong          Student Number: 1003942629

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 ldaqda.py that contain your code. All these files should be uploaded to Quercus.

# 1 Naïve Bayes Classifier for Spam Filtering

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

$$p_d = \frac{\# \text{ of occurances of word } d \text{ in spam } + 1}{\text{total } \# \text{ of words in spam total distinct } \# \text{ of words in spam } \& \text{ ham}}$$

$$w_d = \frac{\# \text{ of occurances of word } d \text{ in ham} + 1}{\text{total } \# \text{ of words in ham } + \text{total } \# \text{ of distinct words in spam } \& \text{ ham}}$$

(b) Complete function learn_distributions in python file classifier.py based on the expressions. (1 pt)

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

$$y = \arg\max_y P(x|y) = \arg\max_y \frac{(x_1 + \cdots + x_D)!}{(x_1)! (x_2)! (\cdots x_D)!} \prod_{d=1}^{D} P(x_d|y)^{x_d}$$

$$\prod_{d=1}^{D} (p_d)^{x_d} \underset{\text{ham}}{\overset{\text{spam}}{\gtrless}} \prod_{d=1}^{D} (q_d)^{x_d}$$

$r = 1$

(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 ]. (1.5 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 "ratio" parameter. ratio = 1 for previous sections

$$\frac{\prod_{d=1}^{D} (p_d)^{x_d}}{\prod_{d=1}^{D} (q_d)^{x_d}} \underset{\text{ham}}{\overset{\text{spam}}{\gtrless}} r$$

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)

# 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 $\{x_n, y_n\}, n = 1, 2, \ldots, N$. (1 pt)

$$\underline{\mu_m} = \frac{1}{\#\text{of males}} \sum_{i=1}^{\hat{n}} \mathbb{1}\{y_i = 1\} \underline{X_i}$$

$$\underline{\mu_f} = \frac{1}{\#\text{of females}} \sum_{i=1}^{\hat{n}} \mathbb{1}\{y_i = 2\} \underline{X_i}$$

$$\Sigma = \frac{1}{N} \sum_{i=1}^{n} (\underline{X_i} - \underline{\mu_m})(\underline{X_i} - \underline{\mu_m})^T \mathbb{1}(y_i = m) + (\underline{X_i} - \underline{\mu_f})(\underline{X} - \underline{\mu_f})^T \mathbb{1}(y_i = f)$$

$$\Sigma_m = \frac{1}{\#\text{of male}} \sum_{i=1}^{n} (\underline{X_i} - \underline{\mu_m})(\underline{X_i} - \underline{\mu_m})^T \mathbb{1}\{y_i = 1\}$$

$$\Sigma_f = \frac{1}{\#\text{of females}} \sum_{i=1}^{n} (\underline{X_i} - \underline{\mu_f})(\underline{X_i} - \underline{\mu_f})^T \mathbb{1}\{y_i = 2\}$$

(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)

$$\underline{\mu_m}^T \Sigma^{-1} \underline{X} - \frac{1}{2} \underline{\mu_m}^T \Sigma^{-1} \underline{\mu_m} = \underline{\mu_f} \Sigma^{-1} \underline{X} - \frac{1}{2} \underline{\mu_f}^T \Sigma^{-1} \underline{\mu_f}$$

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)

$$-\frac{1}{2} \log |\Sigma_m| - \frac{1}{2}(\underline{X} - \underline{\mu_m})^T \Sigma_m^{-1}(\underline{X} - \underline{\mu_m}) = \frac{1}{2} \log |\Sigma_f| - \frac{1}{2}(\underline{X} - \underline{\mu_f})^T \Sigma_f^{-1}(\underline{X} - \underline{\mu_f})$$

(c) Complete function discrimAnalysis in ldaqda.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 $\boxed{0.118}$ for LDA, and $\boxed{0.05}$ for QDA. (1 pt)

trade off of error type 1 and error type 2

lda

qda