



Department of Electrical Engineering  
School of Science and Engineering

## EE514/CS535 Machine Learning

### HOMEWORK 2

---

**Due Date:** 23:55, Friday, March 26, 2021 (Submit online on LMS)

**Format:** 5 problems, for a total of 80 marks

**Contribution to Final Percentage:** 2.5%

**Instructions:**

- Each student must submit his/her own hand-written assignment, scanned in a **single PDF document**.
  - You are allowed to collaborate with your peers but copying your colleague's solution is strictly prohibited. Anybody found guilty would be subjected to disciplinary action in accordance with the university rules and regulations.
  - Note: Vectors are written in lowercase and bold in the homework, for your written submission kindly use an underline instead. In addition, use capital letters for matrices and lowercase for scalars.
- 

#### Problem 1 (10 marks)

**Bayes' Theorem** - Suppose that you have built an email classifier that filters emails as 'spam' and 'not spam'. It is estimated that 80% of all emails are not spam. You claim that your classifier accurately predicts spam 90% of the time, and only incorrectly predicts as spam 1% of the time.

Using the information provided, deduce the probability of a non-spam email, given that it has been classified as spam by your classifier.

## Problem 2 (15 marks)

**Maximum Likelihood Estimation (MLE) of Poisson Distribution** - The probability mass function (PMF) of a random variable  $X$  that follows the Poisson distribution is given by:

$$p_X(x, \lambda) = \frac{e^{-\lambda} \lambda^x}{x!}$$

where  $x$  is a non-negative integer, i.e.,  $x \in \mathbb{Z}^+$  and  $\lambda$ , which parameterizes the distribution, is the mean of the distribution.

You are provided  $n$  independent observations of this Poisson distribution. Formulate MLE problem for the estimation of parameter  $\lambda$  and prove that it is simply given by the sample mean of the observations, that is:

$$\hat{\lambda} = \frac{1}{n} \sum_{i=1}^n x_i$$

*Hint: Follow the same steps as the example given in slides, take log likelihood function of the PMF, then derivate it with respect to  $\lambda$  and equate to zero.*

**Problem 3 (25 marks)**

**Text Data Classification** - In this question, we will dry-run the Naïve Bayes' (parts (b)-(d)) classification algorithm to predict if the sentiment of the text data provided. You are provided the following reviews and labels:

Sentiment	Text	
Training	Positive	very enjoyable
	Positive	really surprising
	Positive	really fun and enjoyable
	Negative	not surprising
	Negative	very boring and predictable
Test	?	pretty predictable and no fun

**Table 1:** Review data

- (a) [5 marks] First let us use the trained weights provided in Table 2 to predict the sentiment of the test data. We label a positive review as '1' and negative review as '-1', take the decision boundary to be 0. You are provided a description of the features and their respective trained weights in the table below. Use this information to predict the sentiment of the test review. (Make appropriate assumptions about positive and negative words)

$\theta_i$	Feature Description	Weight Value
$\theta_0$	Bias	0.6
$\theta_1$	Count of Positive Words	1.2
$\theta_2$	Count of Negative Words	-3.5
$\theta_3$	log(word count)	0.1

**Table 2:** Features and Trained Weights

- (b) [10 marks] We now use Naïve Bayes to predict the sentiment of the review, use 'Laplace add-1 smoothing' on Table 1 to compute the likelihoods for all words in training data. Populate the table given below:

Word	P(Word  +)	P(Word  -)
very		
enjoyable		
⋮		

**Table 3:** Likelihoods of Training Data

- (c) [5 marks] Using the table in part (b), predict the sentiment of the test data.
- (d) [5 marks] Comment on the usefulness of Naïve Bayes by focusing on problems associated with predicting labels without the assumption of conditional independence. You may choose to use the problem provided to you in this question to augment your statements.

#### Problem 4 (20 marks)

**Multinomial Logistic Regression Formulation** - In this question we will formulate an image classification problem through logistic regression. We will classify fruit images as “orange”, “apple”, “strawberry”, “mango”, or “grape”. Assume that the image data provided to us is in the form of  $n$  grey-scale  $1024 \times 1024$  matrices. We can define our feature matrix for the  $i$ -th image as:  $\mathbf{x}_i \in \mathbb{R}^{1024 \times 1024}$ .

- (a) [3 marks] How will you convert  $\mathbf{x}_i$  feature matrix to a feature vector? What will be the size of this feature vector?
- (b) [5 marks] Define a label to index mapping for this model. Use it to convert the following ordered labels ( $n = 8$ ) into a vector,  $\mathbf{y} \in \mathbb{R}^8$ .

orange, mango, strawberry, apple, apple, grape, mango, strawberry

- (c) [3 marks] What will be the shape of our weights matrix  $\Theta$ ?
- (d) [3 marks] Write down the shape of our prediction matrix, i.e., the product of our feature matrix (with bias) and weights matrix.
- (e) [6 marks] Finally, you pass your prediction matrix through the activation function, softmax. How will you extract the predicted label from this resulting matrix?

**Problem 5** (10 marks)

**Regularized Logistic Regression** - For binary logistic regression, we obtain  $\boldsymbol{\theta}$  by solving the following optimization problem (see lecture notes for the notation adopted here).

$$\underset{\boldsymbol{\theta}}{\text{minimize}} \quad \mathcal{L}(\boldsymbol{\theta}) = - \sum_{i=1}^n y_i \log(h_{\boldsymbol{\theta}}(\mathbf{x}_i)) + (1 - y_i) \log(1 - h_{\boldsymbol{\theta}}(\mathbf{x}_i))$$

We can extend the regularization introduced for linear regression for logistic regression. For example  $L_2$  regularized logistic regression can be formulated by adding  $L_2$  norm of the weight vector  $\boldsymbol{\theta}$  as a penalty term  $\mathcal{R}(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_2^2$  and formulating the optimization problem as

$$\underset{\boldsymbol{\theta}}{\text{minimize}} \quad \mathcal{L}(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_2^2$$

We derived the gradient descent step size in the lectures for the loss function with penalty. Here we require you to derive the gradient of the regularized loss function with respect to  $\boldsymbol{\theta}$ .

— End of Homework —