

Name: \_\_\_\_\_

University of Illinois

Spring 2020

CS 446/ECE 449 Machine Learning  
Homework 4: Multiclass Logistic Regression

Due on Thursday February 27 2020, noon Central Time

1. [16 points] Multiclass Logistic Regression

We are given a dataset  $\mathcal{D} = \left\{ \left( \begin{bmatrix} 1 \\ 0 \end{bmatrix}, 0 \right), \left( \begin{bmatrix} 0 \\ 1 \end{bmatrix}, 1 \right), \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, 2 \right) \right\}$  containing three pairs  $(x, y)$ , where each  $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathbb{R}^2$  denotes a 2-dimensional point and  $y \in \{0, 1, 2\}$ .

We want to train by minimizing the negative log-likelihood the parameters  $w$  (includes bias) of a multi-class logistic regression classifier using

$$\min_w - \sum_{(x,y) \in \mathcal{D}} \log p(y|x) \quad \text{where} \quad p(y|x) = \frac{\exp w_y^\top \begin{bmatrix} x \\ 1 \end{bmatrix}}{\sum_{\hat{y} \in \{0,1,2\}} \exp w_{\hat{y}}^\top \begin{bmatrix} x \\ 1 \end{bmatrix}}. \quad (1)$$

- (a) (2 points) How many parameters do we train, *i.e.*, what's the domain of  $w$ ? Explain what  $w_y$  means and how it relates to  $w$ ?

Your answer:

- (b) (2 points) Alternatively, we can use the equivalent probability model

$$p(y|x) = \frac{\exp w^\top \psi(x, y)}{\sum_{\hat{y} \in \{0,1,2\}} \exp w^\top \psi(x, \hat{y})}.$$

Explain how we need to construct  $\psi(x, y)$  such that  $w^\top \psi(x, y) = w_y^\top \begin{bmatrix} x \\ 1 \end{bmatrix} \forall y \in \{0, 1, 2\}$ .

Your answer:

Name: \_\_\_\_\_

- (c) (3 points) Alternatively, we can use the equivalent probability model

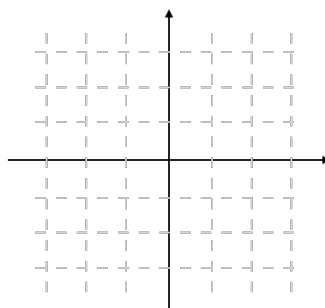
$$p(y|x) = \frac{\exp F(y, w, x)}{\sum_{\hat{y} \in \{0,1,2\}} \exp F(\hat{y}, w, x)} \quad \text{with} \quad F(y, w, x) = [\mathbf{W}x + b]_y,$$

where  $\mathbf{W}$  is a matrix of weights and  $b$  is a vector of biases. The notation  $[a]_y$  extracts the  $y$ -th entry from vector  $a$ . What are the dimensions of  $\mathbf{W}$  and  $b$  and how does  $\mathbf{W}$  and  $b$  related to the originally introduced  $w$ ?

Your answer:

- (d) (6 points) Assume we are given  $\mathbf{W} = \begin{bmatrix} 3 & 0.5 \\ 0 & 1 \\ -1.5 & -1.5 \end{bmatrix}$  and  $b = \begin{bmatrix} 0 \\ 0 \\ 1.5 \end{bmatrix}$ . Draw the datapoints, and the lines  $[\mathbf{W}x + b]_y = 0 \forall y \in \{0, 1, 2\}$  in  $x_1$ - $x_2$ -space and explain whether these weights result in correct prediction for all datapoints in  $\mathcal{D}$ ?

Your answer: **Mark the axis**



Name: \_\_\_\_\_

- (e) (3 points) Complete `A4_Multiclass.py`. After optimizing, what values do you obtain for  $\mathbf{W}$ ,  $b$  and what probability estimates  $p(\hat{y}|x)$  do you obtain for all points  $x \in \mathcal{D}$  in the dataset and for all classes  $\hat{y} \in \{0, 1, 2\}$ . (**Hint:** a total of nine probability estimates are required.)

Your answer: