Boston University
Department of Electrical and Computer Engineering

### ENG EC 414 Introduction to Machine Learning

### HW 4

© 2015 – 2020 Prakash Ishwar
© 2020 Francesco Orabona

**Issued:** Fri 25 Sept 2020          **Due:** 10:00am Fri 2 Oct 2020 in Gradescope (non-code) + Blackboard

**Important:** Before you proceed, please read the documents pertaining to *Homework formatting and submission guidelines* in the Homeworks section of Blackboard. **In particular, for computer assignments you are prohibited from using any online code or built-in MATLAB functions except as indicated in the problem or skeleton code (when provided).**

**Note:** Problem difficulty = number of coffee cups ☕

**Problem 4.1** [15pts] *(Hyperplanes)*
Consider two-class classification with the following training feature vectors: $\mathbf{x}_P = (2,0)^\top, \mathbf{x}_Q = (0,4)^\top, \mathbf{x}_R = (3,3)^\top, \mathbf{x}_S = (7,5)^\top$ with labels $-1, -1, +1, +1$ respectively.

  (a) [2pts] Hand-plot the training set. Proper labeling of axes and key points is needed to receive full credit.

  (b) [1pt] Hand-sketch the hyperplane with parameters $\mathbf{w} = (3,0)^\top$, $b = -3$. Proper labeling of axes and key points is needed to receive full credit.

  (c) [2pts] Hand-compute the $\ell_2$ distance of $\mathbf{x}_P$ from the hyperplane in part (b).

  (d) [2pts] Determine if the hyperplane in part (b) linearly separates the training set. Explanation is needed to receive credit.

  (e) [3pts] Hand-compute the parameters of the hyperplane passing through $\mathbf{x}_P$ and $\mathbf{x}_Q$.

  (f) [5pts] Hand-compute the orthogonal projection of $\mathbf{x}_R$ onto the hyperplane in part (f).

**Problem 4.2** [8pts] *(Logistic regression by hand)*
Let
$$S = \left\{(x_1 = -3, y_1 = -1), (x_2 = 5, y_2 = +1)\right\}$$
be a training set of pairs of scalar features and their labels for training a binary logistic regression classifier with class labels $-1, +1$.

  (a) [2pts] Let $F(w,b) = \frac{1}{m}\sum_{i=1}^m \ln(1 + \exp(-y_i(b + wx_i)))$ be the objective function for one-dimensional logistic regression. Compute the expression of the partial derivative of this function with respect to $w$ and $b$ using the training set above.

  (b) [2pts] With initialization $w_1 = 0$ and $b_1 = 0$, compute $w_2$ and $b_2$, that is the value after one iteration of the gradient descent algorithm for minimizing the objective function, expressing them as functions of the step-size $\eta > 0$.

(c) [4pts] Fix $b = 0$ and find analitically the minimum of the objective function on the above training set and the optimal value of $w$.

**Problem 4.3** [21pts] *(Gradient Descent for Logistic regression)*
In this problem, we will implement the Gradient Descent (GD) Algorithm to learn the parameters of a logistic regression classifier. We will use again the Adult dataset.
*Summary of cost function and algorithm*: Let $\mathcal{S} = \{(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_m, y_m)\}$ be a training set of $m$ examples consisting of feature vector $\mathbf{x}_j \in \mathbb{R}^d$ and label $y_j \in \{-1, 1\}$ for each $j$. Let $\tilde{\mathbf{x}}_j = \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix} \in \mathbb{R}^{d+1}$ denote the augmented representation of the $j$-th feature vector. The pseudocode for implemeting GD for logistic regression is as follows:

---
**Algorithm 1** Gradient Descent for Logistic Regression

---
1: **Input:** Training set $\{(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_m, y_m)\}$, learning rate $\eta > 0$, maximum number of iterations $T$, initial hyperplane $\mathbf{w}_1$, initial bias $b_1$
2: Set $\tilde{\mathbf{w}}_1 = \begin{bmatrix} b_1 \\ \mathbf{w}_1 \end{bmatrix} \in \mathbb{R}^{d+1}$
3: Construct augmented training features: $\tilde{\mathbf{x}}_1, \ldots, \tilde{\mathbf{x}}_m$
4: **for** $t = 1, 2, \ldots, T$ **do**
5:     Calculate value of objective function: $obj_t = \sum_{i=1}^{m} \ln(1 + \exp(-y_i \tilde{\mathbf{w}}_t^\top \tilde{\mathbf{x}}_i))$
6:     Compute gradient: $\tilde{\mathbf{g}}_t = -\sum_{i=1}^{m} \frac{y_i \tilde{\mathbf{x}}_i}{1 + \exp(y_i \tilde{\mathbf{w}}_t^\top \tilde{\mathbf{x}}_i)} \in \mathbb{R}^{d+1}$
7:     Gradient descent step: $\tilde{\mathbf{w}}_{t+1} = \tilde{\mathbf{w}}_t - \eta \tilde{\mathbf{g}}_t$
8: **end for**
9: **return** Output: Extract $\mathbf{w}_{T+1}$ and $b_{T+1}$ from $\tilde{\mathbf{w}}_{T+1}$ and return them

---

(a) [6pts] Complete the skeleton of `train_logistic_regression_gd.m` implementing Algorithm 1 in a function with prototype

```
[w, b, obj] = train_logistic_regression_gd(X, y, eta, T, w1, b1)
```

where $\mathbf{w} \in \mathbb{R}^d$ and $b \in \mathbb{R}$ are the hyperplane vector and the bias found by GD, `obj` is the vector of objective function values, `X` is the matrix of the training samples $\in \mathbb{R}^{m \times d}$, `y` is the vector of labels $\in \mathbb{R}^m$, `eta` is the learning rate, `T` is the maximum number of iterations to run, `w1` $\in \mathbb{R}^d$ is the initialization for $\mathbf{w}$, and `b1` $\in \mathbb{R}$ is the initialization for $b$. Note: It will be faster if you use Matlab vectorized operations to calculate the gradient, but it is also fine to use for loops.

(b) [4pts] In class, we said that using a random initialization to minimize a convex function should be avoided. Here, we will test this claim. Complete the skeleton code in `problem_4_3b.m` to use your implementation of logistic regression on the Adult dataset for $T = 1000$ iterations, $\eta = 1/5000$, and 10 random initializations using Gaussian with zero mean and variance 1 on each coordinate of $\mathbf{w}$ and on $b$ (You can generate Gaussian noise in Matlab with `randn`).

(c) [2pts] Plot the 10 lines of the values of the objective function during training in the point above. What do you observe? Based on these results and what we said in class, why random initialization should be avoided? (To better observe the difference among the lines, use `loglog` that uses a logarithmic scale on both axes.)

(d) [4pts] Complete the skeleton code in `problem_4_3d.m` to use your implementation of GD for logistic regression, starting from the zero vector and zero bias, $T = 1000$ iterations, $\eta = 1/5000$, and test the

obtained solution on the test set. Count the number of prediction mistakes on the test set, predicting for each sample the class with the highest probability.

(e) [5pts] ☞ If the learning rate is too big, you should get a lot of "Inf". Verify it yourself using for example $\eta = 0.1$. Propose a way to fix this issue. Hint: First, locate the <u>cause</u> of the problem, then try to understand <u>why</u> we get this problem numerically. Finally, think how to fix the problem. You can propose approximations.

**Code-submission via Blackboard:** You must prepare 3 files: `train_logistic_regression_gd.m` for Problem 4.3(a), `problem_4_3b.m` for Problem 4.3(b), and `problem_4_3d.m` for Problem 4.3(d). Place them in a **single** directory which should be zipped and uploaded into Blackboard. Your directory must be named as follows: `<yourBUemailID>_hwX` where `X` is the homework number. For example, if your BU email address is `charles500@bu.edu` then for homework number 4 you would submit a single directory named: `charles500_hw4.zip` which contains all the MATLAB code (and only the code).

Three corresponding skeleton code files are provided for your reference. Reach-out to the TAs via Piazza and their office/discussion hours for questions related to coding.