

## Homework 2: Due Friday Nov. 3, 11:59 PM

**Instructions:** upload a PDF report using L<sup>A</sup>T<sub>E</sub>X containing your answers to Canvas (remember to include your name and ID number).

## Problem 1. True or False

Decide whether the following statements are true or false. **Justify your answers.**

- (a) (10 pt) If classifier  $A$  has smaller training error than classifier  $B$ , then classifier  $A$  will have smaller generalization (test) error than classifier  $B$ .
- (b) (10 pt) It is not always good to use model with high complexity.
- (c) (10 pt) Gradient descent needs to decrease the learning rate (step size) in order to converge to the optima.

## Problem 2. Multiple choice questions

Choose the correct answer and **justify your answer.**

- (a) (20 pt) Which of the following is not a possible growth function  $m_{\mathcal{H}}(N)$  for some hypothesis set? (1)  $2^N$   
(2)  $2^{\lfloor \sqrt{N} \rfloor}$  (3) 1 (4)  $N^2 - N + 2$  (5) none of the other choices

## Problem 3. L2-Regularized Logistic Regression

Given a set of instance-label pairs  $(\mathbf{x}_i, y_i)$ ,  $i = 1, \dots, n$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \{+1, -1\}$ , L2-regularized logistic regression estimates the model  $\mathbf{w}$  by solving the following optimization problem:

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) := \left\{ \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i)) \right\} \quad (1)$$

We assume data matrix  $X \in \mathbb{R}^{n \times d}$  is sparse, each column of  $X$  has  $n_j$  nonzero elements, and each row of  $X$  has  $d_i$  nonzero elements. The whole training dataset has  $\text{nnz}(X) := \sum_{j=1}^d n_j = \sum_{i=1}^n d_i$  nonzero elements.

- (a) (20 pt) Derive the gradient and Hessian of  $f(\mathbf{w})$ .
- (b) (5 pt) What is the update rule of gradient descent (using a fixed step size  $\eta$ )
- (c) (5 pt) What is the time complexity of one gradient descent update?

Newton method is a classical second order method for minimizing  $f(\mathbf{w})$ . The update rule for Newton method is:

$$\mathbf{w} \leftarrow \mathbf{w} + \eta \mathbf{d}^* \quad (2)$$

where  $\mathbf{d}^* = -\nabla^2 f(\mathbf{w})^{-1} \nabla f(\mathbf{w})$

- (d) (5 pt) Assume we first form the Hessian matrix  $\nabla^2 f(\mathbf{w})$  and then compute the Newton direction  $(\nabla^2 f(\mathbf{w}))^{-1} \nabla f(\mathbf{w})$ . What is the time complexity of one Newton update (eq. (2)) for L2-regularized logistic regression? (Assume  $n$  is close to  $d$ ).
- (e) (5 pt) The update rule in eq. (2) can also be written as solving the following optimization problem:

$$\mathbf{d}^* = \underset{\mathbf{d}}{\operatorname{argmin}} \left\{ \frac{1}{2} \mathbf{d}^T \nabla^2 f(\mathbf{w}) \mathbf{d} + \nabla f(\mathbf{w})^T \mathbf{d} \right\} := J(\mathbf{d}) \quad (3)$$

Proof the optimal solution of (3) is  $-(\nabla^2 f(\mathbf{w}))^{-1} \nabla f(\mathbf{w})$ .

- (f) (10 pt) Since the matrix inversion would be numerically unstable in certain condition, what is the alternative solution to get  $(\nabla^2 f(\mathbf{w}))^{-1} \nabla f(\mathbf{w})$  without matrix inversion?