

# CS 229, Autumn 2016

## Problem Set #0: Linear Algebra and Multivariable Calculus

---

**Notes:** (1) These questions require thought, but do not require long answers. Please be as concise as possible. (2) If you have a question about this homework, we encourage you to post your question on our Piazza forum, at <https://piazza.com/stanford/autumn2016/cs229>. (3) If you missed the first lecture or are unfamiliar with the collaboration or honor code policy, please read the policy on Handout #1 (available from the course website) before starting work. (4) This specific homework is *not graded*, but we encourage you to solve each of the problems to brush up on your linear algebra. Some of them may even be useful for subsequent problem sets. It also serves as your introduction to using Gradescope for submissions.

If you are scanning your document by cellphone, please check the Piazza forum for recommended cellphone scanning apps and best practices.

### 1. [0 points] Gradients and Hessians

Recall that a matrix  $A \in \mathbb{R}^{n \times n}$  is *symmetric* if  $A^T = A$ , that is,  $A_{ij} = A_{ji}$  for all  $i, j$ . Also recall the gradient  $\nabla f(x)$  of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , which is the  $n$ -vector of partial derivatives

$$\nabla f(x) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(x) \\ \vdots \\ \frac{\partial}{\partial x_n} f(x) \end{bmatrix} \quad \text{where } x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}.$$

The hessian  $\nabla^2 f(x)$  of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is the  $n \times n$  symmetric matrix of twice partial derivatives,

$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2}{\partial x_1^2} f(x) & \frac{\partial^2}{\partial x_1 \partial x_2} f(x) & \cdots & \frac{\partial^2}{\partial x_1 \partial x_n} f(x) \\ \frac{\partial^2}{\partial x_2 \partial x_1} f(x) & \frac{\partial^2}{\partial x_2^2} f(x) & \cdots & \frac{\partial^2}{\partial x_2 \partial x_n} f(x) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial x_n \partial x_1} f(x) & \frac{\partial^2}{\partial x_n \partial x_2} f(x) & \cdots & \frac{\partial^2}{\partial x_n^2} f(x) \end{bmatrix}.$$

- Let  $f(x) = \frac{1}{2}x^T A x + b^T x$ , where  $A$  is a symmetric matrix and  $b \in \mathbb{R}^n$  is a vector. What is  $\nabla f(x)$ ?
- Let  $f(x) = g(h(x))$ , where  $g : \mathbb{R} \rightarrow \mathbb{R}$  is differentiable and  $h : \mathbb{R}^n \rightarrow \mathbb{R}$  is differentiable. What is  $\nabla f(x)$ ?
- Let  $f(x) = \frac{1}{2}x^T A x + b^T x$ , where  $A$  is symmetric and  $b \in \mathbb{R}^n$  is a vector. What is  $\nabla^2 f(x)$ ?
- Let  $f(x) = g(a^T x)$ , where  $g : \mathbb{R} \rightarrow \mathbb{R}$  is continuously differentiable and  $a \in \mathbb{R}^n$  is a vector. What are  $\nabla f(x)$  and  $\nabla^2 f(x)$ ? (*Hint:* your expression for  $\nabla^2 f(x)$  may have as few as 11 symbols, including ' and parentheses.)

### 2. [0 points] Positive definite matrices

A matrix  $A \in \mathbb{R}^{n \times n}$  is *positive semi-definite* (PSD), denoted  $A \succeq 0$ , if  $A = A^T$  and  $x^T A x \geq 0$  for all  $x \in \mathbb{R}^n$ . A matrix  $A$  is *positive definite*, denoted  $A \succ 0$ , if  $A = A^T$  and  $x^T A x > 0$  for

all  $x \neq 0$ , that is, all non-zero vectors  $x$ . The simplest example of a positive definite matrix is the identity  $I$  (the diagonal matrix with 1s on the diagonal and 0s elsewhere), which satisfies  $x^T I x = \|x\|_2^2 = \sum_{i=1}^n x_i^2$ .

- (a) Let  $z \in \mathbb{R}^n$  be an  $n$ -vector. Show that  $A = zz^T$  is positive semidefinite.
- (b) Let  $z \in \mathbb{R}^n$  be a *non-zero*  $n$ -vector. Let  $A = zz^T$ . What is the null-space of  $A$ ? What is the rank of  $A$ ?
- (c) Let  $A \in \mathbb{R}^{n \times n}$  be positive semidefinite and  $B \in \mathbb{R}^{m \times n}$  be arbitrary, where  $m, n \in \mathbb{N}$ . Is  $BAB^T$  PSD? If so, prove it. If not, give a counterexample with explicit  $A, B$ .

### 3. [0 points] Eigenvectors, eigenvalues, and the spectral theorem

The eigenvalues of an  $n \times n$  matrix  $A \in \mathbb{R}^{n \times n}$  are the roots of the characteristic polynomial  $p_A(\lambda) = \det(\lambda I - A)$ , which may (in general) be complex. They are also defined as the values  $\lambda \in \mathbb{C}$  for which there exists a vector  $x \in \mathbb{C}^n$  such that  $Ax = \lambda x$ . We call such a pair  $(x, \lambda)$  an *eigenvector, eigenvalue* pair. In this question, we use the notation  $\text{diag}(\lambda_1, \dots, \lambda_n)$  to denote the diagonal matrix with diagonal entries  $\lambda_1, \dots, \lambda_n$ , that is,

$$\text{diag}(\lambda_1, \dots, \lambda_n) = \begin{bmatrix} \lambda_1 & 0 & 0 & \cdots & 0 \\ 0 & \lambda_2 & 0 & \cdots & 0 \\ 0 & 0 & \lambda_3 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \lambda_n \end{bmatrix}.$$

- (a) Suppose that the matrix  $A \in \mathbb{R}^{n \times n}$  is diagonalizable, that is,  $A = T\Lambda T^{-1}$  for an invertible matrix  $T \in \mathbb{R}^{n \times n}$ , where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$  is diagonal. Use the notation  $t^{(i)}$  for the columns of  $T$ , so that  $T = [t^{(1)} \ \cdots \ t^{(n)}]$ , where  $t^{(i)} \in \mathbb{R}^n$ . Show that  $At^{(i)} = \lambda_i t^{(i)}$ , so that the eigenvalues/eigenvector pairs of  $A$  are  $(t^{(i)}, \lambda_i)$ .

A matrix  $U \in \mathbb{R}^{n \times n}$  is orthogonal if  $U^T U = I$ . The spectral theorem, perhaps one of the most important theorems in linear algebra, states that if  $A \in \mathbb{R}^{n \times n}$  is symmetric, that is,  $A = A^T$ , then  $A$  is *diagonalizable by a real orthogonal matrix*. That is, there are a diagonal matrix  $\Lambda \in \mathbb{R}^{n \times n}$  and orthogonal matrix  $U \in \mathbb{R}^{n \times n}$  such that  $U^T A U = \Lambda$ , or, equivalently,

$$A = U \Lambda U^T.$$

Let  $\lambda_i = \lambda_i(A)$  denote the  $i$ th eigenvalue of  $A$ .

- (b) Let  $A$  be symmetric. Show that if  $U = [u^{(1)} \ \cdots \ u^{(n)}]$  is orthogonal, where  $u^{(i)} \in \mathbb{R}^n$  and  $A = U \Lambda U^T$ , then  $u^{(i)}$  is an eigenvector of  $A$  and  $Au^{(i)} = \lambda_i u^{(i)}$ , where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ .
- (c) Show that if  $A$  is PSD, then  $\lambda_i(A) \geq 0$  for each  $i$ .

# CS 229, Autumn 2016

## Problem Set #1: Supervised Learning

---

**Due Wednesday, October 19 at 11:00 am on Gradescope.**

**Notes:** (1) These questions require thought, but do not require long answers. Please be as concise as possible. (2) If you have a question about this homework, we encourage you to post your question on our Piazza forum, at <http://piazza.com/stanford/autumn2016/cs229>. (3) If you missed the first lecture or are unfamiliar with the collaboration or honor code policy, please read the policy on Handout #1 (available from the course website) before starting work. (4) For problems that require programming, please include in your submission a copy of your code (with comments) and any figures that you are asked to plot. If typing your solutions, include your code as text in your PDF. Do not submit extra files. (5) To account for late days, the due date listed on Gradescope is October 22 at 11 am. If you submit after October 19, you will begin consuming your late days. If you wish to submit on time, submit before October 19 at 11 am.

All students must submit an electronic PDF version. We highly recommend typesetting your solutions via latex. If you are scanning your document by cell phone, please check the Piazza forum for recommended scanning apps and best practices.

### 1. [25 points] Logistic regression

(a) [10 points] Consider the average empirical loss (the risk) for logistic regression:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y^{(i)} \theta^T x^{(i)}}) = -\frac{1}{m} \sum_{i=1}^m \log(h_{\theta}(y^{(i)} x^{(i)}))$$

where  $h_{\theta}(x) = g(\theta^T x)$  and  $g(z) = 1/(1 + e^{-z})$ . Find the Hessian  $H$  of this function, and show that for any vector  $z$ , it holds true that

$$z^T H z \geq 0.$$

*Hint:* You might want to start by showing the fact that  $\sum_i \sum_j z_i x_i x_j z_j = (x^T z)^2 \geq 0$ .

**Remark:** This is one of the standard ways of showing that the matrix  $H$  is positive semi-definite, written “ $H \succeq 0$ .” This implies that  $J$  is convex, and has no local minima other than the global one.<sup>1</sup> If you have some other way of showing  $H \succeq 0$ , you’re also welcome to use your method instead of the one above.

(b) [10 points] We have provided two data files:

- [http://cs229.stanford.edu/ps/ps1/logistic\\_x.txt](http://cs229.stanford.edu/ps/ps1/logistic_x.txt)
- [http://cs229.stanford.edu/ps/ps1/logistic\\_y.txt](http://cs229.stanford.edu/ps/ps1/logistic_y.txt)

These files contain the inputs ( $x^{(i)} \in \mathbb{R}^2$ ) and outputs ( $y^{(i)} \in \{-1, 1\}$ ), respectively for a binary classification problem, with one training example per row. Implement<sup>2</sup> Newton’s method for optimizing  $J(\theta)$ , and apply it to fit a logistic regression model to the data. Initialize Newton’s method with  $\theta = \vec{0}$  (the vector of all zeros). What are the coefficients  $\theta$  resulting from your fit? (Remember to include the intercept term.)

<sup>1</sup>If you haven’t seen this result before, please feel encouraged to ask us about it during office hours.

<sup>2</sup>Write your own version, and do not call a built-in library function.

- (c) [5 points] Plot the training data (your axes should be  $x_1$  and  $x_2$ , corresponding to the two coordinates of the inputs, and you should use a different symbol for each point plotted to indicate whether that example had label 1 or -1). Also plot on the same figure the decision boundary fit by logistic regression. (This should be a straight line showing the boundary separating the region where  $h_\theta(x) > 0.5$  from where  $h_\theta(x) \leq 0.5$ .)