

# CS 229, Autumn 2009

## Problem Set #1: Supervised Learning

---

**Due in class (9:30am) on Wednesday, October 14.**

**Notes:** (1) These questions require thought, but do not require long answers. Please be as concise as possible. (2) When sending questions to `cs229-qa@stanford.edu`, please make sure to write the homework number and the question number in the subject line, such as `Hwk1 Q4`, and send a separate email per question. (3) If you missed the first lecture or are unfamiliar with the class' 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 printout of your code (with comments) and any figure that you are asked to plot.

**SCPD students:** Please fax your solutions to Prof. Ng at (650) 725-1449, and write "ATTN: CS229 (Machine Learning)" on the cover sheet. If you are writing your solutions out by hand, please write clearly and in a reasonably large font using a dark pen to improve legibility.

### 1. [25 points] Logistic regression

- (a) [10 points] Consider the log-likelihood function for logistic regression:

$$\ell(\theta) = \sum_{i=1}^m y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log(1 - h(x^{(i)}))$$

Find the Hessian  $H$  of this function, and show that for any vector  $z$ , it holds true that

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

[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 negative semi-definite, written " $H \leq 0$ ." This implies that  $\ell$  is concave, and has no local maxima other than the global one.<sup>1</sup> If you have some other way of showing  $H \leq 0$ , you're also welcome to use your method instead of the one above.

- (b) [10 points] On the Leland system, the files `/afs/ir/class/cs229/ps/ps1/q1x.dat` and `/afs/ir/class/cs229/ps/ps1/q1y.dat` contain the inputs ( $x^{(i)} \in \mathbb{R}^2$ ) and outputs ( $y^{(i)} \in \{0, 1\}$ ) respectively for a binary classification problem, with one training example per row. Implement<sup>2</sup> Newton's method for optimizing  $\ell(\theta)$ , 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.)
- (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 0). Also plot on the

<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.

same figure the decision boundary fit by logistic regression. (I.e., this should be a straight line showing the boundary separating the region where  $h(x) > 0.5$  from where  $h(x) \leq 0.5$ .)

## 2. [27 points] Locally weighted linear regression

Consider a linear regression problem in which we want to “weight” different training examples differently. Specifically, suppose we want to minimize

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m w^{(i)} \left( \theta^T x^{(i)} - y^{(i)} \right)^2.$$

In class, we worked out what happens for the case where all the weights (the  $w^{(i)}$ 's) are the same. In this problem, we will generalize some of those ideas to the weighted setting, and also implement the locally weighted linear regression algorithm.

- (a) [2 points] Show that  $J(\theta)$  can also be written

$$J(\theta) = (X\theta - \vec{y})^T W (X\theta - \vec{y})$$

for an appropriate diagonal matrix  $W$ , and where  $X$  and  $\vec{y}$  are as defined in class. State clearly what  $W$  is.

- (b) [7 points] If all the  $w^{(i)}$ 's equal 1, then we saw in class that the normal equation is

$$X^T X \theta = X^T \vec{y},$$

and that the value of  $\theta$  that minimizes  $J(\theta)$  is given by  $(X^T X)^{-1} X^T \vec{y}$ . By finding the derivative  $\nabla_{\theta} J(\theta)$  and setting that to zero, generalize the normal equation to this weighted setting, and give the new value of  $\theta$  that minimizes  $J(\theta)$  in closed form as a function of  $X$ ,  $W$  and  $\vec{y}$ .

- (c) [6 points] Suppose we have a training set  $\{(x^{(i)}, y^{(i)}); i = 1 \dots, m\}$  of  $m$  independent examples, but in which the  $y^{(i)}$ 's were observed with differing variances. Specifically, suppose that

$$p(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma^{(i)}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2(\sigma^{(i)})^2}\right)$$

I.e.,  $y^{(i)}$  has mean  $\theta^T x^{(i)}$  and variance  $(\sigma^{(i)})^2$  (where the  $\sigma^{(i)}$ 's are fixed, known, constants). Show that finding the maximum likelihood estimate of  $\theta$  reduces to solving a weighted linear regression problem. State clearly what the  $w^{(i)}$ 's are in terms of the  $\sigma^{(i)}$ 's.

- (d) [12 points] On the Leland computer system, the files `/afs/ir/class/cs229/ps/ps1/q2x.dat` and `/afs/ir/class/cs229/ps/ps1/q2y.dat` contain the inputs  $(x^{(i)})$  and outputs  $(y^{(i)})$  for a regression problem, with one training example per row.
- [2 points] Implement (unweighted) linear regression ( $y = \theta^T x$ ) on this dataset (using the normal equations), and plot on the same figure the data and the straight line resulting from your fit. (Remember to include the intercept term.)
  - [7 points] Implement locally weighted linear regression on this dataset (using the weighted normal equations you derived in part (b)), and plot on the same figure

the data and the curve resulting from your fit. When evaluating  $h(\cdot)$  at a query point  $x$ , use weights

$$w^{(i)} = \exp\left(-\frac{(x - x^{(i)})^2}{2\tau^2}\right),$$

with a bandwidth parameter  $\tau = 0.8$ . (Again, remember to include the intercept term.)

- iii. [3 points] Repeat (ii) four times, with  $\tau = 0.1, 0.3, 2$  and  $10$ . Comment **briefly** on what happens to the fit when  $\tau$  is too small or too large.

3. [18 points] **Poisson regression and the exponential family**

- (a) [5 points] Consider the Poisson distribution parameterized by  $\lambda$ :

$$p(y; \lambda) = \frac{e^{-\lambda} \lambda^y}{y!}.$$

Show that the Poisson distribution is in the exponential family, and clearly state what are  $b(y)$ ,  $\eta$ ,  $T(y)$ , and  $a(\eta)$ .

- (b) [3 points] Consider performing regression using a GLM model with a Poisson response variable. What is the canonical response function for the family? (You may use the fact that a Poisson random variable with parameter  $\lambda$  has mean  $\lambda$ .)
- (c) [10 points] For a training set  $\{(x^{(i)}, y^{(i)}); i = 1, \dots, m\}$ , let the log-likelihood of an example be  $\log p(y^{(i)} | x^{(i)}; \theta)$ . By taking the derivative of the log-likelihood with respect to  $\theta_j$ , derive the stochastic gradient ascent rule for learning using a GLM model with Poisson responses  $y$  and the canonical response function.
- (d) [5 extra credit points] Consider using GLM with a response variable from any member of the exponential family in which  $T(y) = y$ , and the canonical response function for the family. Show that stochastic gradient ascent on the log-likelihood  $\log p(\bar{y} | X, \theta)$  results in the update rule  $\theta_i := \theta_i - \alpha(h(x) - y)x_i$ .

4. [15 points] **Gaussian discriminant analysis**

Suppose we are given a dataset  $\{(x^{(i)}, y^{(i)}); i = 1, \dots, m\}$  consisting of  $m$  independent examples, where  $x^{(i)} \in \mathbb{R}^n$  are  $n$ -dimensional vectors, and  $y^{(i)} \in \{0, 1\}$ . We will model the joint distribution of  $(x, y)$  according to:

$$\begin{aligned} p(y) &= \phi^y (1 - \phi)^{1-y} \\ p(x|y=0) &= \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1} (x - \mu_0)\right) \\ p(x|y=1) &= \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\right) \end{aligned}$$

Here, the parameters of our model are  $\phi$ ,  $\Sigma$ ,  $\mu_0$  and  $\mu_1$ . (Note that while there're two different mean vectors  $\mu_0$  and  $\mu_1$ , there's only one covariance matrix  $\Sigma$ .)

- (a) [5 points] Suppose we have already fit  $\phi$ ,  $\Sigma$ ,  $\mu_0$  and  $\mu_1$ , and now want to make a prediction at some new query point  $x$ . Show that the posterior distribution of the label at  $x$  takes the form of a logistic function, and can be written

$$p(y=1|x; \phi, \Sigma, \mu_0, \mu_1) = \frac{1}{1 + \exp(-\theta^T x)},$$

where  $\theta$  is some appropriate function of  $\phi, \Sigma, \mu_0, \mu_1$ . (Note: To get your answer into the form above, for this part of the problem only, you may have to redefine the  $x^{(i)}$ 's to be  $n + 1$ -dimensional vectors by adding the extra coordinate  $x_0^{(i)} = 1$ , like we did in class.)

- (b) [10 points] For this part of the problem only, you may assume  $n$  (the dimension of  $x$ ) is 1, so that  $\Sigma = [\sigma^2]$  is just a real number, and likewise the determinant of  $\Sigma$  is given by  $|\Sigma| = \sigma^2$ . Given the dataset, we claim that the maximum likelihood estimates of the parameters are given by

$$\begin{aligned}\phi &= \frac{1}{m} \sum_{i=1}^m 1\{y^{(i)} = 1\} \\ \mu_0 &= \frac{\sum_{i=1}^m 1\{y^{(i)} = 0\}x^{(i)}}{\sum_{i=1}^m 1\{y^{(i)} = 0\}} \\ \mu_1 &= \frac{\sum_{i=1}^m 1\{y^{(i)} = 1\}x^{(i)}}{\sum_{i=1}^m 1\{y^{(i)} = 1\}} \\ \Sigma &= \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T\end{aligned}$$

The log-likelihood of the data is

$$\begin{aligned}\ell(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^m p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \\ &= \log \prod_{i=1}^m p(x^{(i)}|y^{(i)}; \mu_0, \mu_1, \Sigma)p(y^{(i)}; \phi).\end{aligned}$$

By maximizing  $\ell$  with respect to the four parameters, prove that the maximum likelihood estimates of  $\phi$ ,  $\mu_0$ ,  $\mu_1$ , and  $\Sigma$  are indeed as given in the formulas above. (You may assume that there is at least one positive and one negative example, so that the denominators in the definitions of  $\mu_0$  and  $\mu_1$  above are non-zero.)

- (c) [5 extra credit points] Without assuming that  $n = 1$ , show that the maximum likelihood estimates of  $\phi$ ,  $\mu_0$ ,  $\mu_1$ , and  $\Sigma$  are as given in the formulas in part (b). [Note: If you're fairly sure that you have the answer to this part right, you don't have to do part (b), since that's just a special case.]

## 5. [12 points] Linear invariance of optimization algorithms

Consider using an iterative optimization algorithm (such as Newton's method, or gradient descent) to minimize some continuously differentiable function  $f(x)$ . Suppose we initialize the algorithm at  $x^{(0)} = \vec{0}$ . When the algorithm is run, it will produce a value of  $x \in \mathbb{R}^n$  for each iteration:  $x^{(1)}, x^{(2)}, \dots$

Now, let some non-singular square matrix  $A \in \mathbb{R}^{n \times n}$  be given, and define a new function  $g(z) = f(Az)$ . Consider using the same iterative optimization algorithm to optimize  $g$  (with initialization  $z^{(0)} = \vec{0}$ ). If the values  $z^{(1)}, z^{(2)}, \dots$  produced by this method necessarily satisfy  $z^{(i)} = A^{-1}x^{(i)}$  for all  $i$ , we say this optimization algorithm is **invariant to linear reparameterizations**.

- (a) [9 points] Show that Newton's method (applied to find the minimum of a function) is invariant to linear reparameterizations. Note that since  $z^{(0)} = \vec{0} = A^{-1}x^{(0)}$ , it is sufficient to show that if Newton's method applied to  $f(x)$  updates  $x^{(i)}$  to  $x^{(i+1)}$ , then Newton's method applied to  $g(z)$  will update  $z^{(i)} = A^{-1}x^{(i)}$  to  $z^{(i+1)} = A^{-1}x^{(i+1)}$ .<sup>3</sup>
- (b) [3 points] Is gradient descent invariant to linear reparameterizations? Justify your answer.

**Reminder:** Please include in your submission a printout of your code and figures for the programming questions.

---

<sup>3</sup>Note that for this problem, you must explicitly prove any matrix calculus identities that you wish to use that are not given in the lecture notes.