

## Instructions

**Submission:** Assignment submission will be via `courses.usciden.net`. By the submission date, there will be a folder set up in which you can submit your files. Please be sure to follow all directions outlined here.

You can submit multiple files, but only the last one submitted counts. That means if you finish some problems and want to submit something now, and then a later file when you finish, that's fine. If I were taking the class, that's what I'd do: that way, if I forget to finish the homework or something happens (remember Murphy's Law), I still get credit for what I finished and turned in. Remember, there are no grace days on problem sets, just on programming assignments!

Problem sets must be typewritten or neatly handwritten when submitted; if the grader cannot read your handwriting on a problem, they may elect to grade it as a zero. If it is handwritten, your submission must still be submitted as a PDF.

It is strongly recommended that you typeset with  $\text{\LaTeX}$  and use that to generate a PDF file.

- The file should be named as `firstname_lastname_USCID.pdf` (e.g., `Jenny_Tutone_8675309.pdf`).
- Do not have any spaces in your file name when uploading it.
- Please include your name and USCID in the header of the report as well.

There are many free integrated  $\text{\LaTeX}$  editors that are convenient to use. Choose the one(s) you like the most. This <http://www.andy-roberts.net/writing/latex> seems like a good tutorial.

**Collaboration:** You may discuss with your classmates. However, you need to write your own solutions and submit separately. Also in your report, you need to list with whom you have discussed for each problem. Please consult the syllabus for what is and is not acceptable collaboration. Review the rules on academic conduct in the syllabus: a single instance of plagiarism can adversely affect you significantly more than you could stand to gain.

## Notes on notation:

- Unless stated otherwise, scalars are denoted by small letter in normal font, vectors are denoted by small letters in bold font and matrices are denoted by capital letters in bold font.
- The bias term is subsumed in the input vector, so the input vector is actually  $x = [x', 1]^T$ , unless mentioned otherwise.
- $\|\cdot\|$  means L2-norm unless specified otherwise i.e.  $\|\cdot\| = \|\cdot\|_2$

## Problem 1 Neural networks (error-backpropagation, initialization, and nonlinearity)

[Recommended maximum time spent: 1 hour]

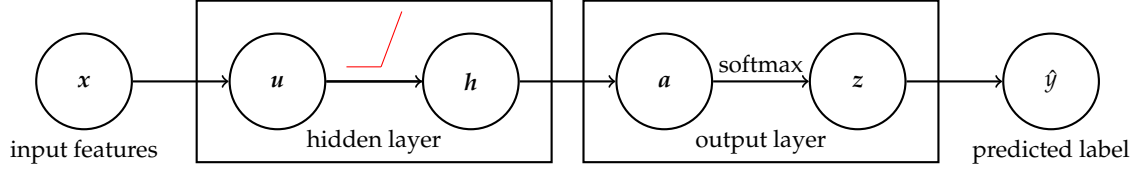


Figure 1: A diagram of a 1 hidden-layer multi-layer perceptron (MLP). The edges mean mathematical operations, and the circles mean variables. Generally we call the combination of a linear (or affine) operation and a nonlinear operation (like element-wise sigmoid or the rectified linear unit (relu) operation as in eq. (3)) as a hidden layer.

In the lecture, we have talked about error-backpropagation, a way to compute partial derivatives (or gradients) w.r.t the parameters of a neural network. We have also mentioned that optimization is challenging and nonlinearity is important for neural networks. In this question, you are going to (Q1.1) practice error-backpropagation, (Q1.2) investigate how initialization affects optimization, and (Q1.3) the importance of nonlinearity.

Specifically, you are given the following 1-hidden layer multi-layer perceptron (MLP) for a  $K$ -class classification problem (see Fig. 1 for illustration and details), and  $(x \in \mathbb{R}^D, y \in \{1, 2, \dots, K\})$  is a labeled instance,

$$x \in \mathbb{R}^D \quad (1)$$

$$u = W^{(1)}x + b^{(1)}, \quad W^{(1)} \in \mathbb{R}^{M \times D} \text{ and } b^{(1)} \in \mathbb{R}^M \quad (2)$$

$$h = \max\{0, u\} = \begin{bmatrix} \max\{0, u_1\} \\ \vdots \\ \max\{0, u_M\} \end{bmatrix} \quad (3)$$

$$a = W^{(2)}h + b^{(2)}, \quad W^{(2)} \in \mathbb{R}^{K \times M} \text{ and } b^{(2)} \in \mathbb{R}^K \quad (4)$$

$$z = \begin{bmatrix} \frac{e^{a_1}}{\sum_k e^{a_k}} \\ \vdots \\ \frac{e^{a_K}}{\sum_k e^{a_k}} \end{bmatrix} \quad (5)$$

$$\hat{y} = \arg \max_k z_k. \quad (6)$$

For  $K$ -class classification problem, one popular loss function for training is the cross-entropy loss,

$$l = - \sum_k \mathbf{1}[y == k] \log z_k, \quad (7)$$

$$\text{where } \mathbf{1}[\text{True}] = 1; \text{ otherwise, } 0. \quad (8)$$

For ease of notation, let us define the one-hot (i.e., 1-of- $K$ ) encoding

$$y \in \mathbb{R}^K \text{ and } y_k = \begin{cases} 1, & \text{if } y = k, \\ 0, & \text{otherwise.} \end{cases} \quad (9)$$

so that

$$l = -\sum_k y_k \log z_k = -\mathbf{y}^T \begin{bmatrix} \log z_1 \\ \vdots \\ \log z_K \end{bmatrix} = -\mathbf{y}^T \log \mathbf{z}. \quad (10)$$

**Q1.1** Assume that you have computed  $\mathbf{u}, \mathbf{h}, \mathbf{a}, \mathbf{z}$ , given  $(\mathbf{x}, \mathbf{y})$ . Please first express  $\frac{\partial l}{\partial \mathbf{u}}$  in terms of  $\frac{\partial l}{\partial \mathbf{a}}, \mathbf{u}, \mathbf{h}$ , and  $\mathbf{W}^{(2)}$ .

$$\frac{\partial l}{\partial \mathbf{u}} = ?$$

Then express  $\frac{\partial l}{\partial \mathbf{a}}$  in terms of  $\mathbf{z}$  and  $\mathbf{y}$ .

$$\frac{\partial l}{\partial \mathbf{a}} = ?$$

Finally, compute  $\frac{\partial l}{\partial \mathbf{W}^{(1)}}$  and  $\frac{\partial l}{\partial \mathbf{b}^{(1)}}$  in terms of  $\frac{\partial l}{\partial \mathbf{u}}$  and  $\mathbf{x}$ . Compute  $\frac{\partial l}{\partial \mathbf{W}^{(2)}}$  in terms of  $\frac{\partial l}{\partial \mathbf{a}}$  and  $\mathbf{h}$ .

$$\begin{aligned} \frac{\partial l}{\partial \mathbf{W}^{(1)}} &= ? \\ \frac{\partial l}{\partial \mathbf{b}^{(1)}} &= ? \\ \frac{\partial l}{\partial \mathbf{W}^{(2)}} &= ? \end{aligned}$$

You are encouraged to use matrix/vector forms to simplify your answers. Note that  $\max\{0, u\}$  is not differentiable w.r.t.  $u$  at  $u = 0$ . Please note that

$$\frac{\partial \max\{0, u\}}{\partial u} = \begin{cases} 1, & \text{if } u > 0, \\ 0, & \text{if } u \leq 0, \end{cases} \quad (11)$$

which stands for the Heaviside step function. You can use

$$\frac{\partial \max\{0, u\}}{\partial u} = H(u) \quad (12)$$

in your derivation of  $\frac{\partial l}{\partial \mathbf{u}}$ .

You can also use  $\cdot *$  to represent element-wise product between two vectors or matrices. For example,

$$\mathbf{v} \cdot * \mathbf{c} = \begin{bmatrix} v_1 \times c_1 \\ \vdots \\ v_I \times c_I \end{bmatrix} \in \mathbb{R}^I, \text{ where } \mathbf{v} \in \mathbb{R}^I \text{ and } \mathbf{c} \in \mathbb{R}^I. \quad (13)$$

What to submit: No more than 5 lines of derivation for each of the 5 partial derivatives.

**(In the nominator layout)**

By the chain rule,

$$\begin{aligned}\frac{\partial l}{\partial \mathbf{h}} &= \frac{\partial l}{\partial \mathbf{a}} \mathbf{W}^{(2)}, \\ \frac{\partial l}{\partial \mathbf{u}} &= \frac{\partial l}{\partial \mathbf{h}} H(\mathbf{u}), \\ \frac{\partial l}{\partial \mathbf{u}} &= \frac{\partial l}{\partial \mathbf{a}} \mathbf{W}^{(2)} H(\mathbf{u}).\end{aligned}$$

$$\begin{aligned}\frac{\partial l}{\partial \mathbf{b}^{(1)}} &= \frac{\partial l}{\partial \mathbf{u}}, \\ \frac{\partial l}{\partial \mathbf{W}^{(1)}} &= \mathbf{x} \frac{\partial l}{\partial \mathbf{u}}, \\ \frac{\partial l}{\partial \mathbf{W}^{(2)}} &= \mathbf{h} \frac{\partial l}{\partial \mathbf{a}}.\end{aligned}$$

By the chain rule,  $\frac{\partial l}{\partial \mathbf{a}_i} = \frac{\partial l}{\partial \mathbf{z}_j} \cdot \frac{\partial \mathbf{z}_j}{\partial \mathbf{a}_i}$ . One can easily verify that

$$\frac{\partial l}{\partial \mathbf{z}_j} = -\delta_{yj} \frac{1}{\mathbf{z}_j},$$

$$\frac{\partial \mathbf{z}_j}{\partial \mathbf{a}_i} = \mathbf{z}_j (\delta_{ij} - \mathbf{z}_i),$$

where  $\delta_{ij} = 1[i == j]$ . Combine these two equations and remove the  $\delta_{ij}, \delta_{jy}$  (by setting  $j = y$ ),

$$\frac{\partial l}{\partial \mathbf{a}_i} = -\delta_{yj} \frac{1}{\mathbf{z}_j} \mathbf{z}_j (\delta_{ij} - \mathbf{z}_i) = -(\delta_{iy} - \mathbf{z}_i),$$

which implies that

$$\frac{\partial l}{\partial \mathbf{a}} = -(\mathbf{y} - \mathbf{z})^T.$$

**(In the denominator layout)**

$$\frac{\partial l}{\partial \mathbf{u}} = H(\mathbf{u}) \mathbf{W}^{(2)T} \frac{\partial l}{\partial \mathbf{a}}.$$

$$\begin{aligned}\frac{\partial l}{\partial \mathbf{b}^{(1)}} &= \frac{\partial l}{\partial \mathbf{u}}, \\ \frac{\partial l}{\partial \mathbf{W}^{(1)}} &= \frac{\partial l}{\partial \mathbf{u}} \mathbf{x}^T, \\ \frac{\partial l}{\partial \mathbf{W}^{(2)}} &= \frac{\partial l}{\partial \mathbf{a}} \mathbf{h}^T. \\ \frac{\partial l}{\partial \mathbf{a}} &= -(\mathbf{y} - \mathbf{z}).\end{aligned}$$

**Q1.2** Suppose we initialize  $W^{(1)}, W^{(2)}, b^{(1)}$  with zero matrices/vectors (i.e., matrices and vectors with all elements set to 0), please first verify that  $\frac{\partial l}{\partial W^{(1)}}, \frac{\partial l}{\partial W^{(2)}}, \frac{\partial l}{\partial b^{(1)}}$  are all zero matrices/vectors, irrespective of  $x, y$  and the initialization of  $b^{(2)}$ .

Now if we perform stochastic gradient descent for learning the neural network using a training set  $\{(x_i \in \mathbb{R}^D, y_i \in \mathbb{R}^K)\}_{i=1}^N$ , please explain with a concise mathematical statement in one sentence why no learning will happen on  $W^{(1)}, W^{(2)}, b^{(1)}$  (i.e., they will not change no matter how many iterations are run). Note that this will still be the case even with weight decay and momentum if the initial velocity vectors/-matrices are set to zero.

*What to submit:* No submission for the verification question. Your concise mathematical statement in one sentence for the explanation question.

Since  $W^{(2)}$  is all zero,  $\frac{\partial l}{\partial u}$  is all zero. So  $\frac{\partial l}{\partial W^{(1)}}, \frac{\partial l}{\partial b^{(1)}}$  are all zero. Since  $W^{(1)}, b^{(1)}$  are all zero,  $h$  is all zero. So  $\frac{\partial l}{\partial W^{(2)}}$  is all zero. In each iteration, all gradients are zero, so no updates will be made.

**Q1.3** As mentioned in the lecture, nonlinearity is very important for neural networks. With nonlinearity (e.g., eq. (3)), the neural network shown in Fig. 1 can be seen as a nonlinear basis function  $\phi$  (i.e.,  $\phi(x) = h$ ) followed by a linear classifier  $f$  (i.e.,  $f(h) = \hat{y}$ ).

Please show that, by removing the nonlinear operation in eq. (3) and setting eq. (4) to be  $a = W^{(2)}u + b^{(2)}$ , the resulting network is essentially a linear classifier. More specifically, you can now represent  $a$  as  $Ux + v$ , where  $U \in \mathbb{R}^{K \times D}$  and  $v \in \mathbb{R}^K$ . Please write down the representation of  $U$  and  $v$  using  $W^{(1)}, W^{(2)}, b^{(1)}$ , and  $b^{(2)}$ .

$$U = ?$$

$$v = ?$$

*What to submit:* No more than 2 lines of derivation for each of the question mark.

By combining the equations, we can get:

$$U = W^{(2)}W^{(1)},$$

$$v = W^{(2)}b^{(1)} + b^{(2)}.$$

## Problem 2 Logistic Regression

[Recommended maximum time spent: 45 minutes]

In the lectures we talked about logistic regression. Recall, that the decision boundary for logistic regression is found by maximizing log-likelihood:

$$\log(P(D)) = \sum_{n=1}^N y_n \log(\sigma(\mathbf{w}^T x_n)) + (1 - y_n) \log(1 - \sigma(\mathbf{w}^T x_n)) \quad (14)$$

Where  $x_n$  are the observations,  $y_n \in \{0, 1\}$  are their true labels. There is no a closed-form solution for logistic regression maximization, so we use iterative methods.

**Q2.1** Does logistic regression always yield a unique solution? If yes, explain why, if no, bring an example and explain why the solution will not be unique.

*What to submit:* Your answer, your less-than-5-lines explanation and an example, if applicable.

No. If the features are not linear-independent, for example  $x_{i1} = 2x_{i2}$ , for all training samples  $x_i$ , then if  $\beta_1 = \alpha_1, \beta_2 = \alpha_2$  is a valid solution,  $\beta'_1 = \alpha_1 + \frac{1}{2}\alpha_2, \beta'_2 = 0$  is a valid solution too.

**Q2.2** Consider a binary classification problem, where each observation  $x_n$  belongs to one of two classes:  $\{0, 1\}$ . However, suppose that the true label  $y_n$  of the observation  $x_n$  is not known, instead, we are given a value  $\pi_n$  that represents the probability that true label  $y_n = 1$ , i.e.  $p(y_n = 1|w) = \pi_n$ . Show that for such data  $D$  the log-likelihood function of logistic regression model with parameters  $w$  will be given by:

$$\log(P(D)) = \sum_{n=1}^N \pi_n \log(\sigma(w^T x_n)) + (1 - \pi_n) \log(1 - \sigma(w^T x_n)) \quad (15)$$

where  $\sigma$  is the sigmoid function:

$$\sigma(a) = \frac{1}{1 + e^{-a}} \quad (16)$$

*What to submit:* Your less-than-5-lines proof of Equation 15

By definition, we know

$$P(y_n = 1|x_n) = \pi_n,$$

$$P(\hat{y}_n = 1|x_n) = \sigma(w^T x_n).$$

So,

$$P(D|w) = \prod_{n=1}^N P(\hat{y}_n = 1|x_n, w)^{y_n} P(\hat{y}_n = 0|x_n, w)^{1-y_n}.$$

$$\log P(D) = \sum_{n=1}^N y_n \log P(\hat{y}_n = 1|x_n, w) + (1 - y_n) \log P(\hat{y}_n = 0|x_n, w)$$

Consider each data point  $(x_n, y_n)$  as a random variable, by taking the expectation over the log likelihood function, we can get

$$\begin{aligned} E[\log P(D)] &= \sum_{n=1}^N E[y_n] \log P(\hat{y}_n = 1|x_n, w) + E[(1 - y_n)] \log P(\hat{y}_n = 0|x_n, w) \\ &= \sum_{n=1}^N \pi_n \log P(\hat{y}_n = 1|x_n, w) + (1 - \pi_n) \log P(\hat{y}_n = 0|x_n, w) \\ &= \sum_{n=1}^N \pi_n \log \sigma(w^T x_n) + (1 - \pi_n) \log(1 - \sigma(w^T x_n)). \end{aligned}$$