

# HINT FILE: Assignment 7 (Fully Connected Neural Network)

## Artificial Intelligence (CSE-241N)

### IIT (BHU), Varanasi

April 14, 2018

## 1 Introduction

In this document we'll describe the equations for forward and backward passes (back-propagation) for a fully connected neural network. (To skip the explanation and see the algorithms themselves, go to Section 4.)

In this assignment Neural Networks will deal with supervised classification. So there are training examples consisting of input-output pairs  $(\mathbf{x}_{i_{f \times 1}}, y_i)$  where  $\mathbf{x}_i$  is the single feature vector corresponding to  $i^{th}$  training example and  $y_i$  is the class label for this feature. Further, for faster processing of data, these input output pairs are clubbed into batches, each consisting of multiple examples. A single batch looks like:

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1b} \\ x_{21} & x_{22} & \cdots & x_{2b} \\ \vdots & \vdots & \ddots & \vdots \\ x_{f1} & x_{f2} & \cdots & x_{fb} \end{pmatrix}_{f \times b} \quad (1)$$

$$\mathbf{Y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_b \end{pmatrix}_{b \times 1} \quad (2)$$

$\mathbf{X}$  can be viewed as a collection of  $b$  vectors:  $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_b)$ .

Neural Networks consist of layers, each of which is composed of a weight matrix and an activation function applied elementwise, which serves to provide non-linearity to the network.

We'll first derive the gradient equations for a three layer neural network. Then we give the equations for a general sized network at the end.

## 2 Forward Pass

Consider a 3 layer neural network, with three weight matrices:  $\mathbf{W}^{(i)}, i = 1, 2, 3$ . Let the input be as  $\mathbf{X}$  and output be  $\mathbf{Y}$  of batchsize  $b$ .

The forward pass for this network would look like:

$$\mathbf{h}^{(0)} = \mathbf{X}_{f \times b} \quad (3a)$$

$$\mathbf{a}_{n \times b}^{(1)} = \mathbf{W}_{n \times f}^{(1)} \mathbf{h}^{(0)} \quad (3b)$$

$$\mathbf{h}_{n \times b}^{(1)} = f(\mathbf{a}^{(1)}) \quad (3c)$$

$$\mathbf{a}_{m \times b}^{(2)} = \mathbf{W}_{m \times n}^{(2)} \mathbf{h}^{(1)} \quad (3d)$$

$$\mathbf{h}_{m \times b}^{(2)} = f(\mathbf{a}^{(2)}) \quad (3e)$$

$$\mathbf{a}_{l \times b}^{(3)} = \mathbf{W}_{l \times m}^{(3)} \mathbf{h}^{(2)} \quad (3f)$$

$$\hat{\mathbf{Y}}_{l \times b} = \mathbf{a}^{(3)} \quad (3g)$$

$$J_k = L(Y_k, \hat{\mathbf{Y}}_k) \quad (3h)$$

$$J = \sum_{i=1}^b J_i \quad (3i)$$

where, dimension  $l$  is the number of class labels in the dataset.  $\mathbf{h}^{(i)}$  is the output of the  $i^{th}$  hidden layer and  $\mathbf{W}^{(i)}$ s are layer's weight matrices.  $f$  is the activation function applied on the matrices elementwise.  $L$  is the loss function which gives a single scalar output and  $\hat{\mathbf{Y}}_k$  is the  $k^{th}$  vector of  $\hat{\mathbf{Y}}$ ;  $k$  corresponds to the  $k^{th}$  example number in the batch.

Note that in the above example, the neural network has:

- Input Layer Size:  $f$
- Hidden Layer 1 Size:  $n$
- Hidden Layer 2 Size:  $m$
- Output Layer Size:  $l$

### 3 Backward Pass

During the backward pass, our aim is to find out the gradients of the final scalar loss  $J$  with respect to each of the weight matrices and apply gradient descent to update them. For example, we need to evaluate

$$\nabla_{\mathbf{W}^{(1)}} J = \begin{pmatrix} \frac{\partial J}{\partial W_{11}^{(1)}} & \frac{\partial J}{\partial W_{12}^{(1)}} & \cdots & \frac{\partial J}{\partial W_{1f}^{(1)}} \\ \frac{\partial J}{\partial W_{21}^{(1)}} & \frac{\partial J}{\partial W_{22}^{(1)}} & \cdots & \frac{\partial J}{\partial W_{2f}^{(1)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial J}{\partial W_{n1}^{(1)}} & \frac{\partial J}{\partial W_{n2}^{(1)}} & \cdots & \frac{\partial J}{\partial W_{nf}^{(1)}} \end{pmatrix}$$

and then update  $\mathbf{W}^{(1)}$  using gradient descent:

$$\mathbf{W}^{(1)} = \mathbf{W}^{(1)} - \alpha \nabla_{\mathbf{W}^{(1)}} J$$

where,  $\alpha$  is the learning rate. We have to perform this same procedure for all other matrices as well.

#### 3.1 Backpropagation Algorithm

We'll now describe how to find out the gradient  $\nabla_{\mathbf{W}^{(i)}} J$  for the three layer neural network considered above using backpropagation algorithm (backprop). Backprop works on the principle of chain rule of calculus. Chain rule says that if you are given  $y = f(x)$  and  $z = g(f(x)) = g(y)$ , then  $\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$ , each of which can be independently calculated. Also, note that for any function  $g$  of variables  $x$  and  $y$ :

$$\nabla g(x, y) = \frac{\partial g}{\partial x} \nabla x + \frac{\partial g}{\partial y} \nabla y \quad (4)$$

Using these ideas we can calculate the value of  $\nabla_{\mathbf{W}^{(3)}} J$ .

### 3.1.1 Calculating $\nabla_{\mathbf{W}^{(3)}} J$

First note that

$$\begin{aligned}
\nabla_{\hat{\mathbf{Y}}} J &= \nabla_{\mathbf{a}^{(3)}} J \\
&= \begin{pmatrix} \frac{\partial J}{\partial a_{11}^{(3)}} & \frac{\partial J}{\partial a_{12}^{(3)}} & \cdots & \frac{\partial J}{\partial a_{1b}^{(3)}} \\ \frac{\partial J}{\partial a_{21}^{(3)}} & \frac{\partial J}{\partial a_{22}^{(3)}} & \cdots & \frac{\partial J}{\partial a_{2b}^{(3)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial J}{\partial a_{l1}^{(3)}} & \frac{\partial J}{\partial a_{l2}^{(3)}} & \cdots & \frac{\partial J}{\partial a_{lb}^{(3)}} \end{pmatrix} \\
&= \begin{pmatrix} \frac{\partial J_1}{\partial a_{11}^{(3)}} & \frac{\partial J_2}{\partial a_{12}^{(3)}} & \cdots & \frac{\partial J_b}{\partial a_{1b}^{(3)}} \\ \frac{\partial J_1}{\partial a_{21}^{(3)}} & \frac{\partial J_2}{\partial a_{22}^{(3)}} & \cdots & \frac{\partial J_b}{\partial a_{2b}^{(3)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial J_1}{\partial a_{l1}^{(3)}} & \frac{\partial J_2}{\partial a_{l2}^{(3)}} & \cdots & \frac{\partial J_b}{\partial a_{lb}^{(3)}} \end{pmatrix} \tag{5}
\end{aligned}$$

where the second equality follows from Eq. 3i. Now using the expression for  $\nabla_{\mathbf{a}^{(3)}} J$  and Eq. 3f, we first find  $\nabla_{\mathbf{W}^{(3)}} \mathbf{a}_{ij}^{(3)}$  and then calculate  $\nabla_{\mathbf{W}^{(3)}} J$ . From Eq. 3f,  $\mathbf{a}_{ij}^{(3)} = \sum_{k=1}^m \mathbf{W}_{ik}^{(3)} \mathbf{h}_{kj}^{(3)}$ , so

$$\nabla_{\mathbf{W}^{(3)}} \mathbf{a}_{ij}^{(3)} = \begin{pmatrix} 0 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ h_{1j}^{(3)} & h_{2j}^{(3)} & \cdots & h_{mj}^{(3)} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \end{pmatrix} i^{th} \text{row}, \tag{6}$$

which is a matrix with a single non-zero row. **(Can you reason out, from Eq. 3f, why  $\nabla_{\mathbf{W}^{(3)}} \mathbf{a}_{ij}^{(3)}$  has a single non-zero row?)** From Eq. 3h,  $J_k$  can be seen as a function of  $\mathbf{a}_{ik}^{(3)}$  s. Combining this fact with Eq. 4 we get:

$$\begin{aligned}
\nabla_{\mathbf{W}^{(3)}} J_k &= \sum_{i=1}^l \left( \frac{\partial J_k}{\partial a_{ik}^{(3)}} \right) \nabla_{\mathbf{W}^{(3)}} \mathbf{a}_{ik}^{(3)} \\
&= \begin{pmatrix} \frac{\partial J_k}{\partial a_{1k}^{(3)}} h_{1k}^{(3)} & \frac{\partial J_k}{\partial a_{1k}^{(3)}} h_{2k}^{(3)} & \cdots & \frac{\partial J_k}{\partial a_{1k}^{(3)}} h_{mk}^{(3)} \\ \frac{\partial J_k}{\partial a_{2k}^{(3)}} h_{1k}^{(3)} & \frac{\partial J_k}{\partial a_{2k}^{(3)}} h_{2k}^{(3)} & \cdots & \frac{\partial J_k}{\partial a_{2k}^{(3)}} h_{mk}^{(3)} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial J_k}{\partial a_{lk}^{(3)}} h_{1k}^{(3)} & \frac{\partial J_k}{\partial a_{lk}^{(3)}} h_{2k}^{(3)} & \cdots & \frac{\partial J_k}{\partial a_{lk}^{(3)}} h_{mk}^{(3)} \end{pmatrix} \\
&= \begin{pmatrix} \frac{\partial J_k}{\partial a_{1k}^{(3)}} \\ \frac{\partial J_k}{\partial a_{2k}^{(3)}} \\ \vdots \\ \frac{\partial J_k}{\partial a_{lk}^{(3)}} \end{pmatrix} \begin{pmatrix} h_{1k}^{(3)} & h_{2k}^{(3)} & \cdots & h_{mk}^{(3)} \end{pmatrix} \\
&= (\nabla_{\mathbf{a}^{(3)}} J)_k \mathbf{h}_k^{(3)\top} \tag{7}
\end{aligned}$$

where  $(\nabla_{\mathbf{a}^{(3)}} J)_k \mathbf{h}_k^{(3)\top}$  is the outer product of  $(\nabla_{\mathbf{a}^{(3)}} J)_k$ , the  $k^{th}$  column of  $\nabla_{\mathbf{a}^{(3)}} J$  and  $\mathbf{h}_k^{(3)\top}$ , the  $k^{th}$  column of  $\mathbf{h}^{(3)}$  transposed. Again, using Eq. 3i, we get:

$$\begin{aligned}
\nabla_{\mathbf{W}^{(3)}} J &= \sum_{k=1}^b \nabla_{\mathbf{W}^{(3)}} J_k \\
&= \sum_{k=1}^b (\nabla_{\mathbf{a}^{(3)}} J)_k \mathbf{h}_k^{(3)\top} \\
&= (\nabla_{\mathbf{a}^{(3)}} J) \mathbf{h}^{(3)\top}
\end{aligned} \tag{8}$$

where the last equality follows from viewing product of the two matrices as the sum of the outer products of their individual columns and rows. Eq. 8 gives the final expression for  $\nabla_{\mathbf{W}^{(3)}} J$  and we are done.

Our next task is to calculate  $\nabla_{\mathbf{W}^{(2)}} J$ , for which we must obtain the expression for  $\nabla_{\mathbf{a}^{(2)}} J$ , and then using Eq. 3d we can proceed as we did in this section. Calculating  $\nabla_{\mathbf{a}^{(2)}} J$  itself requires the knowledge of  $\nabla_{\mathbf{h}^{(2)}} J$ .

### 3.1.2 Calculating $\nabla_{\mathbf{h}^{(2)}} J$

Now we move on to calculate the gradient for  $\mathbf{h}^{(2)}$ . Again, from Eq. 3f:  $a_{ij}^{(3)} = \sum_{k=1}^m \mathbf{W}_{ik}^{(3)} \mathbf{h}_{kj}^{(3)}$ , so

$$\nabla_{\mathbf{h}^{(2)}} a_{ik}^{(3)} = \begin{matrix} k^{th} \text{ column} \\ \begin{pmatrix} 0 & 0 & \cdots & W_{i1}^{(3)} & \cdots & 0 \\ 0 & 0 & \cdots & W_{i2}^{(3)} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & W_{im}^{(3)} & \cdots & 0 \end{pmatrix} \end{matrix} \tag{9}$$

which is a matrix with a single non-zero column. **(Can you see why  $\nabla_{\mathbf{h}^{(2)}} a_{ik}^{(3)}$  has a single non-zero column? Also how is it different from  $\nabla_{\mathbf{W}^{(3)}} a_{ij}^{(3)}$  which had a single non-zero row?)** From Eq. 3h and Eq. 3f,  $J_k$  can be seen as a function of  $\mathbf{h}_{ik}^{(2)}$ s. Combining this fact with Eq. 4 we get:

$$\begin{aligned}
\nabla_{\mathbf{h}^{(2)}} J_k &= \sum_{i=1}^l \left( \frac{\partial J_k}{\partial a_{ik}^{(3)}} \right) \nabla_{\mathbf{h}^{(2)}} a_{ik}^{(3)} \\
&= \begin{pmatrix} 0 & 0 & \cdots & \sum_{i=1}^l W_{i1}^{(3)} \frac{\partial J_k}{\partial a_{ik}^{(3)}} & \cdots & 0 \\ 0 & 0 & \cdots & \sum_{i=1}^l W_{i2}^{(3)} \frac{\partial J_k}{\partial a_{ik}^{(3)}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sum_{i=1}^l W_{im}^{(3)} \frac{\partial J_k}{\partial a_{ik}^{(3)}} & \cdots & 0 \end{pmatrix}
\end{aligned} \tag{10}$$

Again, from Eq. 3i we have  $\nabla_{\mathbf{h}^{(2)}} J = \sum_{k=1}^b \nabla_{\mathbf{h}^{(2)}} J_k$ . Using Eq. 10 with this expression we can write  $(\nabla_{\mathbf{h}^{(2)}} J)_{ik} = \sum_{j=1}^l W_{ij}^{(3)} (\nabla_{\mathbf{a}^{(3)}} J)_{jk}$ . This is equivalent to:

$$\nabla_{\mathbf{h}^{(2)}} J = \mathbf{W}^{(3)\top} (\nabla_{\mathbf{a}^{(3)}} J) \tag{11}$$

### 3.1.3 Calculating $\nabla_{\mathbf{a}^{(2)}} J$

We can now calculate  $\nabla_{\mathbf{a}^{(2)}} J$ . From Eq. 3e,  $\mathbf{h}_{ij}^{(2)} = f(a_{ij}^{(2)})$ . So,

$$\begin{aligned}
\frac{\partial J}{\partial a_{ij}^{(2)}} &= \frac{\partial J}{\partial h_{ij}^{(2)}} \frac{\partial h_{ij}^{(2)}}{\partial a_{ij}^{(2)}} \\
&= \frac{\partial J}{\partial h_{ij}^{(2)}} f'(a_{ij}^{(2)})
\end{aligned}$$

Writing the above equation in the matrix form:

$$\nabla_{\mathbf{a}^{(2)}} J = \nabla_{\mathbf{h}^{(2)}} J \odot f'(\mathbf{a}^{(2)}) \quad (12)$$

where  $\odot$  represents elementwise multiplication.

#### 3.1.4 Calculating $\nabla_{\mathbf{W}^{(2)}} J$ and $\nabla_{\mathbf{W}^{(1)}} J$

Now that we have the expression for  $\nabla_{\mathbf{a}^{(2)}} J$ , we can use Eq. 3d to calculate  $\nabla_{\mathbf{W}^{(2)}} J$ , using a similar procedure as adopted by us in Section 3.1.1. We can then use Eq. 3d and 3c to calculate  $\nabla_{\mathbf{a}^{(1)}} J$ , as done in Section 3.1.3. We keep on repeating this procedure till we have finally evaluated  $\nabla_{\mathbf{W}^{(1)}} J$ , at which point we are done.

## 4 Algorithmic form for forward pass and backward pass

In this section we extend the ideas discussed in the previous sections to an  $l$  layered fully connected neural network. The algorithms below are a modified form of those given in [1].

### 4.1 Forward Pass

---

**Algorithm 1** Forward pass of a fully connected neural network

---

- 1: **Input:** Network depth,  $l$
  - 2: **Input:**  $\mathbf{W}^{(i)}, i \in 1, \dots, l$ , the weight matrices
  - 3: **Input:**  $\mathbf{b}^{(i)}, i \in 1, \dots, l$ , the bias terms
  - 4: **Input:**  $\mathbf{X}_{f \times b}$ , the input to the model, where  $f$  is the number of features in a single sample and  $b$  is the batchsize
  - 5: **Input:**  $\mathbf{Y}_{b \times 1}$ , the output class labels
  - 6:  $\mathbf{h}^{(0)} = \mathbf{X}$
  - 7: **for**  $k = 1, \dots, l - 1$  **do**
  - 8:      $\mathbf{a}^{(k)} = \mathbf{W}^{(k)} \mathbf{h}^{(k-1)} + \mathbf{b}^{(k)}$
  - 9:      $\mathbf{h}^{(k)} = f(\mathbf{a}^{(k)})$
  - 10:  $\mathbf{a}^{(l)} = \mathbf{W}^{(l)} \mathbf{h}^{(l-1)} + \mathbf{b}^{(l)}$
  - 11:  $\hat{\mathbf{Y}} = \mathbf{a}^{(l)}$
  - 12:  $J_k = L(\mathbf{Y}_k, \hat{\mathbf{Y}}_k)$ , for each  $k$  in  $1, \dots, b$
  - 13:  $J = \sum_{k=1}^b J_k$
- 

### 4.2 Backward Pass

---

**Algorithm 2** Backward pass of a fully connected neural network

---

- 1:  $\mathbf{g} \leftarrow \nabla_{\hat{\mathbf{Y}}} J$
  - 2:  $\mathbf{O} = [1 \ 1 \ \dots \ 1]^\top$ , the number of 1s in  $\mathbf{O}$  is equal to the number of columns in  $\mathbf{g}$
  - 3:  $\nabla_{\mathbf{b}^{(l)}} J = \mathbf{g} \mathbf{O}$
  - 4:  $\nabla_{\mathbf{W}^{(l)}} J = \mathbf{g} \mathbf{h}^{(l-1)\top}$
  - 5:  $\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(l-1)}} J = \mathbf{W}^{(l)\top} \mathbf{g}$
  - 6: **for**  $k = l - 1, \dots, 1$  **do**
  - 7:      $\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)})$
  - 8:      $\mathbf{O} = [1 \ 1 \ \dots \ 1]^\top$ , the number of 1s in  $\mathbf{O}$  is equal to the number of columns in  $\mathbf{g}$
  - 9:      $\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} \mathbf{O}$
  - 10:      $\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)\top}$
  - 11:      $\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)\top} \mathbf{g}$
-

## 5 Loss function and the activation function

In this section we would talk about the form of loss function and the activation functions we'll use in the assignment and the procedure for calculating their respective gradients.

The loss function we will work with in the assignment is the **Cross-Entropy Loss**, defined by:

$$J_k = L(Y_k, \hat{Y}_k) = -\log \left( \frac{e^{\hat{Y}_{Y_k, k}}}{\sum_{i=1}^l e^{\hat{Y}_{i, k}}} \right) \quad (13)$$

where  $J_k$  is the loss for the  $k^{th}$  sample,  $l$  is the number of classes and  $Y_k$  is the correct label (a single scalar number) corresponding to the class of the  $k^{th}$  sample. Derivative of the loss function with respect to  $\hat{Y}_{ik}$  is given by:

$$\frac{\partial J_k}{\partial \hat{Y}_{ik}} = \begin{cases} \frac{e^{\hat{Y}_{ik}}}{\sum_{i=1}^l e^{\hat{Y}_{ik}}} - 1 & , i = Y_k \\ \frac{e^{\hat{Y}_{ik}}}{\sum_{i=1}^l e^{\hat{Y}_{ik}}} & , i \neq Y_k \end{cases} \quad (14)$$

For further information, read the section “Softmax classifier” and “Practical issues: Numeric stability” on <http://cs231n.github.io/linear-classify/#softmax>.

One of the popular activation functions is  $\tanh()$ . Again, although you are free to use any activation function of your choice in the assignment, we'll describe just  $\tanh()$  here.

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (15)$$

Its derivative is

$$\frac{\partial \tanh(x)}{\partial x} = 1 - \tanh^2(x) \quad (16)$$

Remember to apply the activation functions elementwise, and also use `numpy.tanh()` function instead of writing your own implementations.

## References

- [1] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.