

National University of Singapore  
School of Computing  
CS2109S: Introduction to AI and Machine Learning  
Semester I, 2023/2024

## Tutorial 8 Back Propagation

These questions will be discussed during the tutorial session. Please be prepared to answer them.

### Summary of Key Concepts

In this tutorial, we will discuss and explore the following key learning points/lessons from lecture:

1. Back propagation
2. Matrix calculus for back propagation
3. Potential issues with back propagation

**Warning:** Lots of math ahead! For "Show that" questions, start from the *left hand side* of the equation

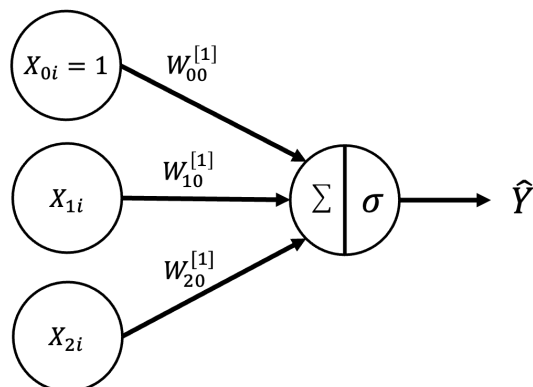
### Problems

#### 1. Back Propagation (Warm Up)

Grace has a wine dataset which comprises of 1100 samples. Each wine sample has a label indicating which cultivar it comes from, and two features: colour intensity and alcohol level.

Now, she wants to build a classifier that can predict which cultivar a wine sample comes from using the two features.

- (a) She decided to train a neural network with the architecture shown in the figure below.



Mathematically, this is given by

$$f^{[1]} = W^{[1]T} X$$

$$\hat{Y} = g^{[1]}(f^{[1]})$$

where  $W^{[1]} = [W_{00}^{[1]} \ W_{10}^{[1]} \ W_{20}^{[1]}]^T$ ,  $X \in \mathbf{R}^{3 \times n}$ , and  $g^{[1]}(s) = \sigma(s) = \frac{1}{1+e^{-s}}$ . Note that here,  $X_{0i} = 1$ , and  $X_{1i}$  and  $X_{2i}$  represent the colour intensity and alcohol level respectively.

She decided to use the following loss function<sup>1</sup>

$$\mathcal{E} = -\frac{1}{n} \sum_{i=0}^{n-1} \left\{ [Y_{0i} \cdot \log(\hat{Y}_{0i})] + [(1 - Y_{0i}) \log(1 - \hat{Y}_{0i})] \right\}$$

where  $\hat{Y} \in \mathbf{R}^{1 \times n}$  and  $Y \in \{0, 1\}^{1 \times n}$  such that  $Y_{0i} = 1$  if the  $i$ th wine sample is from cultivar A and  $Y_{0i} = 0$  if it is from cultivar B, and  $n$  is the number of wine samples (i.e.  $n = 1100$  in this case).

To illustrate, we calculate the loss function for the following sample of 2:

$$Y = [0 \ 1]$$

$$\hat{Y} = [0.2 \ 0.9]$$

Loss function calculation:

$$\begin{aligned} \mathcal{E} &= -\frac{1}{2} \left[ (1 - 0) \cdot \log(1 - 0.2) + 1 \cdot \log(0.9) \right] \\ &= -\frac{1}{2} \left[ \log(0.8) + \log(0.9) \right] = 0.16425 \end{aligned}$$

Note that  $\log(1) = 0$  and  $\log(0) \rightarrow -\infty$ . We assume that  $\log(x) = \ln(x)$  in this tutorial.

For parts (a) and (b), to keep things simple, let us consider the case where  $n = 1$ . **Show that:**

$$\text{i. } \frac{\partial \mathcal{E}}{\partial \hat{Y}} = \left[ -\frac{Y_{00}}{\hat{Y}_{00}} + \frac{1 - Y_{00}}{1 - \hat{Y}_{00}} \right]$$

We provide the answer for this part as a guiding example:

Note that  $\frac{\partial \mathcal{E}}{\partial \hat{Y}} = \left[ \frac{\partial \mathcal{E}}{\partial \hat{Y}_{00}} \right]$  when  $n = 1$

$$\frac{\partial \mathcal{E}}{\partial \hat{Y}_{00}} = -\frac{Y_{00}}{\hat{Y}_{00}} + \frac{1 - Y_{00}}{1 - \hat{Y}_{00}}$$

Therefore,

$$\frac{\partial \mathcal{E}}{\partial \hat{Y}} = \left[ -\frac{Y_{00}}{\hat{Y}_{00}} + \frac{1 - Y_{00}}{1 - \hat{Y}_{00}} \right]$$

Try to plug in the case where (1) the true label  $Y_{00} = 0$  and (2)  $Y_{00} = 1$ . This equation tells us how the change in predicted value  $\hat{Y}$  will be able to influence the change of error/loss value.

---

<sup>1</sup>This loss function is usually known as *log loss* or *binary cross-entropy*.

ii.  $\frac{\partial \mathcal{E}}{\partial f^{[1]}} = \hat{Y} - Y$

Hint: Since  $n = 1$ ,  $\frac{\partial \mathcal{E}}{\partial f^{[1]}} = \left[ \frac{\partial \mathcal{E}}{\partial f_{00}^{[1]}} \right]$ . By chain rule,  $\frac{\partial \mathcal{E}}{\partial f_{00}^{[1]}} = \frac{\partial \mathcal{E}}{\partial \hat{Y}_{00}} \frac{\partial \hat{Y}_{00}}{\partial f_{00}^{[1]}} \dots$

iii.  $\frac{\partial \mathcal{E}}{\partial W_{20}^{[1]}} = \left( \frac{\partial \mathcal{E}}{\partial f^{[1]}} \right)_{00} X_{20}$

Hint: Since  $n = 1$ ,  $f^{[1]} = \left[ f_{00}^{[1]} \right]$ . Thus,  $\frac{\partial f^{[1]}}{\partial W_{20}^{[1]}} = \left[ \frac{\partial f_{00}^{[1]}}{\partial W_{20}^{[1]}} \right] \dots$

**NOTE:**  $\left( \frac{\partial \mathcal{E}}{\partial f^{[1]}} \right)_{00}$  refers to the (0, 0) entry of the matrix  $\frac{\partial \mathcal{E}}{\partial f^{[1]}}$ .

You should use the chain rule to derive your answers. Treat  $f^{[1]}$ ,  $W^{[1]}$ ,  $\hat{Y}$ ,  $Y$  and  $X$  as matrices.

- (b) Using your answer in (a), derive an expression for  $\frac{\partial \mathcal{E}}{\partial W^{[1]}}$ . Recall that Back propagation is the act of nudging/fine-tuning the weights  $W$  during training to minimize the loss  $\mathcal{E}$ . Does your derived expression prove that Back propagation works?
- (c) Let us consider a general case where  $n \in \mathbb{N}$ . Using your answer to (a), find  $\frac{\partial \mathcal{E}}{\partial f^{[1]}}$  when  $n \in \mathbb{N}$ . Check your answer using the Python notebook.
- Hint: Read the Python notebook.
- (d) Let's say that Grace has 100 samples from cultivar A and 1000 samples from cultivar B that make up the 1100 total samples. To deal with the imbalanced data set, she decided to introduce two hyper-parameters  $\alpha$  and  $\beta$  in the loss function, as shown below.

$$\mathcal{E} = -\frac{1}{n} \sum_{i=0}^{n-1} \left\{ \alpha [Y_{0i} \cdot \log(\hat{Y}_{0i})] + \beta [(1 - Y_{0i}) \cdot \log(1 - \hat{Y}_{0i})] \right\}$$

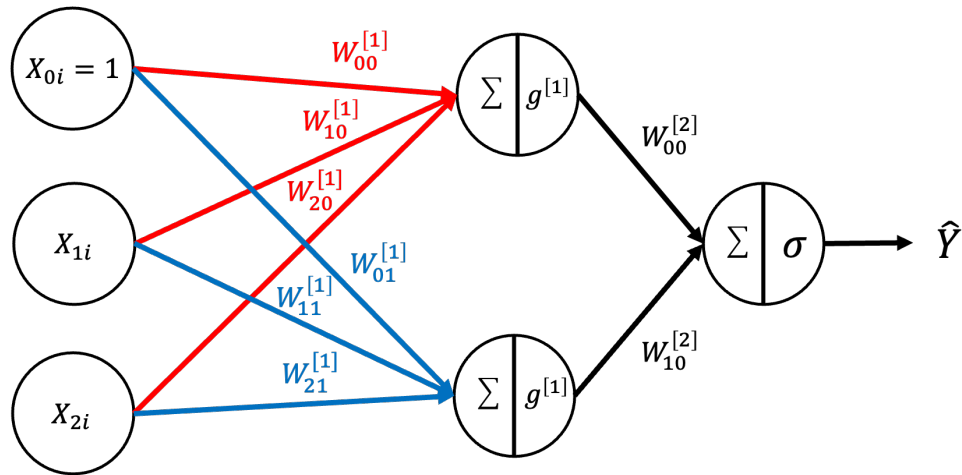
For example, using the same sample of 2 from (a), the loss function can be computed as follows:

$$\begin{aligned} \mathcal{E} &= -\frac{1}{2} \left\{ \beta [(1 - 0) \cdot \log(1 - 0.2)] + \alpha [1 \cdot \log(0.9)] \right\} \\ &= -\frac{1}{2} \left\{ \beta \cdot \log(0.8) + \alpha \cdot \log(0.9) \right\} \end{aligned}$$

Why do you think that she introduced the hyper-parameters  $\alpha$  and  $\beta$ ? How should she set their values?

## 2. Back Propagation for a Deep(er) Network

After training the neural network described in question 1, Grace observed that the training error is high. Thus, she decided to introduce a hidden layer, resulting in the architecture shown in the figure below.



Mathematically, the network is given by

$$f^{[1]} = W^{[1]T} X$$

$$a^{[1]} = g^{[1]}(f^{[1]})$$

$$f^{[2]} = W^{[2]T} a^{[1]}$$

$$\hat{Y} = g^{[2]}(f^{[2]})$$

where  $g^{[1]}(s) = \text{ReLU}(s)$ ,  $g^{[2]}(s) = \sigma(s) = \frac{1}{1+e^{-s}}$ ,  $W^{[1]} \in \mathbb{R}^{3 \times 2}$  and  $W^{[2]} \in \mathbb{R}^{2 \times 1}$ .

The *ReLU* function is defined as follows:

$$\text{ReLU}(s) = \begin{cases} s & \text{if } s > 0 \\ 0 & \text{otherwise} \end{cases}$$

Similar to question 1, let us keep things simple and consider what happens when  $n = 1$ . In addition, assume that the loss function used remains the same.

**Compute**  $\frac{\partial \mathcal{E}}{\partial W_{11}^{[1]}}$ .

Recall again that in Back propagation, we can decrease or increase the loss by changing the weights. But how strong the effect of weight nudging on the loss also depends on the neuron  $X$  associated with the weight. In this question, we introduce hidden layer that has its own neurons activation. Based on your derived expression, how can this hidden layer amplify the effect of weight nudging that helps minimize the training error?

Your answer should not contain any partial derivatives on the right hand side, they should all be evaluated.

*Hint: The results found/observations made in question 1 are likely to be useful here.*

### 3. Potential Issues with Training Deep Neural Networks

Suppose we use  $\sigma$  (the sigmoid function) as our activation function in a neural network with 50 hidden layers, as per the code in the accompanying Python notebook.

- (a) Play around with the code. Notice that when performing back propagation, the gradient magnitudes of the first few layers are extremely small. What do you think causes this problem?
- (b) Based on what we have learnt thus far, how can we **mitigate** this problem? Test out your solution by modifying the code and checking the gradient magnitudes.

*Hint: You don't need to change much.*