

# CS 5480: Deep Learning

## Homework 0

Ques 1

$$\text{Payout} = \begin{cases} +Rs 100, & n=1 \\ -Rs 25, & n \neq 1 \end{cases}$$

where (+) means payout, to ELAN and (-) means payout to a visitor.

Expectation of gaining money (From ELAN's perspective)  $\Rightarrow$

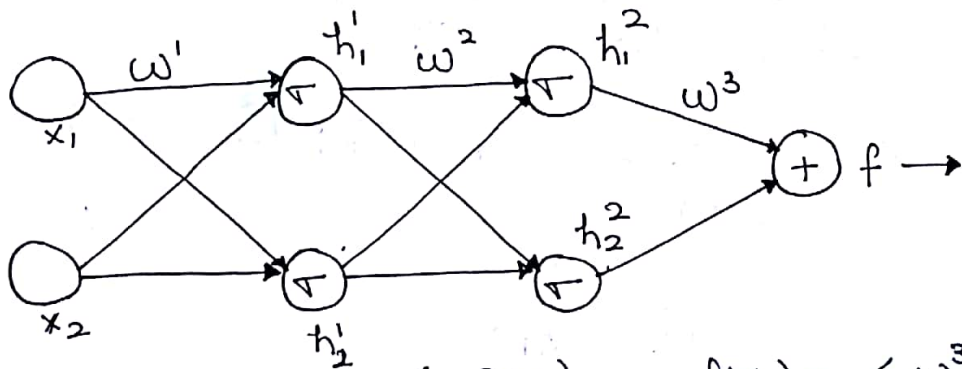
$$E(X) = \left( 100 \times \frac{1}{6} - 25 \times \frac{5}{6} \right) Rs$$

$$E(X) = \frac{-25}{6} = -4.1667 Rs$$

$\rightarrow$  This means ELAN <sup>is expected to</sup> ~~would be~~ giving money to visitors in most cases.

So it's not a good plan.

Ques 2



$$h^1 = \sigma(w^1 x), \quad h^2 = \sigma(w^2 h^1), \quad f(x) = \langle w^3, h^2 \rangle$$

- $\rightarrow h_k^2$  denotes output of  $k^{th}$  node in layer 3
- $\rightarrow h_i^1$  denotes output of  $i^{th}$  node in layer 2
- $\rightarrow x_j$  denotes  $j^{th}$  input in layer 1
- $\rightarrow w_{ij}^1$  denotes weight of  $j^{th}$  input in  $i^{th}$  node in layer 2.

To Compute:

$$\frac{\partial f}{\partial \omega'_{ij}} = \frac{\partial f}{\partial h^2_k} \frac{\partial h^2_k}{\partial h'_i} \frac{\partial h'_i}{\partial \omega'_{ij}} \quad (\text{for all } k)$$

$$\rightarrow \frac{\partial f}{\partial h^2_k} = \omega_k^3 \Rightarrow \begin{matrix} \omega_1^3 \\ \omega_2^3 \end{matrix}$$

$$\rightarrow \frac{\partial h^2_k}{\partial h'_i} \Rightarrow h^2_k = \sigma(\omega_{k1}^2 h'_1 + \omega_{k2}^2 h'_2) = \sigma\left(\sum_{i=1}^2 \omega_{ki}^2 h'_i\right)$$

$$\frac{\partial h^2_k}{\partial h'_i} = (\omega_{ki}^2) h^2_k (1 - h^2_k)$$

$$\rightarrow \frac{\partial h'_i}{\partial \omega'_{ij}} \Rightarrow h'_i = \sigma(\omega_{i1} x_1 + \omega_{i2} x_2) = \sigma\left(\sum_{j=1}^2 \omega_{ij} x_j\right)$$

$$\frac{\partial h'_i}{\partial \omega'_{ij}} = h'_i (1 - h'_i) \cdot x_j$$

$$\frac{\partial f}{\partial \omega'_{ij}} = \frac{\partial f}{\partial h^2_k} \cdot \frac{\partial h^2_k}{\partial h'_i} \cdot \frac{\partial h'_i}{\partial \omega'_{ij}} \quad (\text{for all } k)$$

$$\Rightarrow \frac{\partial f}{\partial \omega'_{ij}} = \omega_k^3 \cdot \omega_{ki}^2 \cdot h^2_k (1 - h^2_k) \cdot h'_i (1 - h'_i) \cdot x_j \quad (\text{for all } k)$$

$$\approx \omega_1^3 \omega_{1i}^2$$

$$\approx \omega_k^3 = \omega_1^3 \cdot \omega_{1,i}^2 \cdot h_1^2 (1 - h_1^2) h'_i (1 - h'_i) \cdot x_j$$

$$+ \omega_2^3 \omega_{2,i}^2 h_2^2 (1 - h_2^2) h'_i (1 - h'_i) \cdot x_j$$

$$\boxed{\frac{\partial f}{\partial \omega'_{ij}} = h'_i (1 - h'_i) x_j \left[ \omega_1^3 \cdot \omega_{1,i}^2 \cdot h_1^2 (1 - h_1^2) + \omega_2^3 \cdot \omega_{2,i}^2 \cdot h_2^2 (1 - h_2^2) \right]}$$

Ques3

$$\Delta_{ij}^{(2)} := \Delta_{ij}^2 + \delta_i^{(3)} * (a_j^{(2)})$$

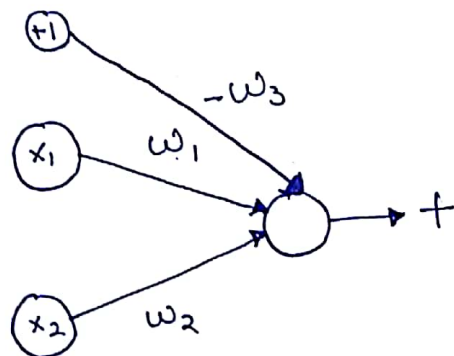
As we can clearly see, while vectorizing this equation we need to add two matrices of same size.

And  $\delta \times a$  should be computed such that  $\delta_i \times a_j$  element is added to  $\Delta_{ij}^2$ .

So we can simply rewrite the above equation as

$$\Delta^2 := \Delta^2 + \delta^{(3)} * [a^{(2)}]^T$$

Ques4



we can see that this perceptron will simply learn a straight line  $w_1 x_1 + w_2 x_2 - w_3 = 0 \Rightarrow w_1 x_1 + w_2 x_2 = w_3$

Now if  $w_3 \neq 0$  we can always make another perceptron of type  $(\frac{w_1}{w_3})x_1 + (\frac{w_2}{w_3})x_2 = 1$ , where fixed bias would be 1.

But in our case  $w_3 = 0$ , so we cannot create another perceptron that will learn same function.

Another approach to look at this would be,

Our first Perceptron will give

$$w_1 x_1 + w_2 x_2 = y$$

Second Perceptron will give  $w_1 x_1 + w_2 x_2 + 1 = y$

For these two outputs to be same

we can  $w_1 x_1 + w_2 x_2 = w'_1 x_1 + w'_2 x_2 + 1$

$$\Rightarrow x_1(w_1 - w'_1) + x_2(w_2 - w'_2) = 1$$

Here we are claiming that our perceptron learn some weights  $(v_1, v_2)$  s.t.  $v_1 x_1 + v_2 x_2 = 1$  for every  $(x_1, x_2)$ , which is clearly not possible

→ So, we cannot create another Perceptron with fixed bias 1 which will learn same as function as a Perceptron with no bias.



Ques 5

$$f(x) = \log \left( \sum_{i=1}^n e^{x_i} \right)$$

$$\begin{aligned} (a) \quad \frac{\partial f}{\partial x_i} &= \frac{1}{\sum_{j=1}^n e^{x_j}} \cdot \frac{\partial \left( \sum_{j=1}^n e^{x_j} \right)}{\partial x_i} \\ &= \frac{1}{\sum_{j=1}^n e^{x_j}} \cdot e^{x_i} = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} \end{aligned}$$

(b) Let's ~~not~~ assume  $s = \sum_{j=1}^n e^{x_j}$

$$\begin{aligned} \text{So, } \nabla f(x) &= \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right) \in \mathbb{R}^n \\ \Rightarrow \nabla f(x) &= \left( \frac{e^{x_1}}{s}, \frac{e^{x_2}}{s}, \dots, \frac{e^{x_n}}{s} \right) \in \mathbb{R}^n \end{aligned}$$

(c) Yes There exist a chain Matrix ~~calculus~~ Calculus equivalent of chain rule from Calculus.

Consider a matrix  $g(x)$  which is acted upon by a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ .  
 $h(x) = f(g(x))$ . Now to calculate  $\frac{dh(x)}{dx}$

we can use the chain rule

$$\frac{dh(x)}{dx} = f'(g(x)) g'(x) \cdot [g'(x) f'(g(x))]$$

The only condition imposed in matrix calculus is that the order should remain same.

## Ques 6

Let's first write Binary Cross entropy function which we will use as the error function in Neural Network.

Here  $y_i$  represent true labels and  $a_i$  represent predictions.

$$BCE(y) = - \sum_{i=1}^n y_i \log a_i + (1 - y_i) \log (1 - a_i)$$

In pure optimization we can just minimize above ~~loss~~ function.  
So, if we use neural network to minimize the BCE loss we will get an optimal solution.

From Probability point of view

$$P(y|\pi) = \prod_{i=1}^n \pi_i^{y_i} (1 - \pi_i)^{1 - y_i}$$

(Bernoulli for Binary Classification)

For given inputs  $x_i$  and given output  $y_i$

$$P(y|x, \theta) = \prod_{i=1}^n p_{\theta}(y|x_i)^{y_i} (1 - p_{\theta}(y|x_i))^{1 - y_i}$$

Since maximizing above function is hard and log is a monotone function, we can instead maximize log-likelihood

$$L(y|x, \theta) = \sum_{i=1}^n y_i \log p_{\theta}(y|x_i) + (1 - y_i) \log (1 - p_{\theta}(y|x_i))$$

Now as we can see maximising log likelihood and minimizing BCE gives same result.

For Multiclass classification one can use softmax cross-entropy as loss function and compare it with Multinoulli likelihood.

### Ques 7

- (i) ~~Based on results obtained from Perceptron and logistic regression are same. Both of them give probabilities of correct output.~~
- (ii) Both the methods ~~are~~ give probabilities of the correct output.
- (i) Solutions obtained by Logistic regression and Perceptron with sigmoid activation function are ~~same~~ more or less the same.
- Perceptron follows online training and Logistic Regression follows offline training (closed form).

### Ques 8

Ques 1 ~ 5 minutes

Ques 2 ~ 15 minutes

Ques 3 ~ 1 minute

Ques 4 ~ 10 minutes

Ques 5 ~ 10 minutes

Ques 6 ~ 20 minutes

Ques 7 ~ 5 minutes

Ques 8 ~ 2 minutes