

# Deep Learning (83882) - Ex 1

Due: 07.12.2024, 11:59pm

## Q1

Consider a multivariate logistic regression problem. Recall that given a vector  $\mathbf{x} \in \mathbb{R}^d$ , and a parameter matrix  $W \in \mathbb{R}^{d \times c}$ , where  $c$  is the number of classes, the model can be written in the following manner using the Softmax function:

$$\text{softmax}(\mathbf{z})_{[i]} = \frac{\exp(\mathbf{z}_i)}{\sum_j \exp(\mathbf{z}_j)},$$

where,  $\mathbf{z} = W^T \mathbf{x}$ .

1. Show that  $\text{softmax}(\mathbf{z}) = \text{softmax}(\mathbf{z} + m \cdot \mathbf{1})$  for every scalar  $m$ , where  $\mathbf{1}$  is a vector of ones of appropriate size.
2. For  $c = 2$ , show that the Sigmoid function is equivalent to the Softmax function.
3. Present an alternative to the Sigmoid function which also maps from the real line to the  $[0, 1]$  interval (any valid solution will be acceptable here).

## Q2

Let  $\mathbf{x} \in \{0, 1\}^2$  be an input vector. Consider the following model (scalar function):

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{w}^T \mathbf{h} + b_2 \\ \mathbf{h} &= \max(U^T \mathbf{x} + \mathbf{b}_1, 0) \end{aligned}$$

Where  $U \in \mathbb{R}^{2 \times 2}$ ,  $\mathbf{b}_1 \in \mathbb{R}^2$ ,  $\mathbf{w} \in \mathbb{R}^2$ ,  $b_2 \in \mathbb{R}$ , and the  $\max$  is taken element-wise.

Suppose we would like to represent with  $f(\mathbf{x})$  the XOR function, defined as:

$$\begin{aligned} \text{XOR}(0, 0) &= 0, \\ \text{XOR}(0, 1) &= 1, \\ \text{XOR}(1, 0) &= 1, \\ \text{XOR}(1, 1) &= 0, \end{aligned}$$

using the rule  $\text{sign}(f(\mathbf{x}))$ , that is, the answer is 1 if  $f(\mathbf{x}) \geq 0$  and 0 if  $f(\mathbf{x}) < 0$ .

1. Find a suitable set of parameters for this task. A guess is fine, but show that indeed it solves the above task.
2. Will it be possible to represent the XOR function if we replace the  $\max$  function with the identity function (i.e.,  $\mathbf{h} = U^T \mathbf{x}$ )? If so, show how. If not, explain why not.
3. What happens if we use  $\text{ReLU}(U^T \mathbf{x})$  instead of  $\max$ ? Analyze whether the XOR function can still be represented.

## Q3

Using Numpy only, implement the model described in Q2 and learn an optimal set of parameters using the Gradient Descent (GD) algorithm.

1. Create a dataset consisting of the following 4 examples  $\{(0, 0), (0, 1), (1, 0), (1, 1)\}$ , and assign them the following labels  $\{-1, 1, 1, -1\}$  correspondingly.
2. We will use the squared loss to optimize the parameters:  $\mathcal{L}(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$ .

3. Plot the loss value as a function of the number of epochs.

4. Additionally, randomly choose two trainable parameters, and plot their evolution epochs.

Note that you may need several random initializations to converge to the optimal solution (the one that correctly classifies all examples).

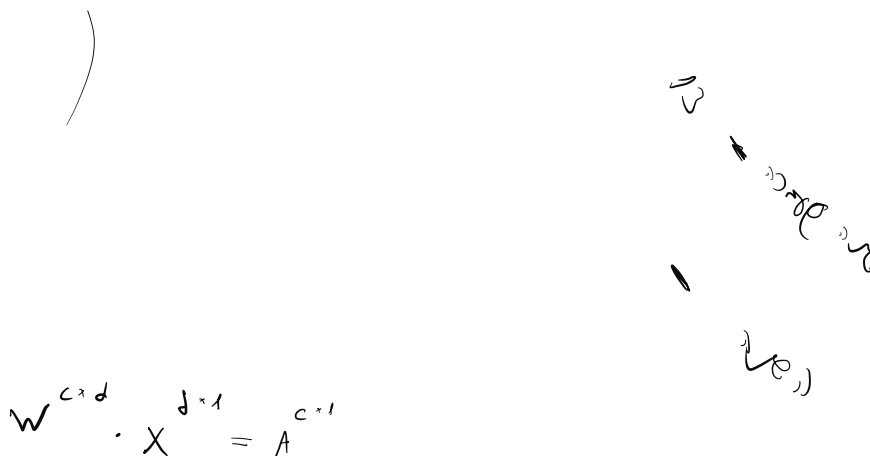
## Submission Instructions

Please submit a report with the answers to Q1 and Q2, the plots for Q3, and your code. Add to the report clear instructions how to run your code. Additionally:

1. Code should include comments explaining key steps of the implementation.
2. The report should include a short section on the challenges you faced during the implementation.
3. Code should be submitted in .py or .ipynb file format only.
4. Add to the report your name and ID.
5. Pack your submission files with your favorite file archiver (e.g., .rar, .zip).
6. The archive name should be your ID. If the exercise is done in pairs, the name of the file should be in the following format ID1\_ID2. Only one member needs to submit the solution.

Good Luck

Ran



Handwritten notes and a diagram. At the top left is a large closing parenthesis ')'. To the right is a diagram showing a vector  $\vec{w}$  and a vector  $\vec{v}$  originating from the same point, with a vector  $\vec{u}$  pointing from  $\vec{w}$  to  $\vec{v}$ . Below the diagram is the equation  $\vec{w}^{C \times D} \cdot \vec{X}^{D \times 1} = \vec{A}^{C \times 1}$ .

Q1

Consider a multivariate logistic regression problem. Recall that given a vector  $x \in \mathbb{R}^d$ , and a parameter matrix  $W \in \mathbb{R}^{d \times c}$ , where  $c$  is the number of classes, the model can be written in the following manner using the Softmax function:

$$\text{softmax}(z)_{[i]} = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

$$\text{softmax}(z)_{[i]} = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

where,  $z = W^T x$ .

1. Show that  $\text{softmax}(z) = \text{softmax}(z + m \cdot \mathbf{1})$  for every scalar  $m$ , where  $\mathbf{1}$  is a vector of ones of appropriate size.
2. For  $c = 2$ , show that the Sigmoid function is equivalent to the Softmax function.
3. Present an alternative to the Sigmoid function which also maps from the real line to the  $[0, 1]$  interval (any valid solution will be acceptable here).

④

$$W^{c \times d} \cdot X^{d \times 1} = A^{c \times 1}$$

①

we have  $c$  classes.  $\rightarrow (\exp(z_1), \exp(z_2), \dots, \exp(z_n)) \rightarrow \left( \frac{\exp(z_1)}{\sum \exp(z_i)}, \frac{\exp(z_2)}{\sum \exp(z_i)}, \dots, \frac{\exp(z_n)}{\sum \exp(z_i)} \right)$

$$\frac{\exp(z_i)}{\sum_j \exp(z_j)} = \frac{\exp(z_i + m)}{\sum_j \exp(z_j + m)} = \frac{e^{z_i} \cdot e^m}{\sum_j (e^{z_j} \cdot e^m)} = e^m \cdot \frac{e^{z_i}}{\sum_j e^{z_j} \cdot e^m} = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

The sigmoid function  $\sigma(x) = \frac{e^x}{1 + e^x} = \frac{1}{1 + e^{-x}} \left( \frac{e^{z_1}}{e^{z_1} + e^{z_2}}, \frac{e^{z_2}}{e^{z_1} + e^{z_2}} \right)$

$c=2$  ⑤

$$\text{softmax}(z_1) = \frac{e^{z_1}}{e^{z_1} + e^{z_2}} \quad z' = z_1 - z_2 \rightarrow z_1 = z' + z_2 \quad e^{z_1} = e^{z'} \cdot e^{z_2}$$

$$\text{softmax}(z_1) = \frac{e^{z'} \cdot e^{z_2}}{e^{z'} \cdot e^{z_2} + e^{z_2}} = \frac{e^{z'} \cdot e^{z_2}}{e^{z_2} (1 + e^{z'})} = \frac{e^{z'}}{1 + e^{z'}} = \frac{e^{z_1 - z_2}}{1 + e^{z_1 - z_2}} = \frac{e^{z_1}}{1 + e^{z_1 - z_2}} = \frac{1}{1 + e^{-z_1}}$$

$$\text{softmax}(z_2) = \frac{e^{z_2}}{e^{z_1} + e^{z_2}} \quad z_2 = z_1 - z' \Rightarrow \frac{e^{z_1 - z'}}{e^{z_1} + e^{z_1 - z'}} = \frac{e^{-z'}}{1 + e^{-z'}} = \frac{1}{1 + e^{z'}}$$

$$\frac{e^{-z'} e^{z_1}}{e^{z_1} (1 + e^{-z'})} = \frac{e^{-z'}}{1 + e^{-z'}} = \frac{1}{1 + e^{z'}}$$

We need to prove that  $\sigma(z') = \text{softmax}(z_1)$  and  $\sigma(z_1) = \text{softmax}(z_2)$

$$\text{softmax}(z_1) = \frac{e^{z_1}}{e^{z_1} + e^{z_2}} \quad \sigma(z_1) = \frac{1}{1 + e^{-z_1}} = \frac{e^{z_1}}{1 + e^{z_1 - z_2}} \Rightarrow$$

$$\Rightarrow \frac{1}{1 + e^{-(z_1 - z_2)}} = \frac{1}{\frac{e^{z_1}}{e^{z_1}} + \frac{e^{z_2}}{e^{z_1}}} = \frac{1}{\frac{e^{z_1} + e^{z_2}}{e^{z_1}}} = \frac{e^{z_1}}{e^{z_1} + e^{z_2}} = \text{softmax}(z_1)$$

$z_1$  is unbiased

$$\text{softmax}(z_i) = \frac{e^{z_i}}{e^{z_1} + e^{z_2}}$$

$$\sigma(z_2) = \frac{1}{1 + e^{z_1 - z_2}} = \frac{1}{1 + e^{z_1 - z_2}} = \frac{1}{1 + \frac{e^{z_1}}{e^{z_2}}} = \frac{1}{\frac{e^{z_1} + e^{z_2}}{e^{z_2}}} = \frac{e^{z_2}}{e^{z_1} + e^{z_2}}$$

also  $\text{softmax}(z_2) = \sigma(z_2)$  is obvious

We can use this  $f(x)$

$$f(x) = \begin{cases} |x| & (x < 1 \wedge x > 0) \vee (x < 0 \wedge x > -1) \\ |x| & x \geq 1 \vee x \leq -1 \end{cases}$$



Q2

Let  $x \in \{0, 1\}^2$  be an input vector. Consider the following model (scalar function):

$$f(x) = w^T h + b_2$$

$$h = \max(U^T x + b_1, 0)$$

Where  $U \in \mathbb{R}^{2 \times 2}$ ,  $b_1 \in \mathbb{R}^2$ ,  $w \in \mathbb{R}^2$ ,  $b_2 \in \mathbb{R}$ , and the  $\max$  is taken element-wise.

Suppose we would like to represent with  $f(x)$  the XOR function, defined as:

$$\begin{aligned} \text{XOR}(0, 0) &= 0, \\ \text{XOR}(0, 1) &= 1, \\ \text{XOR}(1, 0) &= 1, \\ \text{XOR}(1, 1) &= 0, \end{aligned}$$

using the rule  $\text{sign}(f(x))$ , that is, the answer is 1 if  $f(x) \geq 0$  and 0 if  $f(x) < 0$ .

1. Find a suitable set of parameters for this task. A guess is fine, but show that indeed it solves the above task.
2. Will it be possible to represent the XOR function if we replace the  $\max$  function with the identity function (i.e.,  $h = U^T x$ )? If so, show how. If not, explain why not.
3. What happens if we use  $\text{ReLU}(U^T x)$  instead of  $\max$ ? Analyze whether the XOR function can still be represented.

$$U = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$$

$$b_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$b_2 = -0.5$$

$$w = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$f(x) = w^T \cdot \max(U^T \cdot x + b_1, 0)$$

$$f([1, 0]) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}^T \cdot \max \left( \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}^T \begin{bmatrix} 1 \\ 0 \end{bmatrix}, 0 \right) - 0.5 = 0.5 \Rightarrow 1$$

$$f([0, 1]) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}^T \cdot \max \left( \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}^T \begin{bmatrix} 0 \\ 1 \end{bmatrix}, 0 \right) - 0.5 = 0.5 \Rightarrow 1$$

$$f([0, 0]) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}^T \cdot \max \left( \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}^T \begin{bmatrix} 0 \\ 0 \end{bmatrix}, 0 \right) - 0.5 \Rightarrow 0$$

$$f([1, 1]) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}^T \cdot \max \left( \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}^T \begin{bmatrix} 1 \\ 1 \end{bmatrix}, 0 \right) - 0.5 \Rightarrow 0$$

②

can we remove the  $\max$  and find solution?

$$f(x) = \underbrace{w^{1 \times 2}}_{1 \times 2} \cdot \underbrace{U^{2 \times 2}}_{2 \times 2} \cdot x + b = \underbrace{[a_1, a_2]}_{1 \times 2} \cdot \underbrace{\begin{bmatrix} x_1 & x_2 \end{bmatrix}}_{2 \times 1}^T + b = a_1 x_1 + a_2 x_2 + b$$

$$(x_1, x_2) = (0, 0) \rightarrow 0 \quad a_1 \cdot 0 + a_2 \cdot 0 + b_2 = 0 \Rightarrow b_2 = 0$$

$$(x_1, x_2) = (1, 0) \rightarrow a_1 + 0 + 0 = 1 \rightarrow a_1 = 1$$

$$(x_1, x_2) = (0, 1) \rightarrow 0 + a_2 + 0 = 1 \rightarrow a_2 = 1$$

$$(x_1, x_2) = (1, 1) \rightarrow 1 + 1 + 0 \neq 1 \rightarrow \text{לא מתאים}$$

יש להוסיף עוד תנאי

max  $\{0, x\}$

$$\textcircled{3} \quad f(x) = w^T \cdot \text{ReLU}(u^T \cdot x) - b \quad \text{ReLU} = \max(0, x)$$

$$f(x) = [w_1, w_2]^{1 \times 2} \cdot \max(u^T x) - b$$

If we take the max element wise then it is the same as section ①, but if we take the max by how has higher norm then we get the case in section ②