

Solution for Homework 3

This solution is based on the work submitted by some of your colleagues.

Question 1

A) We define the sign function as :

$$\text{sign}(x) = \begin{cases} 1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases}$$

Using only indicator functions, the sign function can be written as :

$$\text{sign}(x) = \mathbb{1}_{x>0}(x) - \mathbb{1}_{x<0}(x)$$

B) The derivative of the relu function can be expressed as the following :

$$g'(x) = \begin{cases} 1, & \text{when } x > 0 \\ 0, & \text{when } x < 0 \end{cases}$$

$$g'(x) = \text{Heaviside}(x)$$

Note that the derivative does not exist at $x=0$.

C) By fundamental theorem of calculus :

$$\text{ReLU}(x) = \int H(x') dx'$$

Also :

$$\begin{aligned} \text{ReLU}(x) &= \begin{cases} xH(x) & \text{if } x \neq 0 \\ 0 & \text{if } x = 0 \end{cases} \\ \text{ReLU}(x) &= \begin{cases} x - xH(-x) & \text{if } x \neq 0 \\ 0 & \text{if } x = 0 \end{cases} \end{aligned}$$

D)

$$\text{For } x > 0 : \lim_{k \rightarrow \infty} a(x) = \lim_{k \rightarrow \infty} \frac{1}{1 + e^{-kx}} = \frac{1}{1 + e^{-\infty}} = 1$$

$$\text{For } x = 0 : \lim_{k \rightarrow \infty} a(0) = \lim_{k \rightarrow \infty} \frac{1}{1 + e^{-k \cdot 0}} = \frac{1}{1 + 1} = \frac{1}{2}$$

$$\text{For } x < 0 : \lim_{k \rightarrow \infty} a(x) = \lim_{k \rightarrow \infty} \frac{1}{1 + \underbrace{e^{-k \cdot x}}_{k \cdot x < 0}} = \frac{1}{1 + \infty} = 0$$

$$\lim_{k \rightarrow \infty} a(x) = \begin{cases} 1 & \text{si } x > 0 \\ 1/2 & \text{si } x = 0 \\ 0 & \text{si } x < 0 \end{cases} = H(x)$$

E)

Knowing that the sigmoid function is : $\sigma(x) = \frac{1}{1+e^{-x}}$, we can calculate its derivative.

$$\begin{aligned} \frac{\partial \sigma(x)}{\partial x} &= \frac{\partial \frac{1}{1+e^{-x}}}{\partial x} = \frac{\partial (1 + e^{-x})^{-1}}{\partial x} \\ &= -(1 + e^{-x})^{-2} \times -e^{-x} \quad (\text{Chain Rule}) \\ &= \frac{e^{-x}}{(1 + e^{-x})^2} \\ &= \frac{1}{(1 + e^{-x})} \frac{e^{-x}}{(1 + e^{-x})} \\ &= \sigma(x) \frac{e^{-x}}{(1 + e^{-x})} \\ &= \sigma(x) \left(\frac{1 + e^{-x} - 1}{(1 + e^{-x})} \right) \\ &= \sigma(x) \left(\frac{1 + e^{-x}}{(1 + e^{-x})} - \frac{1}{(1 + e^{-x})} \right) \\ &= \sigma(x)(1 - \sigma(x)) \end{aligned}$$

let $i, j \in \llbracket 1, n \rrbracket \times \llbracket 1, n \rrbracket$

$$\begin{aligned} \frac{d\sigma(x)_i}{dx_j} &= \frac{d}{dx_j} \left(\frac{1}{1 + e^{-x_i}} \right) \\ &= \frac{-\delta_{ij} e^{-x_i}}{(1 + e^{-x_i})^2} \\ &= -\delta_{ij} \frac{1 + e^{-x_i} - 1}{(1 + e^{-x_i})^2} \\ &= \delta_{ij} \sigma(x)_i (1 - \sigma(x)_i) \end{aligned}$$

So,

$$\forall i, j \in \llbracket 1, n \rrbracket \times \llbracket 1, n \rrbracket, i \neq j \implies \frac{d\sigma(x)_i}{dx_j} = 0$$

thus

$$\forall i \in \llbracket 1, n \rrbracket, \quad \frac{d\sigma(x)_i}{dx_i} = \sigma(x)_i (1 - \sigma(x)_i)$$

the Jacobian matrix of σ is :

$$\frac{\partial \sigma(x)}{\partial x} = \text{diag}(\sigma(x)) (I_n - \text{diag}(\sigma(x)))$$

F)

Showing that $\ln(\sigma(x)) = -\text{softplus}(-x)$

$$\begin{aligned} \ln(\sigma(x)) &= \ln\left(\frac{1}{1 + e^{-x}}\right) \\ &= \ln(1) - \ln(1 + e^{-x}) \quad (\text{Log properties}) \\ &= 0 - \ln(1 + e^{-x}) = -\text{softplus}(-x) \end{aligned}$$

G)

Showing that $\text{softplus}(x) - \text{softplus}(-x) = x$:

$$\begin{aligned} \text{softplus}(x) - \text{softplus}(-x) &= \ln(1 + e^x) - \ln(1 + e^{-x}) \\ &= \ln(1 + e^x) - \ln(e^{-x}(1 + e^x)) \\ &= \ln\left(\frac{1 + e^x}{1 + e^x} \frac{1}{e^{-x}}\right) \\ &= \ln\left(\frac{1}{e^{-x}}\right) \\ &= \ln(e^x) \\ &= x \end{aligned}$$

And,

$$\begin{aligned} \frac{d}{dx}[\text{softplus}(x) - \text{softplus}(-x)] &= \frac{d}{dx} \cdot x \\ \frac{d}{dx} \text{softplus}(x) &= 1 + \frac{d}{dx} \text{softplus}(-x) \\ &= 1 - \frac{d}{dx}[-\text{softplus}(-x)] \\ &= 1 - \frac{d}{dx} \ln(\sigma(x)) \\ &= 1 - \frac{1}{\sigma(x)} \cdot \sigma(x)[1 - \sigma(x)] \\ \frac{d}{dx} \text{softplus}(x) &= \sigma(x) \end{aligned}$$

H)

Showing that the Softmax is translation-invariant :

$$s(x+c)_i = \frac{e^{x_i+c}}{\sum_j e^{x_j+c}} = \frac{e^{x_i} e^c}{\sum_j e^{x_j} e^c} = \frac{e^c e^{x_i}}{e^c \sum_j e^{x_j}} = \frac{e^{x_i}}{\sum_j e^{x_j}} = s(x)_i$$

Therefore, the softmax function is translation invariant.

I) The softmax function is not invariant under multiplication, because when $-1 < x < 1$, multiplication will shrink $|x|$, whereas when $x > 1$ and $x < -1$, multiplication will increase $|x|$. Multiplication would have the effect of making the softmax more certain of the most likely output, i.e., it would sharpen the distribution of probabilities computed by the softmax.

J)

Let's first calculate the partial derivative when $i=j$ (is a diagonal element) :

$$\begin{aligned} \frac{\partial \frac{e^{x_i}}{\sum_k e^{x_k}}}{\partial x_j} &= \frac{e^{x_i} \sum_k e^{x_k} - e^{x_j} e^{x_i}}{(\sum_k e^{x_k})^2} \quad (\text{Quotient rule}) \\ &= \frac{e^{x_i}}{\sum_k e^{x_k}} \frac{\sum_k e^{x_k} - e^{x_j}}{\sum_k e^{x_k}} \\ &= S(x)_i \left(\frac{\sum_k e^{x_k}}{\sum_k e^{x_k}} - \frac{e^{x_j}}{\sum_k e^{x_k}} \right) \\ &= S(x)_i (1 - S(x)_j) \\ &= S(x)_i - S(x)_i S(x)_j \end{aligned}$$

Now, let's calculate the partial derivative when $i \neq j$ (is a off-diagonal element) :

$$\begin{aligned} \frac{\partial \frac{e^{x_i}}{\sum_k e^{x_k}}}{\partial x_j} &= \frac{0 \sum_k e^{x_k} - e^{x_j} e^{x_i}}{(\sum_k e^{x_k})^2} \quad (\text{Quotient rule}) \\ &= \frac{-e^{x_i} e^{x_j}}{\sum_k e^{x_k}} \\ &= -\frac{e^{x_i}}{\sum_k e^{x_k}} \frac{e^{x_j}}{\sum_k e^{x_k}} \\ &= -S(x)_i S(x)_j \end{aligned}$$

Therefore, $\frac{\partial S(x)_i}{\partial x_j} = S(x)_i \delta_{i=j} - S(x)_i S(x)_j$

k)

$$\begin{aligned}
\frac{\partial S(\mathbf{x})}{\partial \mathbf{x}} &= \begin{bmatrix} \frac{\partial S(\mathbf{x})_1}{\partial \mathbf{x}} \\ \frac{\partial S(\mathbf{x})_2}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial S(\mathbf{x})_n}{\partial \mathbf{x}} \end{bmatrix} = \begin{bmatrix} \frac{\partial S(\mathbf{x})_1}{\partial x_1} & \frac{\partial S(\mathbf{x})_1}{\partial x_2} & \cdots & \frac{\partial S(\mathbf{x})_1}{\partial x_n} \\ \frac{\partial S(\mathbf{x})_2}{\partial x_1} & \frac{\partial S(\mathbf{x})_2}{\partial x_2} & \cdots & \frac{\partial S(\mathbf{x})_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial S(\mathbf{x})_n}{\partial x_1} & \frac{\partial S(\mathbf{x})_n}{\partial x_2} & \cdots & \frac{\partial S(\mathbf{x})_n}{\partial x_n} \end{bmatrix} \\
&= \begin{bmatrix} S(\mathbf{x})_1(1 - S(\mathbf{x})_1) & S(\mathbf{x})_1(0 - S(\mathbf{x})_2) & \cdots & S(\mathbf{x})_1(0 - S(\mathbf{x})_n) \\ S(\mathbf{x})_2(0 - S(\mathbf{x})_1) & S(\mathbf{x})_2(1 - S(\mathbf{x})_2) & \cdots & S(\mathbf{x})_2(0 - S(\mathbf{x})_n) \\ \vdots & \vdots & \ddots & \vdots \\ S(\mathbf{x})_n(0 - S(\mathbf{x})_1) & S(\mathbf{x})_n(0 - S(\mathbf{x})_2) & \cdots & S(\mathbf{x})_n(1 - S(\mathbf{x})_n) \end{bmatrix} \\
&= \underbrace{\text{diag}(S(\mathbf{x}))}_{n \times n} - \underbrace{S(\mathbf{x})}_{n \times 1} \underbrace{S(\mathbf{x})^T}_{1 \times n} \\
&\quad \underbrace{\hspace{10em}}_{n \times n}
\end{aligned}$$

We assume that vectors are by default column vectors.

l)

$$\begin{aligned}
\nabla_u \log S(\mathbf{x}(\mathbf{u}))_i &= \frac{\partial}{\partial u} \log \left(\frac{e^{\mathbf{x}(u)_i}}{\sum_{j=1}^n e^{\mathbf{x}(u)_j}} \right) \\
&= \frac{\partial}{\partial u} \log e^{\mathbf{x}(u)_i} - \frac{\partial}{\partial u} \log \left(\sum_{j=1}^n e^{\mathbf{x}(u)_j} \right) \\
&= \frac{\partial}{\partial u} \mathbf{x}(u)_i - \frac{\frac{\partial}{\partial u} \sum_{k=1}^n e^{\mathbf{x}(u)_k}}{\sum_{j=1}^n e^{\mathbf{x}(u)_j}} \\
&= \frac{\partial}{\partial u} \mathbf{x}(u)_i - \sum_{k=1}^n \frac{e^{\mathbf{x}(u)_k}}{\sum_{j=1}^n e^{\mathbf{x}(u)_j}} \frac{\partial}{\partial u} \mathbf{x}(u)_k \\
&= \nabla_u \mathbf{x}(u)_i - \sum_{k=1}^n S(\mathbf{x}(u))_k \cdot \nabla_u \mathbf{x}(u)_k \\
&= \nabla_u \mathbf{x}(u)_i - \mathbb{E}_j [\nabla_u \mathbf{x}(u)_j]
\end{aligned}$$

m)

$$\begin{aligned}
\nabla_u L(x, c) &= \nabla_u \sum_{i=1}^K -c_i \log y_i \\
&= -\nabla_u \log y_k \quad \text{with } c_k = 1 \\
&= -\nabla_u x_k + \mathbb{E}_j \nabla_u x_j \\
&= \mathbb{E}_j \nabla_u x_j - \nabla_u x_k \\
\nabla_u L(x, c) &= \sum_{j=1}^k y_j \nabla_u x_j - \nabla_u x_k
\end{aligned}$$

Question 2

A) The dimension of $\mathbf{b}^{(1)}$ a vector of dimension $d_h \times 1$ which is equal to the number of neurons of the hidden layer.

The formula for the pre-activation vector can be expressed as (matrix form) - (where $\mathbf{P}^{(0)}(\mathbf{x})$ is the post activation of the first (input layer) layer) :

$$\mathbf{h}^a = \mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}$$

To calculate the specific element \mathbf{h}_j^a , we get :

$$\mathbf{h}_j^a = \mathbf{b}_j^{(1)} + \sum_i \mathbf{W}_{j,i}^{(1)} x_i$$

And

$$\mathbf{h}^s = \phi(\mathbf{h}^a)$$

B) The dimension of $\mathbf{W}^{(2)}$ is of dimension $m \times d_h$, which is the dimension of the output layer times the dimension of the hidden layer. The dimension of $\mathbf{b}^{(2)}$ is a vector of size $m \times 1$.

The output layer before activation (pre-activation) is calculated (in matrix form) as :

$$\mathbf{o}^a = \mathbf{W}^{(2)} \mathbf{h}^s + \mathbf{b}^{(2)}$$

We can also write it in detailed form for \mathbf{o}_k^a as :

$$\mathbf{o}_k^a = b_k^{(2)} + \sum_j \mathbf{W}_{k,j}^{(2)} h_j^s$$

C) We can write o_k^s as a function of o_k^a and o_j^a .

$$\begin{aligned} o_k^s &= \text{softmax}(\mathbf{o}^a)_k \\ &= \frac{\exp(o_k^a)}{\sum_{j=1}^m \exp(o_j^a)} \end{aligned}$$

All o_k^s are positive because the exponential function is strictly positive : $e^x > 0 \quad \forall x \in \mathbb{R}$. The softmax() function is a quotient of a exponential function (positive) by a sum of exponential functions (sum of positive values results in a positive value). The quotient of two positive values is also positive. Now let's prove that $\sum_k o_k^s = 1$.

$$\begin{aligned} \sum_k o_k^s &= \sum_k \text{softmax}(\mathbf{o}^a)_k \\ &= \sum_k \frac{\exp(o_k^a)}{\sum_j \exp(o_j^a)} \\ &= \frac{\exp(o_1^a)}{\sum_j \exp(o_j^a)} + \frac{\exp(o_2^a)}{\sum_j \exp(o_j^a)} + \dots + \frac{\exp(o_m^a)}{\sum_j \exp(o_j^a)} \\ &= \frac{\sum_j \exp(o_j^a)}{\sum_j \exp(o_j^a)} \\ &= 1 \end{aligned}$$

The fact that $\sum_k o_k^s = 1$ is useful because it acts as a probability distribution over the different possible outcomes. By probability theory the sum of a discrete probability function must equals to 1.

D)

Let \mathbf{x} be a 2-dimensional vector. One can represent a 2-class categorical probability using softmax :

$$\begin{aligned} P(\mathbf{x}_1) &= S(\mathbf{x}_1) = \frac{e^{x_1}}{e^{x_1} + e^{x_2}} = \frac{1}{1 + e^{x_2 - x_1}} \\ P(\mathbf{x}_2) &= S(\mathbf{x}_2) = \frac{e^{x_2}}{e^{x_1} + e^{x_2}} = \frac{e^{x_2 - x_1}}{1 + e^{x_2 - x_1}} \end{aligned}$$

If we denote z as a scalar function of \mathbf{x} , i.e. $z = \mathbf{x}_2 - \mathbf{x}_1$, then

$$S(\mathbf{x}_1) = \frac{1}{1 + e^z} = \sigma(z)$$

$$S(\mathbf{x}_2) = \frac{e^z}{1 + e^z} = 1 - \sigma(z)$$

Hence $S(\mathbf{x}) = [\sigma(z), 1 - \sigma(z)]^\top$, where $z = \mathbf{x}_2 - \mathbf{x}_1$

E)

$$\begin{aligned}
L_{MSE}(\sigma(\mathbf{o}^a), y) &= \frac{1}{n} \sum_{i=1}^n \left(\sigma(\mathbf{o}_i^a) - y^{(i)} \right)^2 \\
\frac{\partial L_{MSE}(\sigma(\mathbf{o}^a), y)}{\partial \mathbf{o}^a} &= \frac{\partial \left(\frac{1}{n} \sum_{i=1}^n \left(\sigma(\mathbf{o}^{a(i)}) - y^{(i)} \right)^2 \right)}{\partial \mathbf{o}^a} \\
&= \frac{1}{n} \sum_{i=1}^n \frac{\partial \left(\left(\sigma(\mathbf{o}^{a(i)}) - y^{(i)} \right)^2 \right)}{\partial \mathbf{o}^{a(i)}} \\
&= \frac{1}{n} \sum_{i=1}^n \frac{\partial \left(\sigma(\mathbf{o}^{a(i)}) - y^{(i)} \right)^2}{\partial \left(\sigma(\mathbf{o}^{a(i)}) - y^{(i)} \right)} \frac{\partial \left(\sigma(\mathbf{o}^{a(i)}) - y^{(i)} \right)}{\partial \sigma(\mathbf{o}^{a(i)})} \frac{\partial \sigma(\mathbf{o}^{a(i)})}{\partial \mathbf{o}^{a(i)}} \\
&= \frac{1}{n} \sum_{i=1}^n 2 \left(\sigma(\mathbf{o}^{a(i)}) - y^{(i)} \right) \frac{\partial \sigma(\mathbf{o}^{a(i)})}{\partial \mathbf{o}^{a(i)}}
\end{aligned}$$

And :

$$\begin{aligned}
\frac{\partial \sigma(\mathbf{o}^{a(i)})}{\partial \mathbf{o}^{a(i)}} &= \partial \frac{\frac{1}{1 + \exp(-\mathbf{o}^{a(i)})}}{\partial \mathbf{o}^{a(i)}} \\
&= \frac{\exp(-\mathbf{o}^{a(i)})}{(1 + \exp(-\mathbf{o}^{a(i)}))^2} \\
&= \frac{1}{1 + \exp(-\mathbf{o}^{a(i)})} \frac{\exp(-\mathbf{o}^{a(i)})}{1 + \exp(-\mathbf{o}^{a(i)})} \\
&= \frac{1}{1 + \exp(-\mathbf{o}^{a(i)})} \frac{1 + \exp(-\mathbf{o}^{a(i)}) - 1}{1 + \exp(-\mathbf{o}^{a(i)})} \\
&= \sigma(\mathbf{o}^{a(i)}) (1 - \sigma(\mathbf{o}^{a(i)}))
\end{aligned}$$

Finally :

$$\frac{\partial L_{MSE}(\sigma(\mathbf{o}^a), y)}{\partial \mathbf{o}^a} = \frac{1}{n} \sum_{i=1}^n 2 \left(\sigma(\mathbf{o}^{a(i)}) - y^{(i)} \right) \sigma(\mathbf{o}^{a(i)}) (1 - \sigma(\mathbf{o}^{a(i)}))$$

F)

$$L_{CE}(\sigma(\mathbf{o}^a), y) = \sum_{i=1}^n - \left(y^{(i)} \log \left(\sigma(\mathbf{o}^{a(i)}) \right) + (1 - y^{(i)}) \log \left(1 - \sigma(\mathbf{o}^{a(i)}) \right) \right)$$

$$\begin{aligned}
\frac{\partial L_{CE}(\sigma(\mathbf{o}^a), y)}{\partial \mathbf{o}^a} &= \frac{\partial (\sum_{i=1}^n - (y^{(i)} \log(\sigma(\mathbf{o}^{a(i)})) + (1 - y^{(i)}) \log(1 - \sigma(\mathbf{o}^{a(i)}))))}{\partial \mathbf{o}^a} \\
&= \sum_{i=1}^n \frac{\partial (- (y^{(i)} \log(\sigma(\mathbf{o}^{a(i)})) + (1 - y^{(i)}) \log(1 - \sigma(\mathbf{o}^{a(i)}))))}{\partial \mathbf{o}^{a(i)}} \\
&= \sum_{i=1}^n \frac{\partial (-y^{(i)} \log(\sigma(\mathbf{o}^{a(i)})) - (1 - y^{(i)}) \log(1 - \sigma(\mathbf{o}^{a(i)})))}{\partial \mathbf{o}^{a(i)}} \\
&= \sum_{i=1}^n \frac{\partial (-y^{(i)} \log(\sigma(\mathbf{o}^{a(i)})))}{\partial \mathbf{o}^{a(i)}} - \frac{\partial ((1 - y^{(i)}) \log(1 - \sigma(\mathbf{o}^{a(i)})))}{\partial \mathbf{o}^{a(i)}} \\
&= \sum_{i=1}^n \left(-y^{(i)} (1 - \sigma(\mathbf{o}^{a(i)})) + (1 - y^{(i)}) \sigma(\mathbf{o}^{a(i)}) \right) \\
&= \sum_{i=1}^n \left(-y^{(i)} + \sigma(\mathbf{o}^{a(i)}) y^{(i)} + \sigma(\mathbf{o}^{a(i)}) - y^{(i)} \sigma(\mathbf{o}^{a(i)}) \right) \\
&= \sum_{i=1}^n \left(\sigma(\mathbf{o}^{a(i)}) - y^{(i)} \right)
\end{aligned}$$

G) The derivative of the MSE is equivalent to the derivative of CE times $\sigma(1-\sigma)$. This term does not depend on the target, and vanishes when the output σ is near 0 or 1. In other words, when the predicted probability is around the extreme values (0 and 1), and especially when it does not correspond to the desired target, the MSE gradient will be very small which will hinder and considerably limit any useful correction towards the right target, corrections reflected in the term $\sigma - y$. Hence, the CE is a more appropriate choice for binary classification since it does not exhibit the same issue. H)

$$L(x, y) = -\log o_y^s(x) = -\log \left(\frac{e^{o_y^a}}{\sum_{i=1}^{d_o} d_o e^{o_i^a}} \right) = -o_y^a + \log \left(\sum_{i=1}^{d_o} e^{o_i^a} \right)$$

I)

$$\hat{R}(\theta) = \sum_D L(x^{(i)}, y^{(i)})$$

The parameters θ are the two sets of W and b connecting the input layer to the hidden layer and the hidden layer to the output layer. In total, there are $n_{theta} = d_h d + d_h + d_o d_h + d_o$. The optimization problem is

$$\arg \min_{\theta} \hat{R}(\theta)$$

$$\hat{R}(\theta) = \sum_D L(x^{(i)}, y^{(i)})$$

The parameters θ are the two sets of W and b connecting the input layer to the hidden layer and the hidden layer to the output layer. In total, there are $n_{theta} = d_h d + d_h + d_o d_h + d_o$. The optimization problem is

$$\arg \min_{\theta} \hat{R}(\theta)$$

J)

$$\theta_{t+1} = \theta_t - \eta \nabla \hat{R}(\theta)$$

K)

$$\theta_{t+1} = \theta_t - \eta [\nabla \hat{R}(\theta) + 2\lambda\theta]$$

L) For this problem, we have to show that :

$$\frac{\partial L}{\partial \mathbf{o}^a} = \mathbf{o}^s - \text{onehot}_m(y)$$

Using the expression of L as a function of \mathbf{o}_k^a We have :

$$L(\mathbf{o}^a, y) = -\log(e^{\mathbf{o}_y^a}) + \log(\sum_j e^{\mathbf{o}_j^a}) = \log(\sum_j e^{\mathbf{o}_j^a}) - \mathbf{o}_y^a$$

When $k \neq y$:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{o}_k^a} &= \frac{1}{\sum_j e^{\mathbf{o}_j^a}} \frac{\partial \sum_j e^{\mathbf{o}_j^a}}{\partial \mathbf{o}_k^a} - 0 \\ &= \frac{e^{\mathbf{o}_k^a}}{\sum_j e^{\mathbf{o}_j^a}} = o_k^s \end{aligned}$$

When $k = y$:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{o}_k^a} &= \frac{1}{\sum_j e^{\mathbf{o}_j^a}} \frac{\partial \sum_j e^{\mathbf{o}_j^a}}{\partial \mathbf{o}_k^a} - \frac{\partial \mathbf{o}_y^a}{\partial \mathbf{o}_y^a} \\ &= \frac{e^{\mathbf{o}_k^a}}{\sum_j e^{\mathbf{o}_j^a}} - 1 = o_k^s - 1 \end{aligned}$$

Therefore, we can see that there is a minus one only when $k = y$, which means we can express :

$$\frac{\partial L}{\partial \mathbf{o}^a} = \mathbf{o}^s - \text{onehot}_m(y)$$

M)

$$\frac{\partial L}{\partial \mathbf{W}_{kj}^{(2)}} = \frac{\partial L}{\partial \mathbf{o}_k^a} \mathbf{h}_j^s$$

$$\frac{\partial L}{\partial \mathbf{b}_k^{(2)}} = \frac{\partial L}{\partial \mathbf{o}_k^a}$$

N)

$$\frac{\partial L}{\partial \mathbf{W}^{(2)}} = \frac{\partial L}{\partial \mathbf{o}^a} (\mathbf{h}^s)^T$$

$$\frac{\partial L}{\partial \mathbf{b}^{(2)}} = \frac{\partial L}{\partial \mathbf{o}^a}$$

where $\frac{\partial L}{\partial \mathbf{o}^a}$ is a $d_o \times 1$ vector and \mathbf{h}^s is a $d_h \times 1$ vector.

```
grad\_w2 = grad\_oa * hs.T \\  
grad\_b2 = grad\_oa
```

O)

The partial derivative of the loss L with respect to the output of the neurons at the hidden layer is :

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{h}_j^s} &= \sum_{k=1}^m \frac{\partial L}{\partial \mathbf{o}_k^a} \frac{\partial \mathbf{o}_k^a}{\partial \mathbf{h}_j^s} \\ &= \sum_{k=1}^m \frac{\partial L}{\partial \mathbf{o}_k^a} \frac{\partial \left(b_k^{(2)} + \sum_j \mathbf{W}_{k,j}^{(2)} h_j^s \right)}{\partial \mathbf{h}_j^s} \\ &= \sum_{k=1}^m \frac{\partial L}{\partial \mathbf{o}_k^a} \mathbf{W}_{k,j}^{(2)} \end{aligned}$$

P)

$$\frac{\partial L}{\partial \mathbf{h}_j^s} = (\mathbf{W}^{(2)})^T \frac{\partial L}{\partial \mathbf{o}^a}$$

```
grad_hs = w2.T * grad_oa
```

where the dimension of $(\mathbf{W}^{(2)})^T$ is of size $d_h \times m$, $\frac{\partial L}{\partial \mathbf{o}^a}$ is also a vector of size $m \times 1$ and $\frac{\partial L}{\partial \mathbf{h}^s}$ is a vector of size $d_h \times 1$.

Q)

Now, calculating the partial derivative with respect to the activation of the neurons at the hidden layer :

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{h}_j^a} &= \frac{\partial L}{\partial \mathbf{h}_j^s} \frac{\partial \mathbf{h}_j^s}{\partial \mathbf{h}_j^a} \\ &= \frac{\partial L}{\partial \mathbf{h}_j^s} \frac{\partial \phi(\mathbf{h}_j^a)}{\partial \mathbf{h}_j^a} \\ &= \frac{\partial L}{\partial \mathbf{h}_j^s} \sigma(\mathbf{h}_j^a)\end{aligned}$$

R) The gradient of the last equation can be written as :

$$\frac{\partial L}{\partial \mathbf{h}^a} = \frac{\partial L}{\partial \mathbf{h}^s} \cdot \text{diag}(\sigma(\mathbf{h}^a))$$

The dimension of $\frac{\partial L}{\partial \mathbf{h}^s}$ is a vector of size $d_h \times 1$, (calculated above).

The dimension of $\frac{\partial L}{\partial \mathbf{h}^a}$ is a vector of size $d_h \times 1$.

The dimension of $\text{diag}(\sigma(\mathbf{h}^a))$ is a matrix of size $d_h \times d_h$.

Question 3

a) Using the formula :

$$\text{out} = \left\lfloor \frac{\text{in} + 2p - d(k-1) - 1}{s} \right\rfloor + 1$$

Input layer :

$3 \times 128 \times 128$

Output 1st layer :

out - 1 : **$32 \times 64 \times 64$**

Output 2nd layer :

out - 2 : **$32 \times 32 \times 32$**

Output 3rd layer :

out - 3 : **$32 \times 64 \times 64$**

b) the total number of parameters needed for last layer = $3 \times 3 \times 32 \times 64 = 18432$

c)

$$\begin{aligned} [1, 1, 4, 4, 4, 1, 1] *_{\nu} [1/4, 1/4, 1/4] &= [1, 1, 4, 4, 4, 1, 1] * [1/4, 1/4, 1/4] \\ &= [3/2, 9/4, 3, 9/4, 3/2]. \end{aligned}$$

$$\begin{aligned} [1, 1, 4, 4, 4, 1, 1] *_{\epsilon} [1/4, 1/4, 1/4] &= [0, 0, 1, 1, 4, 4, 4, 1, 1, 0, 0] * [1/4, 1/4, 1/4] \\ &= [1/4, 1/2, 3/2, 9/4, 3, 9/4, 3/2, 1/2, 1/4]. \end{aligned}$$

$$\begin{aligned} [1, 1, 4, 4, 4, 1, 1] *_{\delta} [1/4, 1/4, 1/4] &= [0, 1, 1, 4, 4, 4, 1, 1, 0] * [1/4, 1/4, 1/4] \\ &= [1/2, 3/2, 9/4, 3, 9/4, 3/2, 1/2]. \end{aligned}$$

d) The previous kernel smoothes the input : there are no significant jumps after the application of this convolution.