

IFT6135-H2019 - REPRESENTATION LEARNING

ASSIGNMENT 1, THEORETICAL PART

SANAE LOTFI

NUM DE MATRICULE (POLY) : 1968682
NUM DE MATRICULE (UDEM) : 20147309

17th Februray, 2019

Due Date : February 16th, 2019

Instructions

- For all questions, show your work!
- Use a document preparation system such as LaTeX.
- Submit your answers electronically via Gradescope.

Question 1 (4-4-4-2). Using the following definition of the derivative and the definition of the Heaviside step function :

$$\frac{d}{dx}f(x) = \lim_{\epsilon \rightarrow 0} \frac{f(x + \epsilon) - f(x)}{\epsilon} \quad H(x) = \begin{cases} 1 & \text{if } x > 0 \\ \frac{1}{2} & \text{if } x = 0 \\ 0 & \text{if } x < 0 \end{cases}$$

1. Show that the derivative of the rectified linear unit $g(x) = \max\{0, x\}$, **wherever it exists**, is equal to the Heaviside step function.
2. Give two alternative definitions of $g(x)$ using $H(x)$.
3. Show that $H(x)$ can be well approximated by the sigmoid function $\sigma(x) = \frac{1}{1+e^{-kx}}$ asymptotically (i.e for large k), where k is a parameter.
- *4. Although the Heaviside step function is not differentiable, we can define its **distributional derivative**. For a function F , consider the functional $F[\phi] = \int_{\mathbb{R}} F(x)\phi(x)dx$, where ϕ is a smooth function (infinitely differentiable) with compact support ($\phi(x) = 0$ whenever $|x| \geq A$, for some $A > 0$).

Show that whenever F is differentiable, $F'[\phi] = -\int_{\mathbb{R}} F(x)\phi'(x)dx$. Using this formula as a definition in the case of non-differentiable functions, show that $H'[\phi] = \phi(0)$. ($\delta[\phi] \doteq \phi(0)$ is known as the Dirac delta function.)

Answer 1. Write your answer here.

1. We know that the rectified linear unit, defined as : $g(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$, is continuous and differentiable on \mathbb{R}_+^* and \mathbb{R}_-^* .

Its derivative on these two intervals is : $g'(x) = \begin{cases} 1 = H(x) & \text{if } x > 0 \\ 0 = H(x) & \text{if } x < 0 \end{cases}$.

g is differentiable in $x = 0$ if its left derivative and right derivative in $x = 0$ have the same value. We have :

$$\frac{d_+}{dx}g(x) = \lim_{\epsilon \rightarrow 0^+} \frac{g(0 + \epsilon) - g(0)}{\epsilon} = \lim_{\epsilon \rightarrow 0^+} \frac{\epsilon - 0}{\epsilon} = 1$$

and

$$\frac{d_-}{dx}g(x) = \lim_{\epsilon \rightarrow 0^-} \frac{g(0 + \epsilon) - g(0)}{\epsilon} = \lim_{\epsilon \rightarrow 0^-} \frac{0 - 0}{\epsilon} = 0$$

thus, g is not differentiable in $x = 0$.

We conclude that g differentiable on \mathbb{R}_+^* and \mathbb{R}_-^* and its derivative is equal to the Heaviside step function on those two intervals.

2. We can define g using H as follows : $\forall x \in \mathbb{R}, g(x) = x H(x)$ or $\forall x \in \mathbb{R}, g(x) = \int_{-\infty}^x H(x)$.
3. To show that $H(x)$ can be well approximated by the sigmoid function $\sigma(x) = \frac{1}{1+e^{-kx}}$ for large k , it's enough to study the value of σ when k takes a large value :
 - For $x = 0$, $\sigma(x) = \frac{1}{2} = H(x)$,
 - For $x > 0$, $\sigma(x) \approx \frac{1}{1+0} \approx 1 = H(x)$ (because $\lim_{k \rightarrow \infty} \exp(-kx) = 0$, $x > 0$).
 - For $x < 0$, $\sigma(x) \approx 0 = H(x)$ (because $\lim_{k \rightarrow \infty} \exp(-kx) = +\infty$ when $x < 0$).
4. By definition, we have : $F'[\phi] = \int_{\mathbb{R}} F'(x)\phi(x)dx$
We use the integration by parts to write, for given a and b in \mathbb{R} such that $a < b$:

$$\begin{aligned} \int_a^b F'(x)\phi(x)dx &= [F(x)\phi(x)]_a^b - \int_a^b F(x)\phi'(x)dx \\ &= F(b)\phi(b) - F(a)\phi(a) - \int_a^b F(x)\phi'(x)dx \end{aligned} \quad (1)$$

Since ϕ has a compact support, we can calculate the limit of the above equation when a goes to $-\infty$ and b goes to $+\infty$. This gives us (since ϕ is null for negative and positive high values) :

$$\begin{aligned} F'[\phi] &= \int_{\mathbb{R}} F'(x)\phi(x)dx = \lim_{a \rightarrow -\infty} \lim_{b \rightarrow +\infty} \int_a^b F'(x)\phi(x)dx \\ &= 0 - \int_{\mathbb{R}} F(x)\phi'(x)dx = - \int_{\mathbb{R}} F(x)\phi'(x)dx \end{aligned} \quad (2)$$

If we use this formula to calculate $H'[\phi]$, then we have (since ϕ is null for negative and positive high values) :

$$\begin{aligned} H'[\phi] &= - \int_{\mathbb{R}} H(x)\phi'(x)dx \\ &= - \int_{\mathbb{R}_+} \phi'(x)dx \\ &= - \left[\phi(x) \right]_0^{+\infty} \\ &= \phi(0) \end{aligned}$$

Question 2 (5-8-5-5). Let x be an n -dimensional vector. Recall the softmax function : $S : \mathbf{x} \in \mathbb{R}^n \mapsto S(\mathbf{x}) \in \mathbb{R}^n$ such that $S(\mathbf{x})_i = \frac{e^{\mathbf{x}_i}}{\sum_j e^{\mathbf{x}_j}}$; the diagonal function : $\text{diag}(\mathbf{x})_{ij} = \mathbf{x}_i$ if $i = j$ and $\text{diag}(\mathbf{x})_{ij} = 0$ if $i \neq j$; and the Kronecker delta function : $\delta_{ij} = 1$ if $i = j$ and $\delta_{ij} = 0$ if $i \neq j$.

1. Show that the derivative of the softmax function is $\frac{dS(\mathbf{x})_i}{d\mathbf{x}_j} = S(\mathbf{x})_i (\delta_{ij} - S(\mathbf{x})_j)$.
2. Express the Jacobian matrix $\frac{\partial S(\mathbf{x})}{\partial \mathbf{x}}$ using matrix-vector notation. Use $\text{diag}(\cdot)$.
3. Compute the Jacobian of the sigmoid function $\sigma(\mathbf{x}) = 1/(1 + e^{-\mathbf{x}})$.

4. Let \mathbf{y} and \mathbf{x} be n -dimensional vectors related by $\mathbf{y} = f(\mathbf{x})$, L be an unspecified differentiable loss function. According to the chain rule of calculus, $\nabla_{\mathbf{x}} L = (\frac{\partial \mathbf{y}}{\partial \mathbf{x}})^\top \nabla_{\mathbf{y}} L$, which takes up $\mathcal{O}(n^2)$ computational time in general. Show that if $f(\mathbf{x}) = \sigma(\mathbf{x})$ or $f(\mathbf{x}) = S(\mathbf{x})$, the above matrix-vector multiplication can be simplified to a $\mathcal{O}(n)$ operation.

Answer 2. 1. Let i and j be in $\{1, \dots, n\}$. The function $S_i : \mathbf{x} \in \mathbb{R}^n \mapsto S_i(\mathbf{x}) \in \mathbb{R}$, such that $S_i(\mathbf{x}) = S(\mathbf{x})_i = \frac{e^{\mathbf{x}_i}}{\sum_k e^{\mathbf{x}_k}}$, is differentiable with respect to x_j and we have :

$$\begin{aligned} \frac{dS(\mathbf{x})_i}{d\mathbf{x}_j} &= \frac{\frac{de^{\mathbf{x}_i}}{d\mathbf{x}_j} \sum_k e^{\mathbf{x}_k} - e^{\mathbf{x}_i} \frac{d\sum_k e^{\mathbf{x}_k}}{d\mathbf{x}_j}}{(\sum_k e^{\mathbf{x}_k})^2} \\ &= \frac{\delta_{ij} e^{\mathbf{x}_i} \sum_k e^{\mathbf{x}_k} - e^{\mathbf{x}_i} e^{\mathbf{x}_j}}{(\sum_k e^{\mathbf{x}_k})^2} \\ &= \frac{e^{\mathbf{x}_i}}{\sum_k e^{\mathbf{x}_k}} \left(\delta_{ij} - \frac{e^{\mathbf{x}_j}}{\sum_k e^{\mathbf{x}_k}} \right) \\ &= S(\mathbf{x})_i (\delta_{ij} - S(\mathbf{x})_j) \end{aligned} \quad (3)$$

2. Let i and j be in $\{1, \dots, n\}$. We have :

$$\begin{aligned} \left(\frac{\partial S(\mathbf{x})}{\partial \mathbf{x}} \right)_{ij} &= \frac{dS(\mathbf{x})_i}{d\mathbf{x}_j} \\ &= S(\mathbf{x})_i (\delta_{ij} - S(\mathbf{x})_j) \\ &= \begin{cases} S(\mathbf{x})_i S(\mathbf{x})_j & \text{if } i \neq j \\ S(\mathbf{x})_i S(\mathbf{x})_i - S(\mathbf{x})_i & \text{if } i = j \end{cases} \end{aligned} \quad (4)$$

Thus, we have :

$$\begin{aligned} \frac{\partial S(\mathbf{x})}{\partial \mathbf{x}} &= \begin{bmatrix} \frac{dS(\mathbf{x})_1}{d\mathbf{x}_1} & \dots & \frac{dS(\mathbf{x})_1}{d\mathbf{x}_n} \\ \vdots & \ddots & \vdots \\ \frac{dS(\mathbf{x})_n}{d\mathbf{x}_1} & \dots & \frac{dS(\mathbf{x})_n}{d\mathbf{x}_n} \end{bmatrix} \\ &= \begin{bmatrix} -S(\mathbf{x})_1 S(\mathbf{x})_1 & \dots & -S(\mathbf{x})_1 S(\mathbf{x})_n \\ \vdots & \ddots & \vdots \\ -S(\mathbf{x})_n S(\mathbf{x})_1 & \dots & -S(\mathbf{x})_n S(\mathbf{x})_n \end{bmatrix} + \begin{bmatrix} S(\mathbf{x})_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & S(\mathbf{x})_n \end{bmatrix} \\ &= -S(\mathbf{x}) S^\top(\mathbf{x}) + \text{diag}(S(\mathbf{x})) \end{aligned} \quad (5)$$

3. Let i and j be in $\{1, \dots, n\}$. We have :

$$\begin{aligned}
 \left(\frac{\partial \sigma(\mathbf{x})}{\partial \mathbf{x}}\right)_{ij} &= \frac{d(1/(1+e^{-x_i}))}{d\mathbf{x}_j} \\
 &= \frac{-de^{-x_i}/d\mathbf{x}_j}{(1+e^{-x_i})^2} \\
 &= \frac{\delta_{ij}e^{-x_i}}{(1+e^{-x_i})^2} \\
 &= \delta_{ij}\sigma(x_i)(1-\sigma(x_i)) \\
 &= \begin{cases} \sigma(x_i)(1-\sigma(x_i)) & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}
 \end{aligned} \tag{6}$$

Thus, we have :

$$\begin{aligned}
 \frac{\partial \sigma(\mathbf{x})}{\partial \mathbf{x}} &= \begin{bmatrix} \sigma(x_1)(1-\sigma(x_1)) & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma(x_n)(1-\sigma(x_n)) \end{bmatrix} \\
 &= \text{diag}(\sigma(\mathbf{x})) \text{diag}(\mathbf{1}_n - \sigma(\mathbf{x}))
 \end{aligned} \tag{7}$$

where $\mathbf{1}_n$ is the vector of length n with 1 in all its rows.

4. - If $f(\mathbf{x}) = \sigma(\mathbf{x})$, then

$$\nabla_{\mathbf{x}} L = \left(\frac{\partial \sigma(\mathbf{x})}{\partial \mathbf{x}}\right)^\top \nabla_{\mathbf{y}} L = \text{diag}(\sigma(\mathbf{x})) \text{diag}(\mathbf{1}_n - \sigma(\mathbf{x})) \nabla_{\mathbf{y}} L$$

The multiplication $\text{diag}(\mathbf{1}_n - \sigma(\mathbf{x})) \nabla_{\mathbf{y}} L$ is a $\mathcal{O}(n)$ operation, because it can be seen as an element-wise multiplication between $\sigma(\mathbf{x}) - 1$ and $\nabla_{\mathbf{y}} L$. This will result a vector that we can name \mathbf{v} , then $\nabla_{\mathbf{x}} L = \text{diag}(\sigma(\mathbf{x})) \mathbf{v}$ can be seen as an element-wise product as well between $\sigma(\mathbf{x})$ and \mathbf{v} , which is also a $\mathcal{O}(n)$ operation. Thus, the entire multiplication is simplified to a $\mathcal{O}(n)$ operation.

- If $f(\mathbf{x}) = S(\mathbf{x})$, then

$$\nabla_{\mathbf{x}} L = \left(\frac{\partial S(\mathbf{x})}{\partial \mathbf{x}}\right)^\top \nabla_{\mathbf{y}} L = (-S(\mathbf{x})S^\top(\mathbf{x}) + \text{diag}(S(\mathbf{x}))) \nabla_{\mathbf{y}} L = \text{diag}(S(\mathbf{x})) \nabla_{\mathbf{y}} L - S(\mathbf{x})(S^\top(\mathbf{x}) \nabla_{\mathbf{y}} L)$$

The multiplication $\text{diag}(S(\mathbf{x})) \nabla_{\mathbf{y}} L$ is a $\mathcal{O}(n)$ operation, because it can be seen as an element-wise multiplication between $S(\mathbf{x})$ and $\nabla_{\mathbf{y}} L$. The multiplication $(S^\top(\mathbf{x}) \nabla_{\mathbf{y}} L)$ is also a $\mathcal{O}(n)$ operation since its a multiplication between two vectros. This multiplication gives a scalar. Thus the rest is a product between a scalar and a vector. In conclusion, the matrix-vector multiplication is simplified to a $\mathcal{O}(n)$ operation.

Question 3 (3-3-3-3). Recall the definition of the softmax function : $S(\mathbf{x})_i = e^{x_i} / \sum_j e^{x_j}$.

1. Show that softmax is translation-invariant, that is : $S(\mathbf{x} + c) = S(\mathbf{x})$, where c is a scalar constant.
2. Show that softmax is not invariant under scalar multiplication. Let $S_c(\mathbf{x}) = S(c\mathbf{x})$ where $c \geq 0$. What are the effects of taking c to be 0 and arbitrarily large ?

3. Let \mathbf{x} be a 2-dimentional vector. One can represent a 2-class categorical probability using softmax $S(\mathbf{x})$. Show that $S(\mathbf{x})$ can be reparameterized using sigmoid function, i.e. $S(\mathbf{x}) = [\sigma(z), 1 - \sigma(z)]^\top$ where z is a scalar function of \mathbf{x} .
4. Let \mathbf{x} be a K -dimentional vector ($K \geq 2$). Show that $S(\mathbf{x})$ can be represented using $K - 1$ parameters, i.e. $S(\mathbf{x}) = S([0, y_1, y_2, \dots, y_{K-1}]^\top)$ where y_i is a scalar function of \mathbf{x} for $i \in \{1, \dots, K - 1\}$.

Answer 3. 1. Let c be in \mathbb{R} and i in $\{1, \dots, n\}$, we have :

$$\begin{aligned}
 S(\mathbf{x} + c)_i &= e^{\mathbf{x}_i + c} / \sum_j e^{\mathbf{x}_j + c} \\
 &= e^{\mathbf{x}_i} e^c / \sum_j e^{\mathbf{x}_j} e^c \\
 &= e^{\mathbf{x}_i} / \sum_j e^{\mathbf{x}_j} \\
 &= S(\mathbf{x})_i
 \end{aligned} \tag{8}$$

We conclude that $S(\mathbf{x} + c) = S(\mathbf{x})$, thus softmax is translation-invariant.

2. To prove this, it's enough to find a counterexample.

In fact, for $n = 2$, let's take $x_1 = 1$ and $x_2 = 0$ with $c = 2$. We have $S_c(\mathbf{x})_1 = S(2\mathbf{x})_1 = e^{2 \cdot 1} / (e^{2 \cdot 1} + e^{2 \cdot 0}) = e^2 / (e^2 + 1) \approx 0.88$. However, $S(\mathbf{x})_1 = e^1 / (e^1 + e^0) = e / (e + 1) \approx 0.71$.

It is clear that $S(c\mathbf{x}) \neq S(\mathbf{x})$ in general (the following part shows that it is not true in general for any n and $c = 0$ or arbitrarily large).

Let's find the expression of $S(c\mathbf{x})_i$ for a given i in $\{1, \dots, n\}$ and $c \geq 0$:

$$\begin{aligned}
 S(c\mathbf{x})_i &= e^{c\mathbf{x}_i} / \sum_j e^{c\mathbf{x}_j} \\
 &= e^{c\mathbf{x}_i} / \sum_j e^{c\mathbf{x}_j} \\
 &= 1 / \sum_j e^{c(\mathbf{x}_j - \mathbf{x}_i)}
 \end{aligned} \tag{9}$$

- If $c = 0$, then $S(c\mathbf{x})_i = 1/n$. In this case, all the classes have the same probability and the transformation becomes not interesting (since our goal is to be able to find the class with the highest probability).
- If c is arbitrarily large, let $L = \operatorname{argmax}\{\mathbf{x}_j, j \in \{1, \dots, n\}\}$ be the set of the arguments all the elements of the vector \mathbf{x} that share the highest value. We put $k = \operatorname{card}(L)$. We have that :

$$S(c\mathbf{x})_i = \begin{cases} \frac{1}{k} & \text{if } i \in L \\ 0 & \text{if } i \notin L \end{cases} \tag{10}$$

We notice then that in this case, any little difference between the values of the vector \mathbf{x} is translated by big changes in terms of probability. In case $k = 1$, we obtain that the element of the vector \mathbf{x} with the highest value gets a probability of 1 while all other elements get a probability of 0. This makes us lose a lot of information concerning to which extent the element with the highest value is superior to other elements.

3. Let \mathbf{x} be a 2-dimentional vector. One can represent a 2-class categorical probability using softmax $S(\mathbf{x})$. We have :

$$\begin{aligned} S(\mathbf{x}) &= \begin{bmatrix} e^{\mathbf{x}_1}/e^{\mathbf{x}_1} + e^{\mathbf{x}_2} \\ e^{\mathbf{x}_2}/e^{\mathbf{x}_1} + e^{\mathbf{x}_2} \end{bmatrix} \\ &= \begin{bmatrix} 1/(1 + e^{\mathbf{x}_2 - \mathbf{x}_1}) \\ 1 - 1/(1 + e^{\mathbf{x}_2 - \mathbf{x}_1}) \end{bmatrix} \\ &= \begin{bmatrix} \sigma(\mathbf{x}_1 - \mathbf{x}_2) \\ 1 - \sigma(\mathbf{x}_1 - \mathbf{x}_2) \end{bmatrix} \\ &= [\sigma(z), 1 - \sigma(z)]^\top \end{aligned} \tag{11}$$

where $z = \mathbf{x}_1 - \mathbf{x}_2$.

4. Let \mathbf{x} be a K -dimentional vector ($K \geq 2$). We can write :

- $S(\mathbf{x})_1 = e^{\mathbf{x}_1} / \sum_{j=1}^K e^{\mathbf{x}_j} = 1/1 + \sum_{j=2}^K e^{\mathbf{x}_j - \mathbf{x}_1}$,
- $\forall i \in \{2, \dots, K\}, S(\mathbf{x})_i = e^{\mathbf{x}_i} / \sum_{j=1}^K e^{\mathbf{x}_j} = e^{\mathbf{x}_i - \mathbf{x}_1} / 1 + \sum_{j=2}^K e^{\mathbf{x}_j - \mathbf{x}_1}$

If we put $y_i = \mathbf{x}_{i+1} - \mathbf{x}_1, \forall i \in \{1, \dots, K-1\}$, then we have :

- $S(\mathbf{x})_1 = e^0/e^0 + \sum_{j=1}^{K-1} e^{y_j}$,
- $\forall i \in \{2, \dots, K\}, S(\mathbf{x})_i = e^{y_{i-1}}/e^0 + \sum_{j=1}^{K-1} e^{y_j}$

We conclude that : $S(\mathbf{x}) = S([0, y_1, y_2, \dots, y_{K-1}]^\top)$.

Question 4 (15). Consider a 2-layer neural network $y : \mathbb{R}^D \rightarrow \mathbb{R}^K$ of the form :

$$y(x, \Theta, \sigma)_k = \sum_{j=1}^M \omega_{kj}^{(2)} \sigma \left(\sum_{i=1}^D \omega_{ji}^{(1)} x_i + \omega_{j0}^{(1)} \right) + \omega_{k0}^{(2)}$$

for $1 \leq k \leq K$, with parameters $\Theta = (\omega^{(1)}, \omega^{(2)})$ and logistic sigmoid activation function σ . Show that there exists an equivalent network of the same form, with parameters $\Theta' = (\tilde{\omega}^{(1)}, \tilde{\omega}^{(2)})$ and tanh activation function, such that $y(x, \Theta', \tanh) = y(x, \Theta, \sigma)$ for all $x \in \mathbb{R}^D$, and express Θ' as a function of Θ .

Answer 4. Let x be in \mathbb{R} , we have :

$$\begin{aligned} \tanh(x) &= (e^x - e^{-x})/(e^x + e^{-x}) \\ &= 2 \times 1/(1 + e^{-2x}) - 1 \\ &= 2\sigma(2x) - 1 \end{aligned} \tag{12}$$

Thus, $\sigma(x) = \frac{1}{2} \tanh(\frac{x}{2}) + \frac{1}{2}$.

Let $\mathbf{x} \in \mathbb{R}^D$. We write then :

$$\begin{aligned}
y(x, \Theta, \sigma)_k &= \sum_{j=1}^M \omega_{kj}^{(2)} \sigma \left(\sum_{i=1}^D \omega_{ji}^{(1)} x_i + \omega_{j0}^{(1)} \right) + \omega_{k0}^{(2)} \\
&= \sum_{j=1}^M \left[\frac{\omega_{kj}^{(2)}}{2} \tanh \left(\sum_{i=1}^D \frac{\omega_{ji}^{(1)}}{2} x_i + \frac{\omega_{j0}^{(1)}}{2} \right) + \frac{\omega_{kj}^{(2)}}{2} \right] + \omega_{k0}^{(2)} \\
&= \sum_{j=1}^M \left[\frac{\omega_{kj}^{(2)}}{2} \tanh \left(\sum_{i=1}^D \frac{\omega_{ji}^{(1)}}{2} x_i + \frac{\omega_{j0}^{(1)}}{2} \right) \right] + \sum_{j=1}^M \frac{\omega_{kj}^{(2)}}{2} + \omega_{k0}^{(2)}
\end{aligned} \tag{13}$$

We conclude then that there exists an equivalent network of the same form, with parameters $\Theta' = (\tilde{\omega}^{(1)}, \tilde{\omega}^{(2)})$ and tanh activation function, such that $y(x, \Theta', \tanh) = y(x, \Theta, \sigma)$ for all $x \in \mathbb{R}^D$, where :

- $\tilde{\omega}^{(1)} = \frac{1}{2}\omega^{(1)}$
- $\tilde{\omega}^{(2)}$ is defined such as : $\tilde{\omega}_{kj}^{(2)} = \begin{cases} \frac{1}{2}\omega_{kj}^{(2)} & \text{if } j \neq 0 \\ \sum_{l=1}^M \frac{\omega_{kl}^{(2)}}{2} + \omega_{k0}^{(2)} = \sum_{l=0}^M \frac{\omega_{kl}^{(2)}}{2} + \frac{\omega_{k0}^{(2)}}{2} & \text{if } j = 0 \end{cases}$

We can write Θ' in a more compact way using Θ :

$$\Theta' = \frac{1}{2}\Theta + \frac{1}{2}(Z, \omega^{(2)}B)$$

where Z is a zero matrix of the same dimensions as $\omega^{(1)}$ and B is a matrix of the same dimensions as $\omega^{(2)\top}$ with $B_{ij} = \begin{cases} 1 & \text{if } j = 1 \\ 0 & \text{otherwise} \end{cases}$ (the first column contains values equal to one and the rest of the matrix has values equal to 0).

An even more compact way consists of writing Θ' as follows :

$$\Theta' = \frac{1}{2}\Theta + \frac{1}{2}\Theta\tilde{B}$$

Where \tilde{B} is a $(M + D + 2) \times (K + M)$ matrix, defined as follows :

$$\tilde{B}_{ij} = \begin{cases} 1 & \text{if } i \in \{D + 2, \dots, D + M + 2\} \text{ and } j = 1 \\ 0 & \text{otherwise} \end{cases}$$

Question 5 (2-2-2-2). Given $N \in \mathbb{Z}^+$, we want to show that for any $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and any sample set $\mathcal{S} \subset \mathbb{R}^n$ of size N , there is a set of parameters for a two-layer network such that the output $y(\mathbf{x})$ matches $f(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{S}$. That is, we want to interpolate f with y on any finite set of samples \mathcal{S} .

1. Write the generic form of the function $y : \mathbb{R}^n \rightarrow \mathbb{R}^m$ defined by a 2-layer network with $N - 1$ hidden units, with linear output and activation function ϕ , in terms of its weights and biases $(\mathbf{W}^{(1)}, \mathbf{b}^{(1)})$ and $(\mathbf{W}^{(2)}, \mathbf{b}^{(2)})$.

2. In what follows, we will restrict $\mathbf{W}^{(1)}$ to be $\mathbf{W}^{(1)} = [\mathbf{w}, \dots, \mathbf{w}]^T$ for some $\mathbf{w} \in \mathbb{R}^n$ (so the rows of $\mathbf{W}^{(1)}$ are all the same). Show that the interpolation problem on the sample set $\mathcal{S} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\} \subset \mathbb{R}^n$ can be reduced to solving a matrix equation : $\mathbf{M}\tilde{\mathbf{W}}^{(2)} = \mathbf{F}$, where $\tilde{\mathbf{W}}^{(2)}$ and \mathbf{F} are both $N \times m$, given by

$$\tilde{\mathbf{W}}^{(2)} = [\mathbf{W}^{(2)}, \mathbf{b}^{(2)}]^\top \quad \mathbf{F} = [f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(N)})]^\top$$

Express the $N \times N$ matrix \mathbf{M} in terms of \mathbf{w} , $\mathbf{b}^{(1)}$, ϕ and $\mathbf{x}^{(i)}$.

- *3. **Proof with Relu activation.** Assume $\mathbf{x}^{(i)}$ are all distinct. Choose \mathbf{w} such that $\mathbf{w}^\top \mathbf{x}^{(i)}$ are also all distinct (Try to prove the existence of such a \mathbf{w} , although this is not required for the assignment - See Assignment 0). Set $\mathbf{b}_j^{(1)} = -\mathbf{w}^\top \mathbf{x}^{(j)} + \epsilon$, where $\epsilon > 0$. Find a value of ϵ such that \mathbf{M} is triangular with non-zero diagonal elements. Conclude. (Hint : assume an ordering of $\mathbf{w}^\top \mathbf{x}^{(i)}$.)
- *4. **Proof with sigmoid-like activations.** Assume ϕ is continuous, bounded, $\phi(-\infty) = 0$ and $\phi(0) > 0$. Decompose \mathbf{w} as $\mathbf{w} = \lambda \mathbf{u}$. Set $\mathbf{b}_j^{(1)} = -\lambda \mathbf{u}^\top \mathbf{x}^{(j)}$. Fixing \mathbf{u} , show that $\lim_{\lambda \rightarrow +\infty} \mathbf{M}$ is triangular with non-zero diagonal elements. Conclude. (Note that doing so preserves the distinctness of $\mathbf{w}^\top \mathbf{x}^{(i)}$.)

Answer 5. 1. Let's write the generic form of the function $y : \mathbb{R}^n \rightarrow \mathbb{R}^m$ defined by a 2-layer network with $N - 1$ hidden units, with linear output and activation function ϕ , in terms of its weights and biases $(\mathbf{W}^{(1)}, \mathbf{b}^{(1)})$ and $(\mathbf{W}^{(2)}, \mathbf{b}^{(2)})$:

$$y(x)_k = \sum_{j=1}^{N-1} \omega_{kj}^{(2)} \phi \left(\sum_{i=1}^n \omega_{ji}^{(1)} x_i + b_j^{(1)} \right) + b_k^{(2)}$$

We can write it in a compact way :

$$y(\mathbf{x})_k = \mathbf{w}_k^{(2)\top} \phi(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}) + b_k^{(2)}$$

where $\phi(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)})_j = \phi((\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)})_j) = \phi\left(\sum_{i=1}^n \omega_{ji}^{(1)} x_i + b_j^{(1)}\right)$, for $j \in \{1, \dots, N - 1\}$ and $\mathbf{W}^{(2)} = [\mathbf{w}_1^{(2)}, \dots, \mathbf{w}_m^{(2)}]^\top$.

Thus, we can write (using the same convention for ϕ applied to a vector) :

$$y(\mathbf{x}) = \mathbf{W}^{(2)} \phi(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}$$

2. The interpolation problem on the sample set $\mathcal{S} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\} \subset \mathbb{R}^n$ is about solving the equation :

$$\mathbf{F} = [f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(N)})]^\top = [y(\mathbf{x}^{(1)}), \dots, y(\mathbf{x}^{(N)})]^\top$$

.

Thus, we need to solve the equations :

$$\begin{aligned} f(\mathbf{x}^{(i)}) &= \mathbf{W}^{(2)} \phi(\mathbf{W}^{(1)} \mathbf{x}^{(i)} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}, \quad \forall i \in \{1, \dots, N\} \\ \implies f^\top(\mathbf{x}^{(i)}) &= \phi(\mathbf{W}^{(1)} \mathbf{x}^{(i)} + \mathbf{b}^{(1)})^\top \mathbf{W}^{(2)\top} + \mathbf{b}^{(2)\top}, \quad \forall i \in \{1, \dots, N\} \\ \implies f^\top(\mathbf{x}^{(i)}) &= [\phi(\mathbf{W}^{(1)} \mathbf{x}^{(i)} + \mathbf{b}^{(1)})^\top, 1] \begin{bmatrix} \mathbf{W}^{(2)\top} \\ \mathbf{b}^{(2)\top} \end{bmatrix}, \quad \forall i \in \{1, \dots, N\} \end{aligned} \tag{14}$$

where $\phi(\mathbf{W}^{(1)}\mathbf{x}^{(i)} + \mathbf{b}^{(1)})_j = \phi((\mathbf{W}^{(1)}\mathbf{x}^{(i)} + \mathbf{b}^{(1)})_j) = \phi(\mathbf{w}^\top \mathbf{x}^{(i)} + b_j^{(1)})$, for $j \in \{1, \dots, N-1\}$.

Thus, we need to solve the equation :

$$\mathbf{F} = [f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(N)})]^\top = \begin{bmatrix} \sigma(\mathbf{W}^{(1)}\mathbf{x}^{(1)} + \mathbf{b}^{(1)})^\top & 1 \\ \vdots & \vdots \\ \sigma(\mathbf{W}^{(1)}\mathbf{x}^{(N)} + \mathbf{b}^{(1)})^\top & 1 \end{bmatrix} \begin{bmatrix} \mathbf{W}^{(2)\top} \\ \mathbf{b}^{(2)\top} \end{bmatrix} \quad (15)$$

We put :

$$\tilde{\mathbf{W}}^{(2)} = [\mathbf{W}^{(2)}, \mathbf{b}^{(2)}]^\top \quad \text{and} \quad \mathbf{M} = \begin{bmatrix} \phi(\mathbf{W}^{(1)}\mathbf{x}^{(1)} + \mathbf{b}^{(1)})^\top & 1 \\ \vdots & \vdots \\ \phi(\mathbf{W}^{(1)}\mathbf{x}^{(N)} + \mathbf{b}^{(1)})^\top & 1 \end{bmatrix}$$

We can write the : $\mathbf{M}\tilde{\mathbf{W}}^{(2)} = \mathbf{F}$. Where \mathbf{M} is an $N \times N$ matrix defined by :

$$M_{ij} = \begin{cases} \phi(\mathbf{W}^{(1)}\mathbf{x}^{(i)} + \mathbf{b}^{(1)})_j = \phi(\mathbf{w}^\top \mathbf{x}^{(i)} + b_j^{(1)}), & \text{if } i \in \{1, \dots, N\} \text{ and } j \in \{1, \dots, N-1\} \\ 1, & \text{if } i \in \{1, \dots, N\} \text{ and } j = N \end{cases}$$

Conclusion : We showed that the interpolation problem on the sample set $\mathcal{S} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\} \subset \mathbb{R}^n$ can be reduced to solving a matrix equation : $\mathbf{M}\tilde{\mathbf{W}}^{(2)} = \mathbf{F}$, where $\tilde{\mathbf{W}}^{(2)}$ and \mathbf{M} are given by the expressions above.

3. Proof with Relu activation.

Let's assume that $\mathbf{x}^{(i)}$ are all distinct. and choose \mathbf{w} such that $\mathbf{w}^\top \mathbf{x}^{(i)}$ are also all distinct (the existence of \mathbf{w} was proved in assignment 0). Let $\mathbf{b}_j^{(1)} = -\mathbf{w}^\top \mathbf{x}^{(j)} + \epsilon$, where $\epsilon > 0$. Let's find a value of ϵ such that \mathbf{M} is triangular with non-zero diagonal elements.

We have from the previous question that :

$$\begin{aligned} M_{ij} &= \begin{cases} \phi(\mathbf{w}^\top \mathbf{x}^{(i)} + b_j^{(1)}), & \text{if } i \in \{1, \dots, N\} \text{ and } j \in \{1, \dots, N-1\} \\ 1, & \text{if } i \in \{1, \dots, N\} \text{ and } j = N \end{cases} \\ &= \begin{cases} \max(\mathbf{w}^\top \mathbf{x}^{(i)} + b_j^{(1)}, 0), & \text{if } i \in \{1, \dots, N\} \text{ and } j \in \{1, \dots, N-1\} \\ 1, & \text{if } i \in \{1, \dots, N\} \text{ and } j = N \end{cases} \\ &= \begin{cases} \max(\mathbf{w}^\top \mathbf{x}^{(i)} - \mathbf{w}^\top \mathbf{x}^{(j)} + \epsilon, 0), & \text{if } i \in \{1, \dots, N\} \text{ and } j \in \{1, \dots, N-1\} \\ 1, & \text{if } i \in \{1, \dots, N\} \text{ and } j = N \end{cases} \end{aligned}$$

Without loss of generality, we suppose an order on our vectors : $\mathbf{w}^\top \mathbf{x}^{(1)} > \dots > \mathbf{w}^\top \mathbf{x}^{(N)}$ (the inequalities are strict since the values are distinct by choice of \mathbf{w}).

We put $\epsilon = \min\{\mathbf{w}^\top \mathbf{x}^{(i)} - \mathbf{w}^\top \mathbf{x}^{(j)}, i < j, i \text{ and } j \text{ in } \{1, \dots, N\}\}$. Using our hypothesis, we see that $\epsilon > 0$. Furthermore, by definition of epsilon, we have for $i > j$ that : $\alpha_{ij} = \mathbf{w}^\top \mathbf{x}^{(i)} - \mathbf{w}^\top \mathbf{x}^{(j)} + \epsilon = \epsilon - (\mathbf{w}^\top \mathbf{x}^{(j)} - \mathbf{w}^\top \mathbf{x}^{(i)}) \leq 0$. Thus, $M_{ij} = \max(\alpha_{ij}, 0) = 0$ for $i > j$ and $j \neq N$. The matrix \mathbf{M} is upper triangular. The diagonal elements are given by :

$$M_{ii} = \begin{cases} \max(\mathbf{w}^\top \mathbf{x}^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)} + \epsilon, 0) = \max(\epsilon, 0) = \epsilon & \text{if } i \neq N \\ 1 & \text{if } i = N \end{cases} \quad (16)$$

Which means that for all $i \in \{1, \dots, N\}$, $M_{ii} > 0$. In other words, \mathbf{M} is triangular with non-zero diagonal elements. Thus \mathbf{M} is invertible and the matrix equation $\mathbf{M}\tilde{\mathbf{W}}^{(2)} = \mathbf{F}$ can be solved using the inverse of \mathbf{M} . Conclusion : the subset \mathcal{S} being fixed as well as the vector \mathbf{w} , we can find the weights matrix $\tilde{\mathbf{W}}^{(2)} = \mathbf{M}^{-1}\mathbf{F}$ such that we interpolate f with y .

4. Proof with sigmoid-like activations.

Let's assume ϕ is continuous, bounded, $\phi(-\infty) = 0$ and $\phi(0) > 0$. We decompose \mathbf{w} as $\mathbf{w} = \lambda\mathbf{u}$ and set $\mathbf{b}_j^{(1)} = -\lambda\mathbf{u}^\top \mathbf{x}^{(j)}$. Fixing \mathbf{u} , let's show that $\lim_{\lambda \rightarrow +\infty} \mathbf{M}$ is triangular with non-zero diagonal elements. We have :

$$\begin{aligned} M_{ij} &= \begin{cases} \phi(\mathbf{w}^\top \mathbf{x}^{(i)} + b_j^{(1)}) , & \text{if } i \in \{1, \dots, N\} \text{ and } j \in \{1, \dots, N-1\} \\ 1, & \text{if } i \in \{1, \dots, N\} \text{ and } j = N \end{cases} \\ &= \begin{cases} \phi(\lambda(\mathbf{u}^\top \mathbf{x}^{(i)} - \mathbf{u}^\top \mathbf{x}^{(j)})) , & \text{if } i \in \{1, \dots, N\} \text{ and } j \in \{1, \dots, N-1\} \\ 1, & \text{if } i \in \{1, \dots, N\} \text{ and } j = N \end{cases} \end{aligned}$$

The diagonal elements are given by :

$$M_{ii} = \begin{cases} \phi(0) & \text{if } i \neq N \\ 1 & \text{if } i = N \end{cases} \quad (17)$$

Which means that for all $i \in \{1, \dots, N\}$, $M_{ii} > 0$.

We take $\lambda > 0$. Just like the previous question, we assume (without loss of generality) an order on our vectors : $\mathbf{u}^\top \mathbf{x}^{(1)} > \dots > \mathbf{u}^\top \mathbf{x}^{(N)}$ (the inequalities are strict since the values are distinct by choice of $\lambda\mathbf{u}$ and since the factor $\lambda > 0$, this doesn't change the inequality).

Thus, for $i > j$ and $j \neq N$, $\lim_{\lambda \rightarrow +\infty} M_{ij} = \phi(-\infty) = 0$, using the order hypothesis.

From the other side, we assume that the value of $\phi(+\infty) = l$ exists (some continuous and bounded functions don't have a finite limit in $+\infty$). We have, for $i < j$ and $j \neq N$, $\lim_{\lambda \rightarrow +\infty} M_{ij} = \phi(+\infty) = l$.

Conclusion : $\lim_{\lambda \rightarrow +\infty} \mathbf{M}$ is triangular with non-zero diagonal elements. Thus it is invertible and the matrix equation $\mathbf{M}\tilde{\mathbf{W}}^{(2)} = \mathbf{F}$ can be solved using the inverse of the limit matrix of \mathbf{M} because

$\tilde{\mathbf{W}}^{(2)}$ and \mathbf{F} do not depend on λ . In other words, the subset \mathcal{S} being fixed as well as the vector \mathbf{w} , we can find the weights matrix $\tilde{\mathbf{W}}^{(2)} = \mathbf{M}^{-1}\mathbf{F}$ such that we interpolate f with y .

Question 6 (6). Compute the *full*, *valid*, and *same* convolution (with kernel flipping) for the following 1D matrices : $[1, 2, 3, 4] * [1, 0, 2]$

Answer 6. To compute the convolution, we will use the expression $(x*k)_{ij} = \sum_{p,q} x_{i+p,j+q} k_{r_1-p,r_2-q}$ where $r_1 \times r_2$ is the size of the kernel.

- *full convolution* : consists of adding the maximum zero-padding such that the convolution product with the kernel still takes into account elements from the original matrix. In this case, the maximum possible zero-padding is 2. We obtain the following 1D matrix : $[1, 2, 5, 8, 6, 8]$.
- *valid convolution* : consists of adding no zero-padding and performing the classic convolution product between the original matrix and the kernel. We obtain the following 1D matrix : $[5, 8]$.
- *same convolution* : consists of adding the enough zero-padding such that the output of the convolution product has the same dimension as the original matrix. In this case, the necessary zero-padding is 1. We obtain the following 1D matrix : $[2, 5, 8, 6]$.

Question 7 (5-5). Consider a convolutional neural network. Assume the input is a colorful image of size 256×256 in the RGB representation. The first layer convolves 64 8×8 kernels with the input, using a stride of 2 and no padding. The second layer downsamples the output of the first layer with a 5×5 non-overlapping max pooling. The third layer convolves 128 4×4 kernels with a stride of 1 and a zero-padding of size 1 on each border.

1. What is the dimensionality (scalar) of the output of the last layer ?
2. Not including the biases, how many parameters are needed for the last layer ?

Answer 7. 1. The dimensionality of the feature map after the first layer is given by the relation : $o = \lfloor \frac{i+2p-(d(k-1)+1)}{s} \rfloor + 1$, where i is the size of the input, p is the padding, d is the dilation, k is the size of the kernel and s is the stride.

Thus : $o = \lfloor \frac{256+2 \times 0 - (1 \times (8-1)+1)}{2} \rfloor + 1 = 125$ and the output shape of the first layer is $(64, 125, 125)$. In the same way we calculate the output of the next layers :

- The output shape of the second layer is $(64, 25, 25)$ (because $d = \lfloor \frac{125+2 \times 0 - 5}{1} \rfloor + 1 = 25$).
- The output shape of the third layer is $(128, 24, 24)$ (because $d = \lfloor \frac{25+2 \times 1 - 4}{1} \rfloor + 1 = 24$).

Thus the dimensionality (scalar) of the output of the last layer is : $128 \times 24 \times 24 = 73728$.

2. The number of parameters that are needed for the last layer is : $64 \times 4 \times 4 \times 128 = 131072$ parameters.

Question 8 (4-4-4). Assume we are given data of size $3 \times 64 \times 64$. In what follows, provide the correct configuration of a convolutional neural network layer that satisfies the specified assumption. Answer with the window size of kernel (k), stride (s), padding (p), and dilation (d , with convention $d = 1$ for no dilation). Use square windows only (e.g. same k for both width and height).

1. The output shape of the first layer is $(64, 32, 32)$.

- (a) Assume $k = 8$ without dilation.
- (b) Assume $d = 7$, and $s = 2$.
- 2. The output shape of the second layer is $(64, 8, 8)$. Assume $p = 0$ and $d = 1$.
 - (a) Specify k and s for pooling with non-overlapping window.
 - (b) What is output shape if $k = 8$ and $s = 4$ instead?
- 3. The output shape of the last layer is $(128, 4, 4)$.
 - (a) Assume we are not using padding or dilation.
 - (b) Assume $d = 2$, $p = 2$.
 - (c) Assume $p = 1$, $d = 1$.

Answer 8. We base the answers of all the following questions on the following expression of the output shape of a convolutional layer : $o = \lfloor \frac{i+2p-(d(k-1)+1)}{s} \rfloor + 1$, where i is the size of the input, p is the padding, d is the dilation, k is the size of the kernel and s is the stride. The only reason we didn't give the details of the calculations is because they are very repetitive and simple

Here are the answers :

- 1. The output shape of the first layer is $(64, 32, 32)$.
 - (a) Assuming $k = 8$ without dilation ($d = 1$), we can take : $p = 3$ and $s = 2$.
 - (b) Assuming $d = 7$, and $s = 2$, we can take : $k = 2$ and $p = 3$.
- 2. The output shape of the second layer is $(64, 8, 8)$. Assume $p = 0$ and $d = 1$.
 - (a) For pooling with non-overlapping window, we should have $k \leq s$, we can take : $k = 4$ and $s = 4$
 - (b) The output shape if $k = 8$ and $s = 4$ is : $(64, 7, 7)$.
- 3. The output shape of the last layer is $(128, 4, 4)$.
 - (a) Assuming we are not using padding ($p = 0$) or dilation ($d = 1$), we can take : $k = 5$ and $s = 1$.
 - (b) Assuming $d = 2$, $p = 2$, we can take : $k = 5$ and $s = 1$.
 - (c) Assuming $p = 1$, $d = 1$, we can take : $k = 4$ and $s = 2$.