

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. 1. For $x > 0$:

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

For $x < 0$ (such that $x+\epsilon < 0$ as $\epsilon \rightarrow 0$) :

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

For $x = 0$, approaching from the left gives 0, while approaching from the right gives 1, hence the derivative of $g(x)$ is not defined at $x = 0$.

\therefore The derivative of $g(x)$ for $x \neq 0$ is equal to $H(x)$.

2. (a) $g(x) = x * H(x)$
 (b) $g(x) = x * (1 - H(-x))$

3. For large k , for $x > 0$:

$$\lim_{k \rightarrow +\infty} \sigma(x) = \lim_{k \rightarrow +\infty} \frac{1}{1+e^{-kx}} = \lim_{k \rightarrow +\infty} \frac{1}{1+e^{<large \ negative \ number>}} = 1$$

For large k , for $x < 0$:

$$\lim_{k \rightarrow +\infty} \sigma(x) = \lim_{k \rightarrow +\infty} \frac{1}{1+e^{-kx}} = \lim_{k \rightarrow +\infty} \frac{1}{1+e^{<large \ positive \ number>}} = 0$$

$\therefore H(x)$ can be well approximated by the sigmoid function $\sigma(x)$ asymptotically.

4. $F'[\phi] = \int_{\mathbb{R}} F'(x)\phi(x)dx$

Performing this integration by parts :

$$F'[\phi] = \int_{\mathbb{R}} F'(x)\phi(x)dx = F(x)\phi(x)|_{-\infty}^{+\infty} - \int_{\mathbb{R}} F(x)\phi'(x)dx$$

Since $\phi(x)$ has a compact support, $F(x)\phi(x)|_{-\infty}^{+\infty} = 0 - 0 = 0$.

$$\therefore F'[\phi] = - \int_{\mathbb{R}} F(x)\phi'(x)dx \text{ (whenever } F \text{ is differentiable).}$$

Thus, $H'(\phi) = - \int_{\mathbb{R}} H(x)\phi'(x)dx = - \int_0^{+\infty} \phi'(x)dx = -\phi(x)|_0^{+\infty} = \phi(0)$ (again, since $\phi(x)$ has a compact support).

$$\therefore H'(\phi) = \phi(0)$$

Question 2. (5-8-5-5) Let \mathbf{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. To show that : $\frac{dS(\mathbf{x})_i}{d\mathbf{x}_j} = S(\mathbf{x})_i (\delta_{ij} - S(\mathbf{x})_j)$

Using quotient rule :

$$\frac{dS(\mathbf{x})_i}{d\mathbf{x}_j} = \frac{d}{d\mathbf{x}_j} \left(\frac{e^{\mathbf{x}_i}}{\sum_k e^{\mathbf{x}_k}} \right) = \frac{\sum_k e^{\mathbf{x}_k} \cdot \frac{d}{d\mathbf{x}_j} (e^{\mathbf{x}_i}) - e^{\mathbf{x}_i} \cdot \frac{d}{d\mathbf{x}_j} (\sum_k e^{\mathbf{x}_k})}{(\sum_k e^{\mathbf{x}_k})^2}$$

$$\text{Here, } \frac{d}{d\mathbf{x}_j} (e^{\mathbf{x}_i}) = \begin{cases} e^{\mathbf{x}_i} & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases} \implies \frac{d}{d\mathbf{x}_j} (e^{\mathbf{x}_i}) = \delta_{ij} \cdot e^{\mathbf{x}_i}$$

$$\text{Also, } \frac{d}{d\mathbf{x}_j} (\sum_k e^{\mathbf{x}_k}) = e^{\mathbf{x}_j}$$

Thus,

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

$$\therefore \frac{dS(\mathbf{x})_i}{d\mathbf{x}_j} = S(\mathbf{x})_i (\delta_{ij} - S(\mathbf{x})_j)$$

2. Jacobian of softmax :

$$\begin{aligned}
 \frac{\partial S(\mathbf{x})}{\partial \mathbf{x}} &= \begin{bmatrix} \frac{\partial S(\mathbf{x})_1}{\partial \mathbf{x}_1} & \frac{\partial S(\mathbf{x})_1}{\partial \mathbf{x}_2} & \cdots & \frac{\partial S(\mathbf{x})_1}{\partial \mathbf{x}_n} \\ \frac{\partial S(\mathbf{x})_2}{\partial \mathbf{x}_1} & \frac{\partial S(\mathbf{x})_2}{\partial \mathbf{x}_2} & \cdots & \frac{\partial S(\mathbf{x})_2}{\partial \mathbf{x}_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial S(\mathbf{x})_n}{\partial \mathbf{x}_1} & \frac{\partial S(\mathbf{x})_n}{\partial \mathbf{x}_2} & \cdots & \frac{\partial S(\mathbf{x})_n}{\partial \mathbf{x}_n} \end{bmatrix} \\
 &= \begin{bmatrix} S(\mathbf{x})_1 \cdot (1 - S(\mathbf{x})_1) & -S(\mathbf{x})_1 \cdot S(\mathbf{x})_2 & \cdots & -S(\mathbf{x})_1 \cdot S(\mathbf{x})_n \\ -S(\mathbf{x})_2 \cdot S(\mathbf{x})_1 & S(\mathbf{x})_2 \cdot (1 - S(\mathbf{x})_2) & \cdots & -S(\mathbf{x})_2 \cdot S(\mathbf{x})_n \\ \vdots & \vdots & \ddots & \vdots \\ -S(\mathbf{x})_n \cdot S(\mathbf{x})_1 & -S(\mathbf{x})_n \cdot S(\mathbf{x})_2 & \cdots & S(\mathbf{x})_n \cdot (1 - S(\mathbf{x})_n) \end{bmatrix} \\
 &= \begin{bmatrix} S(\mathbf{x})_1 - S(\mathbf{x})_1 \cdot S(\mathbf{x})_1 & -S(\mathbf{x})_1 \cdot S(\mathbf{x})_2 & \cdots & -S(\mathbf{x})_1 \cdot S(\mathbf{x})_n \\ -S(\mathbf{x})_2 \cdot S(\mathbf{x})_1 & S(\mathbf{x})_2 - S(\mathbf{x})_2 \cdot S(\mathbf{x})_2 & \cdots & -S(\mathbf{x})_2 \cdot S(\mathbf{x})_n \\ \vdots & \vdots & \ddots & \vdots \\ -S(\mathbf{x})_n \cdot S(\mathbf{x})_1 & -S(\mathbf{x})_n \cdot S(\mathbf{x})_2 & \cdots & S(\mathbf{x})_n - S(\mathbf{x})_n \cdot S(\mathbf{x})_n \end{bmatrix} \\
 &= \begin{bmatrix} S(\mathbf{x})_1 & 0 & \cdots & 0 \\ 0 & S(\mathbf{x})_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & S(\mathbf{x})_n \end{bmatrix} - \begin{bmatrix} S(\mathbf{x})_1 \cdot S(\mathbf{x})_1 & S(\mathbf{x})_1 \cdot S(\mathbf{x})_2 & \cdots & S(\mathbf{x})_1 \cdot S(\mathbf{x})_n \\ S(\mathbf{x})_2 \cdot S(\mathbf{x})_1 & S(\mathbf{x})_2 \cdot S(\mathbf{x})_2 & \cdots & S(\mathbf{x})_2 \cdot S(\mathbf{x})_n \\ \vdots & \vdots & \ddots & \vdots \\ S(\mathbf{x})_n \cdot S(\mathbf{x})_1 & S(\mathbf{x})_n \cdot S(\mathbf{x})_2 & \cdots & S(\mathbf{x})_n \cdot S(\mathbf{x})_n \end{bmatrix} \\
 &= \text{diag}(S(\mathbf{x})) - \begin{bmatrix} S(\mathbf{x})_1 \\ S(\mathbf{x})_2 \\ \vdots \\ S(\mathbf{x})_n \end{bmatrix} \cdot [S(\mathbf{x})_1 \quad S(\mathbf{x})_2 \quad \cdots \quad S(\mathbf{x})_n] \\
 &= \text{diag}(S(\mathbf{x})) - S(\mathbf{x}) \cdot S(\mathbf{x})^\top
 \end{aligned}$$

$$\begin{aligned}
 3. \frac{d}{dx_i} \sigma(\mathbf{x})_i &= \frac{d}{dx_i} \left(\frac{1}{1+e^{-x_i}} \right) = -\frac{-e^{-x_i}}{(1+e^{-x_i})^2} = \frac{1}{1+e^{-x_i}} \cdot \frac{e^{-x_i}+1-1}{1+e^{-x_i}} = \frac{1}{1+e^{-x_i}} \cdot \left(\frac{1+e^{-x_i}}{1+e^{-x_i}} - \frac{1}{1+e^{-x_i}} \right) \\
 &= \sigma(\mathbf{x})_i * (1 - \sigma(\mathbf{x})_i) \\
 \frac{d}{dx_j} \sigma(\mathbf{x})_i &= 0
 \end{aligned}$$

\therefore Jacobian of sigmoid :

$$\begin{aligned}
 \frac{\partial \sigma(\mathbf{x})}{\partial \mathbf{x}} &= \begin{bmatrix} \frac{\partial \sigma(\mathbf{x})_1}{\partial \mathbf{x}_1} & \frac{\partial \sigma(\mathbf{x})_1}{\partial \mathbf{x}_2} & \cdots & \frac{\partial \sigma(\mathbf{x})_1}{\partial \mathbf{x}_n} \\ \frac{\partial \sigma(\mathbf{x})_2}{\partial \mathbf{x}_1} & \frac{\partial \sigma(\mathbf{x})_2}{\partial \mathbf{x}_2} & \cdots & \frac{\partial \sigma(\mathbf{x})_2}{\partial \mathbf{x}_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \sigma(\mathbf{x})_n}{\partial \mathbf{x}_1} & \frac{\partial \sigma(\mathbf{x})_n}{\partial \mathbf{x}_2} & \cdots & \frac{\partial \sigma(\mathbf{x})_n}{\partial \mathbf{x}_n} \end{bmatrix} \\
 &= \begin{bmatrix} \sigma(\mathbf{x})_1 * (1 - \sigma(\mathbf{x})_1) & 0 & \cdots & 0 \\ 0 & \sigma(\mathbf{x})_2 * (1 - \sigma(\mathbf{x})_2) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma(\mathbf{x})_n * (1 - \sigma(\mathbf{x})_n) \end{bmatrix} \\
 &= \text{diag}(\sigma(\mathbf{x}) * (1 - \sigma(\mathbf{x})))
 \end{aligned}$$

4. $\nabla_{\mathbf{x}} L = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}}\right)^\top \nabla_{\mathbf{y}} L$

If $\mathbf{y} = \sigma(\mathbf{x})$, $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ is the Jacobian of sigmoid = $\text{diag}(\sigma(\mathbf{x}) * (1 - \sigma(\mathbf{x})))$. Since this is a diagonal matrix, it is sufficient to compute only the n diagonal elements, and multiply each of them with the corresponding element in $\nabla_{\mathbf{y}} L$. Hence, this can be done in $\mathcal{O}(n)$ time.

If $\mathbf{y} = S(\mathbf{x})$, $\left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}}\right)^\top \nabla_{\mathbf{x}} L = (\text{diag}(S(\mathbf{x})) - S(\mathbf{x}).S(\mathbf{x})^\top) \nabla_{\mathbf{x}} L = \text{diag}(S(\mathbf{x})) \nabla_{\mathbf{x}} L - S(\mathbf{x}).S(\mathbf{x})^\top \nabla_{\mathbf{x}} L$
 We just saw that $\text{diag}(S(\mathbf{x})) \nabla_{\mathbf{x}} L$ can be computed in $\mathcal{O}(n)$ time. $S(\mathbf{x})^\top \nabla_{\mathbf{x}} L$ is an inner product of two n -dimensional vectors, and so takes $\mathcal{O}(n)$ time to produce a scalar. $S(\mathbf{x}).<\text{scalar}>$ also takes $\mathcal{O}(n)$ time. Hence, $\nabla_{\mathbf{x}} L$ takes $\mathcal{O}(n)$ time even in case of $\mathbf{y} = S(\mathbf{x})$.

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

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-dimensional 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 -dimensional 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. $S(\mathbf{x} + c)_i = \frac{e^{(\mathbf{x}_i + c)}}{\sum_j e^{(\mathbf{x}_j + c)}} = \frac{e^c * e^{\mathbf{x}_i}}{\sum_j e^c * e^{\mathbf{x}_j}} = \frac{e^c * e^{\mathbf{x}_i}}{e^c * \sum_j e^{\mathbf{x}_j}} = \frac{e^{\mathbf{x}_i}}{\sum_j e^{\mathbf{x}_j}} = S(\mathbf{x})_i$
 $\implies S(\mathbf{x} + c) = S(\mathbf{x})$

2. $S_c(\mathbf{x})_i = S(c\mathbf{x}) = \frac{e^{c\mathbf{x}_i}}{\sum_j e^{c\mathbf{x}_j}} = \frac{(e^{\mathbf{x}_i})^c}{\sum_j (e^{\mathbf{x}_j})^c} \neq S(\mathbf{x})_i$ unless $c = 1$. Hence, softmax is not invariant under scalar multiplication.

$$c = 0 \implies S_c(\mathbf{x})_i = \frac{(e^{\mathbf{x}_i})^0}{\sum_j (e^{\mathbf{x}_j})^0} = \frac{1}{\sum_j 1} = \frac{1}{n}$$

Hence, when $c = 0$, the softmax values in every dimension are equal to $1/n$.

$$c \rightarrow +\infty \implies S_c(\mathbf{x})_i = \lim_{c \rightarrow +\infty} \frac{(e^{\mathbf{x}_i})^c}{\sum_j (e^{\mathbf{x}_j})^c} = \lim_{c \rightarrow +\infty} \frac{1}{\sum_j \left(\frac{e^{\mathbf{x}_j}}{e^{\mathbf{x}_i}}\right)^c}$$

If $\mathbf{x}_i = \mathbf{x}_j$, $\lim_{c \rightarrow +\infty} \left(\frac{e^{\mathbf{x}_j}}{e^{\mathbf{x}_i}}\right)^c = 1$.

If $\mathbf{x}_i > \mathbf{x}_j$, $\lim_{c \rightarrow +\infty} \left(\frac{e^{\mathbf{x}_j}}{e^{\mathbf{x}_i}}\right)^c = 0$.

If $\mathbf{x}_i < \mathbf{x}_j$, $\lim_{c \rightarrow +\infty} \left(\frac{e^{\mathbf{x}_j}}{e^{\mathbf{x}_i}}\right)^c \rightarrow +\infty$.

So, if \mathbf{x}_i is the maximum of all \mathbf{x}_j s, $\lim_{c \rightarrow +\infty} 1 / \sum_j \left(\frac{e^{\mathbf{x}_j}}{e^{\mathbf{x}_i}}\right)^c = 1 / (0 + 0 + \dots + 1 + \dots + 0) = 1$. For any other \mathbf{x}_i , $\lim_{c \rightarrow +\infty} 1 / \sum_j \left(\frac{e^{\mathbf{x}_j}}{e^{\mathbf{x}_i}}\right)^c = 1 / (\infty + \infty + \dots + 1 + \dots + \infty) = 0$.

Hence, when $c \rightarrow +\infty$, the output is 1 at the dimension with the highest value, and 0 in all other dimensions.

3. When \mathbf{x} is 2-dimensional, say $[x_1, x_2]^\top$,
 $S(\mathbf{x})_1 = S(x_1) = \frac{e^{x_1}}{e^{x_1} + e^{x_2}} = \frac{1}{1 + e^{x_2 - x_1}} = \frac{1}{1 + e^{-(x_1 - x_2)}} = \sigma(x_1 - x_2)$

$$S(\mathbf{x})_2 = S(x_2) = \frac{e^{x_2}}{e^{x_1} + e^{x_2}} = \frac{1}{e^{x_1 - x_2} + 1} = \frac{1 + e^{x_1 - x_2} - e^{x_1 - x_2}}{e^{x_1 - x_2} + 1} = 1 - \frac{e^{x_1 - x_2}}{e^{x_1 - x_2} + 1} = 1 - \frac{1}{1 + e^{-(x_1 - x_2)}}$$

$$= 1 - \sigma(x_1 - x_2)$$

Hence, if $z = x_1 - x_2$, $S(\mathbf{x}) = [\sigma(z), 1 - \sigma(z)]^\top$.

4. Let $\mathbf{x} = [x_1, x_2, \dots, x_K]^\top$.

$$S(\mathbf{x}) = [S(x_1), S(x_2), \dots, S(x_K)]^\top$$

$$S(\mathbf{x})_1 = S(x_1) = \frac{e^{x_1}}{e^{x_1} + e^{x_2} + \dots + e^{x_K}} = \frac{1}{1 + e^{(x_2 - x_1)} + \dots + e^{(x_K - x_1)}} = \frac{e^0}{e^0 + e^{(x_2 - x_1)} + \dots + e^{(x_K - x_1)}}$$

$$S(\mathbf{x})_j = S(x_j) = \frac{e^{x_j}}{e^{x_1} + \dots + e^{x_j} + \dots} = \frac{e^{(x_j - x_1)}}{e^0 + \dots + e^{(x_j - x_1)} + \dots}$$

Thus, we can see that $S(\mathbf{x}) = S([x_1, x_2, x_3, \dots, x_K]^\top) = S([0, x_2 - x_1, x_3 - x_1, \dots, x_K - x_1])$.

Hence, $S(\mathbf{x})$ can be represented using $K - 1$ parameters.

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. Recall that $\sigma(x) = \frac{1}{1 + e^{-x}}$, and $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$.

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{1 - e^{-2x}}{1 + e^{-2x}} = \frac{1 + 1 - 1 - e^{-2x}}{1 + e^{-2x}} = 2 \cdot \frac{1}{1 + e^{-2x}} - 1 = 2\sigma(2x) - 1$$

$$\implies \sigma(x) = \frac{1}{2} (\tanh(\frac{x}{2}) + 1)$$

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

\therefore There exists an equivalent network such that $y(x, \Theta', \tanh) = y(x, \Theta, \sigma)$ for all $x \in \mathbb{R}^D$.

$$\text{Here, } \Theta' = \left(\tilde{\omega}^{(1)}, (\tilde{\omega}_{k0}^{(2)}, \tilde{\omega}_{k1}^{(2)}, \tilde{\omega}_{k2}^{(2)}, \dots, \tilde{\omega}_{kM}^{(2)}) \right) = \left(\frac{\omega^{(1)}}{2}, \left(\sum_{j=1}^M \frac{\omega_{kj}^{(2)}}{2} + \omega_{k0}^{(2)}, \frac{1}{2}\omega_{k1}, \frac{1}{2}\omega_{k2}, \dots, \frac{1}{2}\omega_{kM} \right) \right).$$

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 $\mathbf{W}^{(1)} = [\mathbf{w}, \dots, \mathbf{w}]^\top$ 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.

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

2. For each $\mathbf{x}^{(i)}$ in the sample set \mathcal{S} ,

$$\begin{aligned} y(\mathbf{x}^{(i)}) &= \mathbf{W}^{(2)} \cdot \phi(\mathbf{W}^{(1)} \cdot \mathbf{x}^{(i)} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)} = [\mathbf{W}^{(2)}, \mathbf{b}^{(2)}] \cdot [\phi(\mathbf{W}^{(1)} \cdot \mathbf{x}^{(i)} + \mathbf{b}^{(1)}), 1]^\top \\ &= [\phi(\mathbf{x}^{(i)\top} \cdot \mathbf{W}^{(1)\top} + \mathbf{b}^{(1)\top}), 1] \cdot [\mathbf{W}^{(2)}, \mathbf{b}^{(2)}]^\top \end{aligned}$$

Since this is an interpolation problem, we would like $y(\mathbf{x}^{(i)}) = f(\mathbf{x}^{(i)})$. Combining all $\mathbf{x}^{(i)}$ s :

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

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

Thus, the interpolation problem can be reduced to solving a matrix equation : $\mathbf{M} \cdot \tilde{\mathbf{W}}^{(2)} = \mathbf{F}$.
 Here, \mathbf{M} is :

$$\begin{aligned} \mathbf{M} &= \begin{bmatrix} \phi(\mathbf{x}^{(1)\top} \cdot \mathbf{W}^{(1)\top} + \mathbf{b}^{(1)\top}) & 1 \\ \vdots & \vdots \\ \phi(\mathbf{x}^{(N)\top} \cdot \mathbf{W}^{(1)\top} + \mathbf{b}^{(1)\top}) & 1 \end{bmatrix} \\ &= \begin{bmatrix} \phi(\mathbf{W}^{(1)} \mathbf{x}^{(1)} + \mathbf{b}^{(1)})^\top & 1 \\ \vdots & \vdots \\ \phi(\mathbf{W}^{(1)} \mathbf{x}^{(N)} + \mathbf{b}^{(1)})^\top & 1 \end{bmatrix} \\ &= \begin{bmatrix} \phi(\mathbf{w}^\top \mathbf{x}^{(1)} + \mathbf{b}_1^{(1)}) & \phi(\mathbf{w}^\top \mathbf{x}^{(1)} + \mathbf{b}_2^{(1)}) & \cdots & \phi(\mathbf{w}^\top \mathbf{x}^{(1)} + \mathbf{b}_{N-1}^{(1)}) & 1 \\ \phi(\mathbf{w}^\top \mathbf{x}^{(2)} + \mathbf{b}_1^{(1)}) & \phi(\mathbf{w}^\top \mathbf{x}^{(2)} + \mathbf{b}_2^{(1)}) & \cdots & \phi(\mathbf{w}^\top \mathbf{x}^{(2)} + \mathbf{b}_{N-1}^{(1)}) & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \phi(\mathbf{w}^\top \mathbf{x}^{(N-1)} + \mathbf{b}_1^{(1)}) & \phi(\mathbf{w}^\top \mathbf{x}^{(N-1)} + \mathbf{b}_2^{(1)}) & \cdots & \phi(\mathbf{w}^\top \mathbf{x}^{(N-1)} + \mathbf{b}_{N-1}^{(1)}) & 1 \\ \phi(\mathbf{w}^\top \mathbf{x}^{(N)} + \mathbf{b}_1^{(1)}) & \phi(\mathbf{w}^\top \mathbf{x}^{(N)} + \mathbf{b}_2^{(1)}) & \cdots & \phi(\mathbf{w}^\top \mathbf{x}^{(N)} + \mathbf{b}_{N-1}^{(1)}) & 1 \end{bmatrix} \end{aligned}$$

3. When ϕ is ReLU, and $\mathbf{b}_j^{(1)} = -\mathbf{w}^\top \mathbf{x}^{(j)} + \epsilon$ ($\epsilon > 0$), the diagonal elements of \mathbf{M} are :

$$\phi(\mathbf{w}^\top \mathbf{x}^{(j)} + \mathbf{b}_j^{(1)}) = \phi(\epsilon) > 0 \quad (\because \epsilon > 0)$$

For \mathbf{M} to be triangular, all elements below the diagonal must be negative, so then ReLU non-linearity will make them 0.

$$\implies \mathbf{w}^\top \mathbf{x}^{(i)} + \mathbf{b}_j^{(1)} < 0 \text{ for } i > j \implies \mathbf{w}^\top \mathbf{x}^{(i)} - \mathbf{w}^\top \mathbf{x}^{(j)} + \epsilon < 0 \implies \epsilon < \mathbf{w}^\top (\mathbf{x}^{(j)} - \mathbf{x}^{(i)})$$

$\because \epsilon > 0$, we assume that $\mathbf{x}^{(i)}$ are ordered such that $\mathbf{w}^\top (\mathbf{x}^{(j)} - \mathbf{x}^{(i)}) > 0 \implies \mathbf{w}^\top \mathbf{x}^{(j)} > \mathbf{w}^\top \mathbf{x}^{(i)}$ for $i > j$.

$\therefore \epsilon$ can be any real number such that :

$$0 < \epsilon < \min_{i > j} (\mathbf{w}^\top (\mathbf{x}^{(j)} - \mathbf{x}^{(i)}))$$

Conclusion : Since \mathbf{M} is upper triangular, the linear system $\mathbf{M} \tilde{\mathbf{W}}^{(2)} = \mathbf{F}$ can be solved easily for $\tilde{\mathbf{W}}^{(2)}$. The last row of \mathbf{M} is simply $[0, \dots, 0, 1] \implies \text{LastRow}(\mathbf{M} \tilde{\mathbf{W}}^{(2)}) = \mathbf{b}^{(2)\top} = \text{LastRow}(\mathbf{F}) = f(\mathbf{x}^{(N)\top}) \implies \mathbf{b}^{(2)} = f(\mathbf{x}^{(N)})$.

Considering second last row : $[0, \dots, \phi(\mathbf{w}^\top \mathbf{x}^{(N-1)} + \mathbf{b}_{N-1}^{(1)}), 1] = [0, \dots, \phi(\epsilon), 1] \implies \phi(\epsilon) \cdot \mathbf{W}_{(N-1)}^{(2)} + \mathbf{b}^{(2)} = f(\mathbf{x}^{(N-1)}) \implies \mathbf{W}_{(N-1)}^{(2)} = (f(\mathbf{x}^{(N-1)}) - f(\mathbf{x}^{(N)})) / \phi(\epsilon)$. And so on.. Hence, $\tilde{\mathbf{W}}^{(2)}$ can be found easily when \mathbf{M} is triangular.

4. Decomposing $\mathbf{w} = \lambda \mathbf{u}$, and setting $\mathbf{b}_j^{(1)} = -\lambda \mathbf{u}^\top \mathbf{x}^{(j)}$, \mathbf{M} can be formulated as :

$$\text{The diagonal elements of } \mathbf{M} \text{ are : } \phi(\mathbf{w}^\top \mathbf{x}^{(j)} + \mathbf{b}_j^{(1)}) = \phi(\lambda \mathbf{u}^\top \mathbf{x}^{(j)} - \lambda \mathbf{u}^\top \mathbf{x}^{(j)}) = \phi(0) > 0$$

$$\text{The non-diagonal elements are } \phi(\mathbf{w}^\top \mathbf{x}^{(i)} + \mathbf{b}_j^{(1)}) = \phi(\lambda \mathbf{u}^\top \mathbf{x}^{(i)} - \lambda \mathbf{u}^\top \mathbf{x}^{(j)}) = \phi(\lambda \mathbf{u}^\top (\mathbf{x}^{(i)} - \mathbf{x}^{(j)}))$$

\therefore For elements lower than the diagonal, $i > j$. Assuming the same ordering as in 3.,

$$\mathbf{w}^\top (\mathbf{x}^{(j)} - \mathbf{x}^{(i)}) > 0 \implies \lambda \mathbf{u}^\top (\mathbf{x}^{(i)} - \mathbf{x}^{(j)}) < 0$$

$$\therefore \lambda \rightarrow +\infty \implies \text{lower-diagonal elements} \rightarrow 0.$$

$\therefore \lim_{\lambda \rightarrow +\infty} \mathbf{M}$ is triangular with non-zero diagonal elements.

Conclusion : Same as above, as $\lambda \rightarrow +\infty$.

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. *Full* convolution :

$$[1, 2, 3, 4] * [1, 0, 2] = \begin{bmatrix} [2, 0, 1, 0, 0, 0] \cdot [0, 0, 1, 2, 3, 4], \\ [2, 0, 1, 0, 0] \cdot [0, 1, 2, 3, 4], \\ [2, 0, 1, 0] \cdot [1, 2, 3, 4], \\ [0, 2, 0, 1] \cdot [1, 2, 3, 4], \\ [0, 0, 2, 0, 1] \cdot [1, 2, 3, 4, 0], \\ [0, 0, 0, 2, 0, 1] \cdot [1, 2, 3, 4, 0, 0] \end{bmatrix} = [1, 2, 5, 8, 6, 8]$$

Valid convolution : (don't pad the input)

$$[1, 2, 3, 4] * [1, 0, 2] = \begin{bmatrix} [2, 0, 1, 0] \cdot [1, 2, 3, 4], \\ [0, 2, 0, 1] \cdot [1, 2, 3, 4] \end{bmatrix} = [5, 8]$$

Same convolution : (pad input only enough to make output the same size)

$$[1, 2, 3, 4] * [1, 0, 2] = \begin{bmatrix} [2, 0, 1, 0, 0] \cdot [0, 1, 2, 3, 4], \\ [2, 0, 1, 0] \cdot [1, 2, 3, 4], \\ [0, 2, 0, 1] \cdot [1, 2, 3, 4], \\ [0, 0, 2, 0, 1] \cdot [1, 2, 3, 4, 0] \end{bmatrix} = [2, 5, 8, 6]$$

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 $64 \ 8 \times 8$ 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 $128 \ 4 \times 4$ 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. Using the formula : $o = \lfloor \frac{i+2p-k}{s} \rfloor + 1$

$$256 \times 256 \times 3 \xrightarrow[64, \ 8 \times 8, \ s=2, \ p=0]{1} (\lfloor \frac{256-8}{2} \rfloor + 1) \times (\lfloor \frac{256-8}{2} \rfloor + 1) \times 64 = 125 \times 125 \times 64$$

$$125 \times 125 \times 64 \xrightarrow[5 \times 5 \ \text{max pool}]{2} (\lfloor \frac{125-5}{5} \rfloor + 1) \times (\lfloor \frac{125-5}{5} \rfloor + 1) \times 64 = 25 \times 25 \times 64$$

$$25 \times 25 \times 64 \xrightarrow[128, \ 4 \times 4, \ s=1, \ p=1]{3} (\lfloor \frac{25+2-4}{1} \rfloor + 1) \times (\lfloor \frac{25+2-4}{1} \rfloor + 1) \times 128 = 24 \times 24 \times 128$$

2. Last layer convolves $128 \ 4 \times 4$ kernels on a 64-channel input $\implies 128 * (4 * 4 * 64) = 131072$ 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 (i.e. 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. Using the formulae : $o = \lfloor \frac{i+2p-k'}{s} \rfloor + 1$; $k' = k + (d - 1) * (k - 1)$

1. (a) $k = 8$, $s = 2$, $p = 3$, $d = 1$
(b) $k = 1$, $s = 2$, $p = 0$, $d = 7$
2. (a) $k = 4$, $s = 4$
(b) $64 \times 7 \times 7$
3. (a) $k = 2$, $s = 2$, $p = 0$, $d = 1$
(b) $k = 3$, $s = 2$, $p = 2$, $d = 2$
(c) $k = 3$, $s = 1$, $p = 1$, $d = 1$