

Due Date : March 1st (11pm), 2021

Instructions

- For all questions, show your work!
- Use LaTeX and the template we provide when writing your answers. You may reuse most of the notation shorthands, equations and/or tables. See the assignment policy on the course website for more details.
- Submit your answers electronically via Gradescope.
- **TAs for this assignment are Sébastien Lachapelle and Arian Hosseini.**

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

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

1.1 Show that the derivative of the rectified linear unit $g(x) = \max\{0, x\}$, **wherever it exists**, is equal to the Heaviside step function.

The derivative of a function at point x_0 exists if both of the right and left limits of $\frac{d}{dx}f(x)$ exist and are **equal**, i.e. :

1. $\lim_{x \rightarrow x_0^-} \frac{df(x)}{dx}$ exists
2. $\lim_{x \rightarrow x_0^+} \frac{df(x)}{dx}$ exists
3. $\lim_{x \rightarrow x_0^-} \frac{df(x)}{dx} = \lim_{x \rightarrow x_0^+} \frac{df(x)}{dx}$

Thus we consider the existence of this derivative in 3 regions :

- $x_0 < 0$, first the left limit of the derivative :

$$\lim_{\epsilon \rightarrow 0^-} \frac{g(x_0 + \epsilon) - g(x_0)}{\epsilon} \stackrel{g(x)=\max\{0,x\}}{=} \lim_{\epsilon \rightarrow 0^-} \frac{\max\{0, x_0 + \epsilon\} - \max\{0, x_0\}}{\epsilon}$$

But since $x_0 < 0$, we're in the dead region of ReLU and thus $\max\{0, x_0 + \epsilon\} = \max\{0, x_0\} = 0$, and therefore :

$$\lim_{\epsilon \rightarrow 0^-} \frac{\max\{0, x_0 + \epsilon\} - \max\{0, x_0\}}{\epsilon} = \lim_{\epsilon \rightarrow 0^-} \frac{0 - 0}{\epsilon} = 0$$

So the left limit of the derivative exists. Now the right limit :

$$\lim_{\epsilon \rightarrow 0^+} \frac{g(x_0 + \epsilon) - g(x_0)}{\epsilon} \stackrel{g(x)=\max\{0,x\}}{=} \lim_{\epsilon \rightarrow 0^+} \frac{\max\{0, x_0 + \epsilon\} - \max\{0, x_0\}}{\epsilon}$$

But again since $x_0 < 0$, we're in the dead region of ReLU and thus $\max\{0, x_0 + \epsilon\} = \max\{0, x_0\} = 0$, and therefore :

$$\lim_{\epsilon \rightarrow 0^+} \frac{\max\{0, x_0 + \epsilon\} - \max\{0, x_0\}}{\epsilon} = \lim_{\epsilon \rightarrow 0^+} \frac{0 - 0}{\epsilon} = 0$$

Hence, both the right and left limits of the derivative exist and are equal, so the derivative for any $x < 0$ **exists** and is **equal to 0**.

- $x_0 = 0$, first the left limit of the derivative :

$$\lim_{\epsilon \rightarrow 0^-} \frac{g(0 + \epsilon) - g(0)}{\epsilon} \stackrel{g(x)=\max\{0,x\}}{=} \lim_{\epsilon \rightarrow 0^-} \frac{\max\{0, 0 + \epsilon\} - \max\{0, 0\}}{\epsilon}$$

But since $x_0 = 0$, $0 + \epsilon^-$ would be in the dead region of ReLU and thus $\max\{0, 0 + \epsilon^-\} = \max\{0, 0\} = 0$, and therefore :

$$\lim_{\epsilon \rightarrow 0^-} \frac{\max\{0, 0 + \epsilon\} - \max\{0, 0\}}{\epsilon} = \lim_{\epsilon \rightarrow 0^-} \frac{0 - 0}{\epsilon} = 0$$

So the left limit of the derivative exists. Now the right limit :

$$\lim_{\epsilon \rightarrow 0^+} \frac{g(0 + \epsilon) - g(0)}{\epsilon} \stackrel{g(x)=\max\{0,x\}}{=} \lim_{\epsilon \rightarrow 0^+} \frac{\max\{0, 0 + \epsilon\} - \max\{0, 0\}}{\epsilon}$$

But now since $x_0 = 0$, $0 + \epsilon^+$ falls in the active region of ReLU and thus $\max\{0, 0 + \epsilon\} = \epsilon$, $\max\{0, 0\} = 0$, and therefore :

$$\lim_{\epsilon \rightarrow 0^+} \frac{\max\{0, 0 + \epsilon\} - \max\{0, 0\}}{\epsilon} = \lim_{\epsilon \rightarrow 0^+} \frac{\epsilon - 0}{\epsilon} = 1$$

Now both the right and left limits of the derivative exist **but are NOT equal**, so the derivative at $x = 0$ **does NOT exist**.

- $x_0 > 0$, first the left limit of the derivative :

$$\lim_{\epsilon \rightarrow 0^-} \frac{g(x_0 + \epsilon) - g(x_0)}{\epsilon} \stackrel{g(x)=\max\{0,x\}}{=} \lim_{\epsilon \rightarrow 0^-} \frac{\max\{0, x_0 + \epsilon\} - \max\{0, x_0\}}{\epsilon}$$

But since $x_0 > 0$, $x_0 + \epsilon^-$, x_0 fall in the active region of ReLU and thus $\max\{0, x_0 + \epsilon\} = x_0 + \epsilon$, $\max\{0, x_0\} = x_0$, and therefore :

$$\lim_{\epsilon \rightarrow 0^-} \frac{\max\{0, x_0 + \epsilon\} - \max\{0, x_0\}}{\epsilon} = \lim_{\epsilon \rightarrow 0^-} \frac{x_0 + \epsilon - x_0}{\epsilon} = 1$$

So the left limit of the derivative exists. Now the right limit :

$$\lim_{\epsilon \rightarrow 0^+} \frac{g(x_0 + \epsilon) - g(x_0)}{\epsilon} \stackrel{g(x)=\max\{0,x\}}{=} \lim_{\epsilon \rightarrow 0^+} \frac{\max\{0, x_0 + \epsilon\} - \max\{0, x_0\}}{\epsilon}$$

Again since $x_0 > 0$, $x_0 + \epsilon^+$, x_0 fall in the active region of ReLU and thus $\max\{0, x_0 + \epsilon\} = x_0 + \epsilon$, $\max\{0, x_0\} = x_0$, and therefore :

$$\lim_{\epsilon \rightarrow 0^+} \frac{\max\{0, x_0 + \epsilon\} - \max\{0, x_0\}}{\epsilon} = \lim_{\epsilon \rightarrow 0^+} \frac{x_0 + \epsilon - x_0}{\epsilon} = 1$$

Now both the right and left limits of the derivative exist and are equal, so the derivative at any $x > 0$ **exists** and is **equal to 1**.

So to summarize :

$$H(x) = \begin{cases} 1 & \text{if } x > 0 \\ \frac{1}{2} & \text{if } x = 0 \\ 0 & \text{if } x < 0 \end{cases} \quad \frac{dg(x)}{dx} = \begin{cases} 1 & \text{if } x > 0 \\ \text{undefined} & \text{if } x = 0 \\ 0 & \text{if } x < 0 \end{cases}$$

It's clear that wherever the derivative of $g(x)$ exists, it's equal to the Heaviside step function.

1.2 Give two alternative definitions of $g(x)$ using $H(x)$.

From the previous part we've seen that the step function relates to $g(x)$ by a derivative operation, so the natural thing that comes to mind, is integration. So the first alternative would be :

$$g(x) = \int_{-\infty}^x H(x)dx$$

It's clear that for any $x < 0$ this integral amounts to zero. For $x = 0$ again it's zero since a single point's contribution to the integral would be zero. For any $x > 0$ the result of the integral would be

$$\int_{-\infty}^x H(x)dx = \int_0^x H(x)dx = \int_0^x 1dx = x$$

So we have that :

$$\int_{-\infty}^x H(x)dx = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x = 0 \\ 0 & \text{if } x < 0 \end{cases} \quad g(x) = \max\{0, x\} = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x = 0 \\ 0 & \text{if } x < 0 \end{cases}$$

So this alternative is valid.

For the second one, since step function's dead side multiplied by anything results zero, we can think of multiplication and thus come up with $xH(x)$. For all $x < 0$ it'll be zero, for $x = 0$ it'll be zero because of the x itself, and for $x > 0$ it'll be $1 \times x = x$, so :

$$x \times H(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x = 0 \\ 0 & \text{if } x < 0 \end{cases} \quad g(x) = \max\{0, x\} = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x = 0 \\ 0 & \text{if } x < 0 \end{cases}$$

So this second alternative is also valid.

1.3 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.

To show that we define the difference $D(x) = H(x) - \sigma(x)$, evaluate the difference in different regions of x and take the limit when $k \rightarrow \infty$:

$$D(x) = H(x) - \sigma(x) = \begin{cases} 1 - \frac{1}{1+e^{-kx}} = \frac{e^{-kx}}{1+e^{-kx}} & \text{if } x > 0 \\ 0 - \frac{1}{1+e^{-kx}} = -\frac{1}{1+e^{-kx}} & \text{if } x < 0 \\ \frac{1}{2} - \frac{1}{2} = 0 & \text{if } x = 0 \end{cases}$$

$$\lim_{k \rightarrow \infty} D(x) = \begin{cases} \lim_{k \rightarrow \infty} \frac{e^{-kx}}{1+e^{-kx}} & \text{if } x > 0 \\ \lim_{k \rightarrow \infty} -\frac{1}{1+e^{-kx}} & \text{if } x < 0 \\ \lim_{k \rightarrow \infty} 0 & \text{if } x = 0 \end{cases} = \begin{cases} \frac{0}{1+0} & \text{if } x > 0 \\ -\frac{1}{1+\infty} & \text{if } x < 0 \\ 0 & \text{if } x = 0 \end{cases}$$

Thus,

$$\lim_{k \rightarrow \infty} D(x) = H(x) - \sigma(x) = \begin{cases} 0 & \text{if } x > 0 \\ 0 & \text{if } x < 0 \\ 0 & \text{if } x = 0 \end{cases}$$

So the difference asymptotically approaches zero for all $x \in \mathbb{R}$, therefore we can say that $H(x)$ is asymptotically approximated by $\sigma(x)$ (asymptotic in k).

Question 2 (3-4-4-4). Recall the definition of the softmax function : $S(\mathbf{x})_i = e^{x_i} / \sum_j e^{x_j}$.

2.1 Show that softmax is translation-invariant, that is : $S(\mathbf{x} + c) = S(\mathbf{x})$, where c is a scalar constant.

$$S(\mathbf{x} + c)_i = \frac{e^{x_i+c}}{\sum_j e^{x_j+c}} = \frac{\cancel{e^c} e_i^x}{\cancel{e^c} \sum_j e_j^x} = S(\mathbf{x})_i$$

Since this is true for all i , then we have that $S(\mathbf{x}+c) = S(\mathbf{x})$, so softmax is translation-invariant.

2.2 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} .

$$S(\mathbf{x})_1 = \frac{e_1^x}{e_1^x + e_2^x}, S(\mathbf{x})_2 = \frac{e_2^x}{e_1^x + e_2^x}$$

Or equivalently (dividing by e_1^x, e_2^x respectively) :

$$S(\mathbf{x})_1 = \frac{1}{1 + e^{-(x_1-x_2)}}, S(\mathbf{x})_2 = \frac{1}{1 + e^{x_1-x_2}}$$

Now if we define $z := x_1 - x_2$ then we have :

$$S(\mathbf{x}) = \left[\frac{1}{1 + e^{-z}}, \frac{1}{1 + e^z} \right]^\top = [\sigma(z), \sigma(-z)]^\top$$

But note that $\sigma(-z) = \frac{1}{1+e^z} = \frac{e^{-z}}{1+e^{-z}} = \frac{-1+1+e^{-z}}{1+e^{-z}} = 1 - \sigma(z)$, thus we can replace $\sigma(-z)$ by $1 - \sigma(z)$:

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

2.3 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 $i \in \{1, \dots, K - 1\}$.

In 2.1 we showed that softmax is translation-invariant, so if we let $c = -(\mathbf{x})_0$, then we know that due to translation-invariance $S(\mathbf{x}+c) = S(\mathbf{x})$, but $S(\mathbf{x}+c) = S([x_0-x_0, x_1-x_0, \dots, x_{K-1}-x_0]^\top)$ which is equal to $S([0, y_1, \dots, y_{K-1}]^\top)$, where $y_i = x_i - x_0$ are scalar functions of \mathbf{x} for $i \in \{1, \dots, K - 1\}$. Thus we've showed that $S(\mathbf{x})$ can be represented using $K - 1$ parameters.

2.4 Show that the Jacobian of the softmax function $J_{\text{softmax}}(\mathbf{x})$ can be expressed as : $\mathbf{Diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^\top$, where $\mathbf{p} = S(\mathbf{x})$.

If we have a vector valued function $f(\mathbf{x}) = [f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_K(\mathbf{x})]^\top$ like the softmax function, then the Jacobian definition is ($\mathbf{x} = [x_1, \dots, x_m]^\top$) :

$$J_f(\mathbf{x}) = \begin{pmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial x_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_K(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_K(\mathbf{x})}{\partial x_m} \end{pmatrix}$$

And

$$f(\mathbf{x}) = [\frac{e^{x_1}}{\sum_j e^{x_j}}, \frac{e^{x_2}}{\sum_j e^{x_j}}, \dots, \frac{e^{x_m}}{\sum_j e^{x_j}}]^\top$$

Therefore if we plug in the softmax function f into the Jacobian definition we get :

$$J_f(\mathbf{x}) = \begin{pmatrix} \frac{e^{x_1} \sum_j e^{x_j} - e^{x_1} e^{x_1}}{(\sum_j e^{x_j})^2} & \frac{-e^{x_1} e^{x_2}}{(\sum_j e^{x_j})^2} & \dots & \dots & \frac{-e^{x_1} e^{x_m}}{(\sum_j e^{x_j})^2} \\ \frac{-e^{x_1} e^{x_2}}{(\sum_j e^{x_j})^2} & \frac{e^{x_1} \sum_j e^{x_j} - e^{x_2} e^{x_2}}{(\sum_j e^{x_j})^2} & \frac{-e^{x_2} e^{x_3}}{(\sum_j e^{x_j})^2} & \dots & \frac{-e^{x_2} e^{x_m}}{(\sum_j e^{x_j})^2} \\ \vdots & \dots & \ddots & \dots & \vdots \\ \frac{-e^{x_1} e^{x_m}}{(\sum_j e^{x_j})^2} & \dots & \dots & \dots & \frac{e^{x_m} \sum_j e^{x_j} - e^{x_m} e^{x_m}}{(\sum_j e^{x_j})^2} \end{pmatrix}$$

As representatives, consider J_{11}, J_{12} , we'll simplify them and the rest follows the similar procedure :

$$J_{11} = \frac{e^{x_1} \sum_j e^{x_j} - e^{x_1} e^{x_1}}{(\sum_j e^{x_j})^2} = \frac{e^{x_1} \cancel{\sum_j e^{x_j}}}{(\sum_j e^{x_j})^2} - (\frac{e^{x_1}}{\sum_j e^{x_j}})^2 = S(\mathbf{x})_1 - S^2(\mathbf{x})_1$$

$$J_{12} = \frac{-e^{x_1} e^{x_2}}{(\sum_j e^{x_j})^2} = -\frac{e^{x_1}}{\sum_j e^{x_j}} \frac{e^{x_2}}{\sum_j e^{x_j}} = -S(\mathbf{x})_1 S(\mathbf{x})_2$$

Now we can simplify the Jacobian :

$$J_f(\mathbf{x}) = \begin{pmatrix} S(\mathbf{x})_1 - S^2(\mathbf{x})_1 & -S(\mathbf{x})_1 S(\mathbf{x})_2 & \dots & \dots & -S(\mathbf{x})_1 S(\mathbf{x})_m \\ -S(\mathbf{x})_1 S(\mathbf{x})_2 & S(\mathbf{x})_2 - S^2(\mathbf{x})_2 & \dots & \dots & -S(\mathbf{x})_2 S(\mathbf{x})_m \\ \vdots & \dots & \ddots & \dots & \vdots \\ -S(\mathbf{x})_1 S(\mathbf{x})_m & \dots & \dots & \dots & S(\mathbf{x})_m - S^2(\mathbf{x})_m \end{pmatrix}$$

Decomposing it into the subtraction of two matrices :

$$= \begin{pmatrix} S(\mathbf{x})_1 & & & O \\ & S(\mathbf{x})_2 & & \\ & & \ddots & \\ O & & & S(\mathbf{x})_m \end{pmatrix} - \begin{pmatrix} S(\mathbf{x})_1 \\ S(\mathbf{x})_2 \\ \vdots \\ S(\mathbf{x})_m \end{pmatrix} \begin{pmatrix} S(\mathbf{x})_1 & S(\mathbf{x})_2 & \dots & S(\mathbf{x})_m \end{pmatrix} \\ = \text{Diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^\top$$

where $\mathbf{p} = S(\mathbf{x})$.

Question 3 (2-6-6-4). Consider the differentiable functions $f : \mathbb{R}^\ell \rightarrow \mathbb{R}^m$, $g : \mathbb{R}^m \rightarrow \mathbb{R}^n$, and $h : \mathbb{R}^n \rightarrow \mathbb{R}^o$. Let $F : \mathbb{R}^\ell \rightarrow \mathbb{R}^o$ be the composition of these functions, i.e. $F = h \circ g \circ f$. For this question, denote the Jacobian matrix of F evaluated at $x \in \mathbb{R}^\ell$ as $J_F(x) \in \mathbb{R}^{o \times \ell}$, and analogously for any other function.

3.1 Using the chain rule, express $J_F(x)$ using J_f , J_g and J_h . Your answer should be in matrix form. In your expression, make sure it is clear at which point each Jacobian matrix is evaluated.

Let's first take a look at the definition of $J_F(x)$, the Jacobian evaluated at point \mathbf{x} where $F(\mathbf{x})$ is a vector valued function $F(\mathbf{x}) = [F_1(\mathbf{x}), F_2(\mathbf{x}), \dots, F_o(\mathbf{x})]^\top$:

$$J_F(\mathbf{x}) = \begin{pmatrix} \partial F_1 / \partial x_1 & \dots & \partial F_1 / \partial x_l \\ \vdots & \ddots & \vdots \\ \partial F_o / \partial x_1 & \dots & \partial F_o / \partial x_l \end{pmatrix}_{o \times l}$$

Also note that f, g, h are all vector-valued functions :

$$\mathbf{f} = \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{pmatrix}_{m \times 1} \quad \mathbf{g} = \begin{pmatrix} g_1 \\ g_2 \\ \vdots \\ g_n \end{pmatrix}_{n \times 1} \quad \mathbf{h} = \begin{pmatrix} h_1 \\ h_2 \\ \vdots \\ h_o \end{pmatrix}_{o \times 1}$$

For each of these partials $\partial F_u / \partial x_v$ we can use the chain rule. Also note that $\mathbf{F}(\mathbf{x}) = \mathbf{h}(\mathbf{g}(\mathbf{f}(\mathbf{x})))$, (any letter in bold denotes either a vector, or a vector-valued function) so for each partial we have :

$$\frac{\partial F_u}{\partial x_v} \Big|_{\mathbf{x}} = \partial h_u(\mathbf{g}(\mathbf{f}(\mathbf{x}))) / \partial x_v = \sum_{j=1}^n (\nabla_g h_u)_j \Big|_{\mathbf{g}(\mathbf{f}(\mathbf{x}))} \frac{\partial g_j}{\partial x_v} \Big|_{\mathbf{x}}$$

Or in matrix form :

$$J_F(\mathbf{x}) \Big|_{\mathbf{x}} = \underbrace{\begin{pmatrix} (\nabla_g h_1)_1 & (\nabla_g h_1)_2 & \dots & (\nabla_g h_1)_n \\ (\nabla_g h_2)_1 & (\nabla_g h_2)_2 & \dots & (\nabla_g h_2)_n \\ \vdots & \vdots & \ddots & \vdots \\ (\nabla_g h_o)_1 & \dots & \dots & (\nabla_g h_o)_n \end{pmatrix}}_{J_h(\mathbf{g}(\mathbf{f}(\mathbf{x})))} \Big|_{\mathbf{g}(\mathbf{f}(\mathbf{x}))} \times \begin{pmatrix} \partial g_1 / \partial x_1 & \partial g_1 / \partial x_2 & \dots & \partial g_1 / \partial x_l \\ \partial g_2 / \partial x_1 & \partial g_2 / \partial x_2 & \dots & \partial g_2 / \partial x_l \\ \vdots & \vdots & \ddots & \vdots \\ \partial g_n / \partial x_1 & \dots & \dots & \partial g_n / \partial x_l \end{pmatrix} \Big|_{\mathbf{x}}_{n \times l}$$

Now we should simplify the partials $\partial g_u / \partial x_v$ to get them in terms of J_g, J_f . Following a similar procedure we observe that :

$$\frac{\partial(g_u(\mathbf{f}(\mathbf{x})))}{\partial x_v} \Big|_{\mathbf{x}} = \sum_{i=1}^m (\nabla_{\mathbf{f}} g_u)_i \Big|_{\mathbf{f}(\mathbf{x})} \frac{\partial f_i}{\partial x_v} \Big|_{\mathbf{x}}$$

Or in matrix form :

$$J_g(\mathbf{x}) \Big|_{\mathbf{x}} = \underbrace{\begin{pmatrix} (\nabla_{\mathbf{f}} g_1)_1 & (\nabla_{\mathbf{f}} g_1)_2 & \dots & (\nabla_{\mathbf{f}} g_1)_m \\ (\nabla_{\mathbf{f}} g_2)_1 & (\nabla_{\mathbf{f}} g_2)_2 & \dots & (\nabla_{\mathbf{f}} g_2)_m \\ \vdots & \vdots & \ddots & \vdots \\ (\nabla_{\mathbf{f}} g_n)_1 & \dots & \dots & (\nabla_{\mathbf{f}} g_n)_m \end{pmatrix}}_{J_g(\mathbf{f}(\mathbf{x}))} \Big|_{\mathbf{f}(\mathbf{x})}^{n \times m} \times \underbrace{\begin{pmatrix} \partial f_1 / \partial x_1 & \partial f_1 / \partial x_2 & \dots & \partial f_1 / \partial x_l \\ \partial f_2 / \partial x_1 & \partial f_2 / \partial x_2 & \dots & \partial f_2 / \partial x_l \\ \vdots & \vdots & \ddots & \vdots \\ \partial f_m / \partial x_1 & \dots & \dots & \partial f_m / \partial x_l \end{pmatrix}}_{J_f(\mathbf{x})} \Big|_{\mathbf{x}}^{m \times l}$$

Plugging these back to $J_F(\mathbf{x}) \Big|_{\mathbf{x}}$ we get the following matrix multiplication :

$$J_F(\mathbf{x}) \Big|_{\mathbf{x}} = J_h(\mathbf{g}(\mathbf{f}(\mathbf{x}))) J_g(\mathbf{f}(\mathbf{x})) J_f(\mathbf{x})$$

It means that $J_h(\mathbf{g}(\mathbf{f}(\mathbf{x})))$ is the Jacobian of h evaluated at $\mathbf{g}(\mathbf{f}(\mathbf{x}))$, $J_g(\mathbf{f}(\mathbf{x}))$ is the Jacobian of g evaluated at $\mathbf{f}(\mathbf{x})$, $J_f(\mathbf{x})$ is the Jacobian of f evaluated at \mathbf{x} .

- 3.2 Provide a simple pseudo code which computes $J_F(\mathbf{x})$ using *forward mode accumulation* (Section 6.5.9 in DL book) given \mathbf{x} . You can call the functions f, g, h, J_f, J_g and J_h **only once each** (to maximize efficiency). You can call the function `matmul(\cdot, \cdot)` which performs matrix multiplication (no limit on the number of calls). Your pseudo code should also return $F(\mathbf{x})$.

In the forward mode accumulation, we traverse inside out and at each point both the function and the derivative are computed, hence it doesn't require storing any value in the intermediate steps as opposed to the reverse mode where we store the computations of the forward pass for computing the Jacobians. So we start from computing $f(\mathbf{x}), J_f(\mathbf{x})$ and move outwards, similar to multiplying from right to left (numbers below show the ordering in matrix multiplication) :

$$J_F(\mathbf{x}) \Big|_{\mathbf{x}} = J_h(\mathbf{g}(\mathbf{f}(\mathbf{x}))) \underbrace{J_g(\mathbf{f}(\mathbf{x})) J_f(\mathbf{x})}_{\substack{\#1 \\ \#2}}$$

Thus the pseudo code would be the following :

Algorithm 1: Forward Mode Accumulation

Result: Computes the Jacobian and the value of the composite function using forward mode accumulation

Select the point \mathbf{x} at which you want to compute the composite function Jacobian;
Define two temporary variables for storing evaluations and Jacobians called \mathbf{f} and \mathbf{J} ;
call $f(\mathbf{x})$, compute $J_f(\mathbf{x})$, store them temporarily in \mathbf{f}, \mathbf{J} , respectively;
Using the stored values, call $g(f(\mathbf{x}))$ and `matmul($J_g(f), J_f$)` as below :
 $\mathbf{J} \leftarrow \text{matmul}(J_g(\mathbf{f}), \mathbf{J});$
 $\mathbf{f} \leftarrow g(\mathbf{f});$
Using the stored values, call $h(g(f(\mathbf{x})))$ and `matmul($J_h(g), J_g$)` as below :
 $\mathbf{J} \leftarrow \text{matmul}(J_h(\mathbf{f}), \mathbf{J});$
 $\mathbf{f} \leftarrow h(\mathbf{f});$
Return \mathbf{f}, \mathbf{J} as the function and the Jacobian evaluation, respectively;

3.3 Provide a simple pseudo code which computes $J_F(x)$ in *reverse mode accumulation* (Section 6.5.9 in DL book) given x . You can call the functions f, g, h, J_f, J_g and J_h **only once each** (to maximize efficiency). You can call the function `matmul(\cdot, \cdot)` which performs matrix multiplication (no limit on the number of calls). Your pseudo code should also return $F(x)$.

For this mode, we should do the forward pass, but *keep* the values at intermediate steps, i.e. we do the forward pass and keep all values for $f(\mathbf{x}), g(f(\mathbf{x})), h(g(f(\mathbf{x})))$ and do the matrix multiplication from left to right :

$$J_F(x)|_{\mathbf{x}} = \underbrace{J_h(g(f(\mathbf{x}))) J_g(f(\mathbf{x}))}_{\#1} J_f(\mathbf{x})$$

$\underbrace{\hspace{10em}}_{\#2}$

Thus the pseudo code would be the following :

Algorithm 2: Reverse Mode Accumulation

Result: Computes the Jacobian and the value of the composite function using **reverse** mode accumulation

Select the point \mathbf{x} at which you want to compute the composite function Jacobian;

Define temporary variables for storing evaluations of $f(\mathbf{x}), g(f(\mathbf{x})), h(f(\mathbf{x}))$ called $\mathbf{f}, \mathbf{g}, \mathbf{h}$, and \mathbf{J} for the total Jacobian;

Do the forward pass, i.e. $\mathbf{f} \leftarrow f(\mathbf{x})$, then $\mathbf{g} \leftarrow g(\mathbf{f})$, then $\mathbf{h} \leftarrow h(\mathbf{g})$;

Using the stored values, call $J_h(g), J_g(f)$ and compute their matrix product, i.e.

$\mathbf{J} \leftarrow \text{matmul}(J_h(\mathbf{g}), J_g(\mathbf{f}))$;

Using the stored values, call $J_f(\mathbf{x})$ and compute the matrix product `matmul($J_{hg}, J_f(\mathbf{x})$)`, i.e. $\mathbf{J} \leftarrow \text{matmul}(\mathbf{J}, J_f(\mathbf{x}))$;

Return \mathbf{h}, \mathbf{J} as the function and the Jacobian evaluation, respectively;

3.4 Assume evaluating f and J_f cost $O(\ell m)$, evaluating g and J_g cost $O(mn)$ and evaluating h and J_h cost $O(no)$. What is the time complexity of your forward mode pseudo code? your reverse mode pseudo code?

Forward Mode :

Going through the pseudo code, it takes $O(\ell m)$ to evaluate $f(\mathbf{x})$ and $J_f(\mathbf{x})$, $O(mn)$ to evaluate g and J_g (using obtained $f(\mathbf{x}), J_f(\mathbf{x})$), takes $O(nm\ell)$ to compute `matmul($J_g(f), J_f$)`, takes $O(no)$ to evaluate h and J_h (using obtained $g(f(\mathbf{x})), J_g(f(\mathbf{x}))$), takes $O(on\ell)$ to compute `matmul($J_h(g), J_g$)`. So overall :

$$\text{time complexity} = O(\ell m) + O(mn) + O(nm\ell) + O(on) + O(on\ell)$$

Reverse Mode :

Going through the pseudo code, for the forward pass it takes $O(\ell m)$ to evaluate $f(\mathbf{x})$, $O(mn)$ to evaluate g (using obtained $f(\mathbf{x})$), $O(no)$ to evaluate h (using obtained $g(f(\mathbf{x}))$). For the backward pass, it takes $O(no) + O(mn)$ to evaluate $J_h(g), J_g(f)$ (using obtained $g(f(\mathbf{x})), f(\mathbf{x})$), $O(onm)$ to compute $J_{hg} = \text{matmul}(J_h(g), J_g)$, $O(\ell m)$ to evaluate J_f , $O(om\ell)$ to compute $J_F = \text{matmul}(J_{hg}, J_f)$. So overall :

$$\begin{aligned} \text{time complexity} &= O(\ell m) + O(mn) + O(no) + O(no) + O(mn) + O(onm) + O(\ell m) + O(om\ell) \\ &= O(\ell m) + O(mn) + O(no) + O(onm) + O(om\ell) \end{aligned}$$

If we compare the time complexities, we observe the difference is in terms with three elements and those are the dominating factors, for instance consider the forward mode, $O(lm)$, $O(mn)$ appear in $O(nml)$, and $O(on)$ appears in $O(onl)$, so basically these terms with two elements contribute to a factor of $O(1)$ in the terms with three elements, and thus they're absorbed.

So approximately for the forward pass we have $O(nml) + O(onl) = O(nl(m + o))$, and for the backward we have $O(om(l + n))$. In neural networks we often have very large inputs (1000×1000 images) but number of outputs are far less than the input dimension, i.e. on the order of 1000. So $o \ll l$, in that case the time complexity of the forward mode would approximate to $O(nl(m + o)) \approx O(nlm)$, whereas for the reverse mode it would be $O(om(l + n)) \approx O(olm)$ and since $o \ll l$ it justifies why we use some specific kind of reverse mode (backprop) with neural nets. However it should be noted that the memory complexity of the reverse mode is significantly larger than the forward mode, since it requires keeping all intermediate values (as can be seen by the larger number of temporary values defined and kept in alg. 2), whereas in the forward mode we can discard the values *as we pass* and evaluate the next layer outputs and Jacobians.

Question 4 (5). Compute the *full*, *valid*, and *same* convolution (with kernel flipping) for the following 1D matrices : $[5, 6, 7, 8] * [3, 0, 1]$

We first note that here kernels have to be *flipped* before being applied. Second, we know from the class that the output dimension for the result of a convolution is derived from the following :

$$o = \lfloor \frac{i + 2p - k}{s} \rfloor + 1$$

Where i is the input size, p is the padding (determined by the convolution mode), k is the kernel size, s is the stride, which is 1 for all of the cases. Now we can proceed.

For the *valid* convolution which there's no padding, the output size would be $o = \lfloor \frac{i+2p-k}{s} \rfloor + 1 = \lfloor \frac{4+2 \times 0-3}{1} \rfloor + 1 = 2$. Hence we have (note that kernels have been flipped) :

$$\begin{aligned} y[0] &= [5, 6, 7] \cdot [1, 0, 3] = 5 + 21 = 26 \\ y[1] &= [6, 7, 8] \cdot [1, 0, 3] = 6 + 24 = 30 \\ \implies y &= [26, 30]^\top \end{aligned}$$

For the *same* convolution which there's 1 zero padding to each side of the sequence, the output size would be $o = \lfloor \frac{i+2p-k}{s} \rfloor + 1 = \lfloor \frac{4+2 \times 1-3}{1} \rfloor + 1 = 4$. Hence we have (note that kernels have been flipped) :

$$\begin{aligned} y[0] &= [0, 5, 6] \cdot [1, 0, 3] = 18 \\ y[1] &= [5, 6, 7] \cdot [1, 0, 3] = 26 \\ y[2] &= [6, 7, 8] \cdot [1, 0, 3] = 30 \\ y[3] &= [7, 8, 0] \cdot [1, 0, 3] = 7 \\ \implies y &= [18, 26, 30, 7]^\top \end{aligned}$$

For the *full* convolution which there are 2 zero paddings to each side of the sequence (so that all elements of the kernel touch all elements of the input), the output size would be $o = \lfloor \frac{i+2p-k}{s} \rfloor + 1 =$

$\lfloor \frac{4+2 \times 2-3}{1} \rfloor + 1 = 6$. Hence we have (note that kernels have been flipped) :

$$\begin{aligned} y[0] &= [0, 0, 5] \cdot [1, 0, 3] = 15 \\ y[1] &= [0, 5, 6] \cdot [1, 0, 3] = 18 \\ y[2] &= [5, 6, 7] \cdot [1, 0, 3] = 26 \\ y[3] &= [6, 7, 8] \cdot [1, 0, 3] = 30 \\ y[4] &= [7, 8, 0] \cdot [1, 0, 3] = 7 \\ y[5] &= [8, 0, 0] \cdot [1, 0, 3] = 8 \\ \implies y &= [15, 18, 26, 30, 7, 8]^\top \end{aligned}$$

Question 5 (5-5). Consider a convolutional neural network. Assume the input is a colorful image of size 128×128 in the RGB representation. The first layer convolves 32 8×8 kernels with the input, using a stride of 2 and a zero-padding of 3 (three zeros on each side). The second layer downsamples the output of the first layer with a 2×2 non-overlapping max pooling. The third layer convolves 64 3×3 kernels with a stride of 1 and a zero-padding of size 1 on each border.

5.1 What is the dimensionality of the output of the last layer, i.e. the number of scalars it contains ?

we know that the output dimension for the result of a convolution is derived from the following :

$$o = \lfloor \frac{i + 2p - k}{s} \rfloor + 1$$

Where i is the input size, p is the padding (determined by the convolution mode), k is the kernel size, s is the stride. Input layer has 3 channels each of which is 128×128 . Convoluting it with 8×8 kernels with zero padding of 3 and stride 2, gives $o = \lfloor \frac{i+2p-k}{s} \rfloor + 1 = \lfloor \frac{128+2 \times 3-8}{2} \rfloor + 1 = 64$, so in the second layer, we have 32 channels (corresponding to 32 8×8 kernels) each of which being 64×64 . Then after second layer, non-overlapping max pooling of size 2×2 would divide each dimension by 2 so in the second layer we have again 32 channels (pooling doesn't affect the number of channels) each of which now being 32×32 . Then after the third layer's convolution, we have 64 channels of size $o \times o$ where $o = \lfloor \frac{i+2p-k}{s} \rfloor + 1 = \lfloor \frac{32+2 \times 1-3}{1} \rfloor + 1 = 32$. So the dimensionality of the last layer would be $64 \times 32 \times 32 = 65536$

5.2 Not including the biases, how many parameters are needed for the last layer ?

We should note that each kernel in the last layer is fully connected to *all* channels of the layer before it. So each kernel of size 3×3 corresponds to :

$$\#(\text{prev. layer channels}) \times \text{kernel size} \times \text{kernel size}$$

number of parameters. Because there are 64 of these kernels in the third layer, the total number of weights for this layer would be $64 \times 32 \times 3 \times 3 = 18432$

Question 6 (1-4-5-1-4-5). Let us use the notation $*$ and $\tilde{*}$ to denote the valid and full convolution operator **without kernel flipping**, respectively. The operations are defined as

$$\text{valid : } (\mathbf{x} * \mathbf{w})_n = \sum_{j=1}^k x_{n+j-1} w_j \quad (1)$$

$$\text{full : } (\mathbf{x} \tilde{*} \mathbf{w})_n = \sum_{j=1}^k x_{n+j-k} w_j, \quad (2)$$

where k is the size of the kernel \mathbf{w} . As a convention, the value of a vector indexed "out-of-bound" is zero, e.g. if $\mathbf{x} \in \mathbb{R}^d$, then $x_i = 0$ for $i < 1$ and $i > d$. We define the flip operator which reverse the ordering of the components of a vector, i.e. $\text{flip}(\mathbf{w})_j = w_{k-j+1}$.

Consider a convolutional network with 1-D input $\mathbf{x} \in \mathbb{R}^d$. Its first and second convolutional layers have kernel $\mathbf{w}^1 \in \mathbb{R}^{k_1}$ and $\mathbf{w}^2 \in \mathbb{R}^{k_2}$, respectively. Assume $k_1 < d$ and $k_2 < d$. The network is specified as follows :

$$\mathbf{a}^1 \leftarrow \mathbf{x} * \mathbf{w}^1 \text{ (valid convolution)} \quad (3)$$

$$\mathbf{h}^1 \leftarrow \text{ReLU}(\mathbf{a}^1) \quad (4)$$

$$\mathbf{a}^2 \leftarrow \mathbf{h}^1 * \mathbf{w}^2 \text{ (valid convolution)} \quad (5)$$

$$\mathbf{h}^2 \leftarrow \text{ReLU}(\mathbf{a}^2) \quad (6)$$

$$\dots \quad (7)$$

$$L \leftarrow \dots \quad (8)$$

where L is the loss.

6.1 What is the dimensionality of \mathbf{a}^2 ? Denote it by $|\mathbf{a}^2|$.

For valid convolution without flipping, we know that output is defined only where the kernel completely falls in the domain of the input, so we can obtain the dimensionality of $|\mathbf{a}^1|$ and then $|\mathbf{a}^2|$ because ReLU won't change the dimensionality, it just operates on all elements of a vector. So since there's no flipping, output indexes start from 1, and continue until the *last* element of the kernel touches the *last* element of the output. In other words w_{k_1} will touch x_d , so at that specific n we should have $n + k_1 - 1 = d$, So the output's last index is $n = d - k_1 + 1$. We could also achieve this by the formula we saw from the class, since the input size d , the stride is 1, padding is 0, kernel size is k_1 , then $|\mathbf{a}^1| = \lfloor \frac{d+2 \times 0 - k_1}{1} \rfloor + 1 = d - k_1 + 1$. Since ReLU doesn't affect the dimensionality, therefore $|\mathbf{h}^1| = |\mathbf{a}^1|$. We can again use the formula with the second layer to get $|\mathbf{a}^2|$. (Again stride is 1 and padding is zero)

$$\begin{aligned} |\mathbf{a}^2| &= \lfloor \frac{|\mathbf{h}^1| + 2 \times 0 - k_2}{1} \rfloor + 1 = |\mathbf{h}^1| - k_2 + 1 = d - k_1 + 1 - k_2 + 1 \\ &= d - k_1 - k_2 + 2 \end{aligned}$$

6.2 Derive $\frac{\partial a_i^2}{\partial h_n^1}$. Answer for all $i \in \{1, \dots, |\mathbf{a}^2|\}$, given a particular n .

$$\mathbf{a}_n^2 = (\mathbf{h}^1 * \mathbf{w}^2)_n = \sum_{j=1}^{k_1} \mathbf{h}_{n+j-1}^1 \mathbf{w}_j^2$$

Replacing n by i for \mathbf{a}_n^2 , we get :

$$\frac{\partial a_i^2}{\partial h_n^1} = \frac{\partial}{\partial h_n^1} \sum_{j=1}^{k_1} h_{i+j-1}^1 = \mathbb{1}_{n=i+j-1} \times \mathbf{w}_j^2$$

This indicator function is zero everywhere except where the condition holds, and for a given pair of n, i that holds when $j = n - i + 1$, so we have that $\frac{\partial a_i^2}{\partial h_n^1} = w_{n-i+1}^2$. But this amounts to zero if the subscript is not in the bounds of the kernel. To summarize :

$$\frac{\partial a_i^2}{\partial h_n^1} = \begin{cases} w_{n-i+1}^2 & \text{if } n - i + 1 \leq k_2 \\ w_{n-i+1}^2 & \text{if } n - i + 1 \geq 1 \\ 0 & \text{Otherwise} \end{cases} = \begin{cases} w_{n-i+1}^2 & n - k_2 + 1 \leq i \leq n \\ 0 & \text{Otherwise} \end{cases}$$

6.3 Show that $\nabla_{\mathbf{h}^1} L = \nabla_{\mathbf{a}^2} L \tilde{*} \text{flip}(\mathbf{w}^2)$ (full convolution). Start with

$$(\nabla_{\mathbf{h}^1} L)_n = \sum_{i=1}^{|\mathbf{a}^2|} (\nabla_{\mathbf{a}^2} L)_i \frac{\partial a_i^2}{\partial h_n^1}$$

We would first plug in the $\frac{\partial a_i^2}{\partial h_n^1}$ from the previous section :

$$(\nabla_{\mathbf{h}^1} L)_n = \sum_{i=1}^{|\mathbf{a}^2|} (\nabla_{\mathbf{a}^2} L)_i \frac{\partial a_i^2}{\partial h_n^1} = \sum_{i=1}^{|\mathbf{a}^2|} (\nabla_{\mathbf{a}^2} L)_i w_{n-i+1}^2$$

But note that from the previous section, $\frac{\partial a_i^2}{\partial h_n^1} = 0$ if $n - i + 1 \notin [1, \dots, k_2]$ or equivalently for the partial to be non-zero we need that $i \in [n - k_2 + 1, n]$. So we can change the bounds of the summation :

$$(\nabla_{\mathbf{h}^1} L)_n = \sum_{i=1}^{|\mathbf{a}^2|} (\nabla_{\mathbf{a}^2} L)_i \frac{\partial a_i^2}{\partial h_n^1} = \sum_{i=n-k_2+1}^n (\nabla_{\mathbf{a}^2} L)_i w_{n-i+1}^2$$

Next we reorder the sum and define $j = n - i + 1$, then j would run from $1 \rightarrow k_2$:

$$(\nabla_{\mathbf{h}^1} L)_n = \sum_{j=1}^{k_2} (\nabla_{\mathbf{a}^2} L)_{n-j+1} w_j^2$$

We recall that $\text{flip}(\mathbf{w})_j = w_{k-j+1}$, replacing it :

$$(\nabla_{\mathbf{h}^1} L)_n = \sum_{j=1}^{k_2} (\nabla_{\mathbf{a}^2} L)_{n-j+1} w_j^2 = \sum_{j=1}^{k_2} (\nabla_{\mathbf{a}^2} L)_{n-j+1} \text{flip}(\mathbf{w}^2)_{k_2-j+1}$$

Again changing the variables and letting $m = k_2 - j + 1$ then the sum will again run from $1 \rightarrow k_2$. Thus we have :

$$(\nabla_{\mathbf{h}^1} L)_n = \sum_{m=1}^{k_2} (\nabla_{\mathbf{a}^2} L)_{n+m-k_2} \text{flip}(\mathbf{w}^2)_m$$

This summation looks exactly the same as the definition of the full convolution : $u \tilde{*} v = \sum_{j=1}^k u_{n+j-k} v_j$, if we call $u = \nabla_{\mathbf{a}^2} L$ and $v = \text{flip}(\mathbf{w}^2)$, So this concludes the derivation and we have that $\nabla_{\mathbf{h}^1} L = \nabla_{\mathbf{a}^2} L \tilde{*} \text{flip}(\mathbf{w}^2)$.

For the following, assume the convolutions in equations (3) and (5) are **full instead of valid**.

6.4 What is the dimensionality of \mathbf{a}^2 ? Denote it by $|\mathbf{a}^2|$.

For full convolution without flipping, we know that output is defined where even a single element of the kernel touches a single element in the domain of the input, so we can obtain the dimensionality of $|\mathbf{a}^1|$ and then $|\mathbf{a}^2|$ because again ReLU won't change the dimensionality. In the full case, there's $k - 1$ padding to each side of the sequence, and then the kernel is moved across this sequence of length $d + 2 \times (k - 1)$. First index of the output occurs when w_k touches the first element of the input sequence. The last element occurs when w_1 touches the last element of the input sequence, or equivalently w_k touches x_{d+k-1} , hence the length of the output would be obtained by observing how many steps w_k takes which is from $1 \rightarrow d + k - 1$. So the output's length is $n = d + k_1 + 1$. We could also achieve this by the formula we saw from the class, since the input size d , the stride is 1, padding is $k - 1$, kernel size is k_1 , then $|\mathbf{a}^1| = \lfloor \frac{d+2 \times (k_1-1)-k_1}{1} \rfloor + 1 = d + k_1 - 1$. Since ReLU doesn't affect the dimensionality, therefore $|\mathbf{h}^1| = |\mathbf{a}^1|$. We can again use the formula with the second layer to get $|\mathbf{a}^2|$. (Again stride is 1 and padding is $k_2 - 1$ for each side)

$$\begin{aligned} |\mathbf{a}^2| &= \lfloor \frac{|\mathbf{h}^1| + 2 \times (k_2 - 1) - k_2}{1} \rfloor + 1 = |\mathbf{h}^1| + k_2 - 1 = d + k_1 - 1 + k_2 - 1 \\ &= d + k_1 + k_2 - 2 \end{aligned}$$

6.5 Derive $\frac{\partial a_i^2}{\partial h_n^1}$. Answer for all $i \in \{1, \dots, |\mathbf{a}^2|\}$, given a particular n .

$$\mathbf{a}_n^2 = (\mathbf{h}^1 \tilde{*} \mathbf{w}^2)_n = \sum_{j=1}^{k_2} h_{n+j-k_2}^1 w_j^2$$

Replacing n by i for \mathbf{a}_n^2 , we get :

$$\frac{\partial a_i^2}{\partial h_n^1} = \frac{\partial}{\partial h_n^1} \sum_{j=1}^{k_2} h_{i+j-k_2}^1 w_j^2 = \mathbb{1}_{n=i+j-k_2} \times w_j^2$$

This indicator function is zero everywhere except where the condition holds, and for a given pair of n, i that holds when $j = n - i + k_2$, so we have that $\frac{\partial a_i^2}{\partial h_n^1} = w_{n-i+k_2}^2$. But this amounts

to zero if the subscript is not in the bounds of the kernel. To summarize :

$$\frac{\partial a_i^2}{\partial h_n^1} = \begin{cases} w_{n-i+k_2}^2 & \text{if } n-i+k_2 \leq k_2 \\ w_{n-i+k_2}^2 & \text{if } n-i+k_2 \geq 1 \\ 0 & \text{Otherwise} \end{cases} = \begin{cases} w_{n-i+k_2}^2 & n \leq i \leq n+k_2-1 \\ 0 & \text{Otherwise} \end{cases}$$

6.6 Show that $\nabla_{\mathbf{h}^1} L = \nabla_{\mathbf{a}^2} L * \text{flip}(\mathbf{w}^2)$ (valid convolution). Start with

$$(\nabla_{\mathbf{h}^1} L)_n = \sum_{i=1}^{|\mathbf{a}^2|} (\nabla_{\mathbf{a}^2} L)_i \frac{\partial a_i^2}{\partial h_n^1}$$

We would first plug in the $\frac{\partial a_i^2}{\partial h_n^1}$ from the previous section :

$$(\nabla_{\mathbf{h}^1} L)_n = \sum_{i=1}^{|\mathbf{a}^2|} (\nabla_{\mathbf{a}^2} L)_i \frac{\partial a_i^2}{\partial h_n^1} = \sum_{i=1}^{|\mathbf{a}^2|} (\nabla_{\mathbf{a}^2} L)_i w_{n-i+k_2}^2$$

But note that from the previous section, $\frac{\partial a_i^2}{\partial h_n^1} = 0$ if $n-i+k_2 \notin [1, \dots, k_2]$ or equivalently for the partial to be non-zero we need that $i \in [n, n+k_2-1]$. So we can change the bounds of the summation :

$$(\nabla_{\mathbf{h}^1} L)_n = \sum_{i=1}^{|\mathbf{a}^2|} (\nabla_{\mathbf{a}^2} L)_i \frac{\partial a_i^2}{\partial h_n^1} = \sum_{i=n}^{n+k_2-1} (\nabla_{\mathbf{a}^2} L)_i w_{n-i+k_2}^2$$

Next we reorder the sum and define $j = n-i+k_2$, then j would run from $1 \rightarrow k_2$:

$$(\nabla_{\mathbf{h}^1} L)_n = \sum_{j=1}^{k_2} (\nabla_{\mathbf{a}^2} L)_{n-j+k_2} w_j^2$$

We recall that $\text{flip}(\mathbf{w})_j = w_{k-j+1}$, replacing it :

$$(\nabla_{\mathbf{h}^1} L)_n = \sum_{j=1}^{k_2} (\nabla_{\mathbf{a}^2} L)_{n-j+k_2} w_j^2 = \sum_{j=1}^{k_2} (\nabla_{\mathbf{a}^2} L)_{n-j+k_2} \text{flip}(w^2)_{k_2-j+1}$$

Again changing the variables and letting $m = k_2 - j + 1$ then the sum will again run from $1 \rightarrow k_2$. Thus we have :

$$(\nabla_{\mathbf{h}^1} L)_n = \sum_{m=1}^{k_2} (\nabla_{\mathbf{a}^2} L)_{n+m-1} \text{flip}(w^2)_m$$

This summation looks exactly the same as the definition of the valid convolution : $u * v = \sum_{j=1}^k u_{n+j-1} v_j$, if we call $u = \nabla_{\mathbf{a}^2} L$ and $v = \text{flip}(w^2)$, So this concludes the derivation and we have that $\nabla_{\mathbf{h}^1} L = \nabla_{\mathbf{a}^2} L * \text{flip}(\mathbf{w}^2)$.