

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

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.

Answer. In the following, we consider $|x| > |\epsilon|$. Evaluating two one-sided derivatives, we have :

$$\begin{aligned} \frac{d}{dx}g(x)_{x>0} &= \lim_{\substack{\epsilon \rightarrow 0 \\ x>0}} \frac{g(x+\epsilon) - g(x)}{\epsilon} \\ &= \lim_{\substack{\epsilon \rightarrow 0 \\ x>0}} \frac{\max\{0, x+\epsilon\} - \max\{0, x\}}{\epsilon} \\ &= \lim_{\substack{\epsilon \rightarrow 0 \\ x>0}} \frac{x+\epsilon - x}{\epsilon} = \lim_{\substack{\epsilon \rightarrow 0 \\ x>0}} \frac{\epsilon}{\epsilon} = 1 \\ \\ \frac{d}{dx}g(x)_{x<0} &= \lim_{\substack{\epsilon \rightarrow 0 \\ x<0}} \frac{g(x+\epsilon) - g(x)}{\epsilon} \\ &= \lim_{\substack{\epsilon \rightarrow 0 \\ x<0}} \frac{\max\{0, x+\epsilon\} - \max\{0, x\}}{\epsilon} \\ &= \lim_{\substack{\epsilon \rightarrow 0 \\ x<0}} \frac{0 - 0}{\epsilon} = 0 \end{aligned}$$

Since, the left and right derivative at $x = 0$ of the ReLU are not equal, it is not differentiable at this point. We cant then say that

$$\frac{d}{dx}g(x) = H(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x < 0 \end{cases}$$

and that it is undefined at $x = 0$.

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

Answer. ReLU can also be expressed as

$$\begin{aligned} g(x) &= xH(x) \\ g(x) &= \int_{-\infty}^x H(x)dx \end{aligned}$$

- 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.

Answer. For large k and $x > 0$,

$$\begin{aligned}\lim_{k \rightarrow \infty} \sigma(x > 0, k) &= \lim_{k \rightarrow \infty} \frac{1}{1 + e^{-k|x|}} \\ &= \frac{1}{1} = 1\end{aligned}$$

We examine the case where $x < 0$:

$$\begin{aligned}\lim_{k \rightarrow \infty} \sigma(x < 0, k) &= \lim_{k \rightarrow \infty} \frac{1}{1 + e^{k|x|}} \\ &= \frac{1}{\infty} = 0\end{aligned}$$

Finally, for $x = 0$,

$$\begin{aligned}\sigma(x = 0, k) &= \frac{1}{1 + e^0} \\ &= \frac{1}{2}\end{aligned}$$

which is the convention for the midpoint of $H(x)$. We have shown that the Heaviside function is well approximated by the parametrized sigmoid with large 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.

Answer. Defining $e^{\mathbf{x}}$ as a vector of the exponential of each element in \mathbf{x} ,

$$\begin{aligned}S(\mathbf{x} + c) &= \frac{e^{\mathbf{x}+c}}{\sum_j e^{x_j+c}} \\ &= \frac{e^c e^{\mathbf{x}}}{e^c \sum_j e^{x_j}} \\ &= \frac{e^{\mathbf{x}}}{\sum_j e^{x_j}} = S(\mathbf{x})\end{aligned}$$

- 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} .

Answer. With binary classification, a model need only output a single scalar score for one of the classes with which to calculate a class probability, since the probability of the second class

is determined exactly by the one of the first. A sigmoid applied to the output of a model $z(\mathbf{x})$ represents the probability of \mathbf{x} belonging to the "positive" class, which we call A :

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

The complementary probability that \mathbf{x} belongs to the class *not* A is then simply

$$p(\text{not } A) = 1 - \sigma(z(\mathbf{x})) = \frac{e^{-z(\mathbf{x})}}{1 + e^{-z(\mathbf{x})}} = \frac{1}{1 + e^{z(\mathbf{x})}}$$

Alternatively, suppose that the output of the model given a 2-dimensional vector \mathbf{x} was another 2-dimensional vector, with each entry specifying the score of classes A and *not* A in such a way where the first would be $z(\mathbf{x})$ and the second 0. Applying the softmax function to such a vector would result in

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

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

Answer. Let the scalar function for a class k be a linear predictor function that outputs a score for this class from K -dimensional vector \mathbf{x} with a vector of weights β_k :

$$y_k = \beta_k \cdot \mathbf{x}.$$

The probability that \mathbf{x} belongs to class k is given by $S(\mathbf{x})_k$:

$$\Pr(c = k) = S(\mathbf{x})_k = \frac{e^{\beta_k \cdot \mathbf{x}}}{\sum_j e^{\beta_j \cdot \mathbf{x}}}.$$

Since the sum of the probabilities across every class must be 1 by definition, a vector β_k is completely determined once the rest are known. Recalling the result from question 2.1 where it was shown that softmax is translation invariant, a constant vector \mathbf{C} can also be added to all coefficient vectors and still leave the equation unchanged :

$$\frac{e^{(\beta_k + \mathbf{C}) \cdot \mathbf{x}}}{\sum_j e^{(\beta_j + \mathbf{C}) \cdot \mathbf{x}}} = \frac{e^{\beta_k \cdot \mathbf{x}}}{\sum_j e^{\beta_j \cdot \mathbf{x}}}$$

It is then convenient to choose $\mathbf{C} = -\beta_0$, for example. We are then left with the set of scalar functions

$$\begin{aligned} y_0 &= (\beta_0 - \beta_0) \cdot \mathbf{x} = 0 \\ y_1 &= (\beta_1 - \beta_0) \cdot \mathbf{x} = \beta'_1 \cdot \mathbf{x} \\ &\vdots \\ y_{K-1} &= (\beta_{K-1} - \beta_0) \cdot \mathbf{x} = \beta'_{K-1} \cdot \mathbf{x}. \end{aligned}$$

The number of parameters required to represent the softmax has then been reduced to $K - 1$:

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

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

Answer. We evaluate the derivative of the i th component of $S(\mathbf{x})$ with respect to component x_j :

$$\begin{aligned} \frac{\partial}{\partial x_j} S(\mathbf{x})_i &= \frac{\partial}{\partial x_j} \left(\frac{e^{x_i}}{\sum_k e^{x_k}} \right) \\ &= \frac{1}{\sum_k e^{x_k}} \frac{\partial}{\partial x_j} e^{x_i} + e^{x_i} \left(\frac{1}{\sum_k e^{x_k}} \right)' \\ &= \frac{e^{x_i} \delta_{ij}}{\sum_k e^{x_k}} - \frac{e^{x_i}}{(\sum_k e^{x_k})^2} \sum_k \frac{\partial}{\partial x_j} e^{x_k} \\ &= \delta_{ij} S(\mathbf{x})_j - \frac{e^{x_i} e^{x_j}}{(\sum_k e^{x_k})^2} \\ &= \delta_{ij} S(\mathbf{x})_j - S(\mathbf{x})_i S(\mathbf{x})_j. \end{aligned}$$

These two terms can be interpreted as follows : the first only survives for $i = j$, and the second represents the product for every combination of i and j . Extrapolating back to matrix notation, this corresponds to

$$J_{\text{softmax}}(\mathbf{x}) = \mathbf{Diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^\top$$

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.

Answer. Using the chain rule, the Jacobian matrix of F evaluated at $x \in \mathbb{R}^\ell$ is

$$\begin{aligned} J_F(x) &= J_{h \circ g \circ f}(x) \\ &= J_h((g \circ f)(x)) J_{g \circ f}(x) \\ &= J_h(g(f(x))) J_g(f(x)) J_f(x) \end{aligned}$$

Or, in matrix form,

$$\begin{aligned} J_F(x) &= \begin{bmatrix} \frac{\partial h}{\partial g_1} & \dots & \frac{\partial h}{\partial g_n} \end{bmatrix} \begin{bmatrix} \frac{\partial g}{\partial f_1} & \dots & \frac{\partial g}{\partial f_m} \end{bmatrix} \begin{bmatrix} \frac{\partial f}{\partial x_1} & \dots & \frac{\partial f}{\partial x_\ell} \end{bmatrix} \\ &= \begin{bmatrix} \frac{\partial h_1}{\partial g_1} & \dots & \frac{\partial h_1}{\partial g_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_o}{\partial g_1} & \dots & \frac{\partial h_o}{\partial g_n} \end{bmatrix} \begin{bmatrix} \frac{\partial g_1}{\partial f_1} & \dots & \frac{\partial g_1}{\partial f_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_n}{\partial f_1} & \dots & \frac{\partial g_n}{\partial f_m} \end{bmatrix} \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_\ell} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_\ell} \end{bmatrix} \end{aligned}$$

From left to right, the Jacobian matrices are evaluated at $g(f(x)) \in \mathbb{R}^n$, $f(x) \in \mathbb{R}^m$ and $x \in \mathbb{R}^\ell$ respectively.

- 3.2 Provide a simple pseudo code which computes $J_F(x)$ using *forward 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 $F(x)$.

Answer. In the algorithm below, the outputs of functions f , g , h through the forward pass are saved in vectors \mathbf{v}_f , \mathbf{v}_g , F .

Algorithm 1 Forward mode accumulation

Require: $\mathbf{x} \in \mathbb{R}^\ell$, the input

- 1: $\mathbf{v}_f \leftarrow f(\mathbf{x})$
 - 2: $\mathbf{v}_g \leftarrow g(\mathbf{v}_f)$
 - 3: $F \leftarrow h(\mathbf{v}_g)$
 - 4: $\nabla_{\mathbf{x}} \mathbf{v}_f \leftarrow J_f(\mathbf{x})$
 - 5: $\nabla_{\mathbf{x}} \mathbf{v}_g \leftarrow \text{matmul}(J_g(\mathbf{v}_f), \nabla_{\mathbf{x}} \mathbf{v}_f)$
 - 6: $J_F(\mathbf{x}) \leftarrow \text{matmul}(J_h(\mathbf{v}_g), \nabla_{\mathbf{x}} \mathbf{v}_g)$
 - 7: **return** $J_F(\mathbf{x})$ and $F(\mathbf{x})$
-

- 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 $F(x)$.

Answer. The same state variables as in question 3.2 were employed here.

Algorithm 2 Reverse mode accumulation

Require: $\mathbf{x} \in \mathbb{R}^\ell$, the input

- 1: $\mathbf{v}_f \leftarrow f(\mathbf{x})$
 - 2: $\mathbf{v}_g \leftarrow g(\mathbf{v}_f)$
 - 3: $F \leftarrow h(\mathbf{v}_g)$
 - 4: $\nabla_{\mathbf{v}_g} F \leftarrow J_h(\mathbf{v}_g)$
 - 5: $\nabla_{\mathbf{v}_f} F \leftarrow \text{matmul}(\nabla_{\mathbf{v}_g} F, J_g(\mathbf{v}_f))$
 - 6: $J_F(\mathbf{x}) \leftarrow \text{matmul}(\nabla_{\mathbf{v}_f} F, J_f(\mathbf{x}))$
 - 7: **return** $J_F(\mathbf{x})$ and $F(\mathbf{x})$
-

- 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?

Answer. Consider two matrices $A \in \mathbb{R}^{n \times m}$ and $B \in \mathbb{R}^{m \times p}$, and suppose that we use an algorithm for which the time complexity of the matrix multiplication AB is $O(nmp)$. Going through each step of the forward accumulation algorithm, the evaluations of \mathbf{v}_f , \mathbf{v}_g , F , J_f , J_g and J_h have a time complexity of at least $O(2\ell m + 2mn + 2no) = O(\ell m + mn + no)$. At step 5, we have the

product of $J_g(\mathbf{v}_f) \in \mathbb{R}^{n \times m}$ and $\nabla_{\mathbf{x}} \mathbf{v}_f \in \mathbb{R}^{m \times \ell}$, of complexity $O(nm\ell)$, followed by the product of $J_h(\mathbf{v}_g) \in \mathbb{R}^{o \times n}$ and $\nabla_{\mathbf{x}} \mathbf{v}_g \in \mathbb{R}^{n \times \ell}$ at step 6 of complexity $O(on\ell)$. The overall time complexity of the forward accumulation algorithm is then $O(\ell m + mn + no + nm\ell + on\ell) = O(nm\ell + on\ell)$. For the reverse accumulation algorithm, the matrix multiplications differ; at step 5, the product of $\nabla_{\mathbf{v}_g} F \in \mathbb{R}^{o \times n}$ and $J_g(\mathbf{v}_f) \in \mathbb{R}^{n \times m}$ is of complexity $O(onm)$, while that of $\nabla_{\mathbf{v}_f} F \in \mathbb{R}^{o \times m}$ and $J_f(\mathbf{x}) \in \mathbb{R}^{m \times \ell}$ at step six is of $O(om\ell)$. The overall complexity for the reverse pass is then $O(\ell m + mn + no + onm + om\ell) = O(onm + om\ell)$.

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

Answer. We define the input $I = [5, 6, 7, 8]$ and the kernel $K = [3, 0, 1]$. The convolution operation, in the machine learning formalism, would then be computed by flipping the kernel from left to right and taking its dot product with overlapping parts of the input, while sliding it along by increments. The kernel is flipped to ensure the commutativity of the operation. A full convolution is one in which for a kernel of dimension k , we apply zero-padding to I so that every element is visited k times. With $k = 3$ in this case, this implies a zero-padding of 2 : $I' = [0, 0, 5, 6, 7, 8, 0, 0]$. The result of such a convolution would be :

$$\begin{aligned}(I * K)_1 &= 1 \cdot 0 + 0 \cdot 0 + 3 \cdot 5 = 15 \\(I * K)_2 &= 1 \cdot 0 + 0 \cdot 5 + 3 \cdot 6 = 18 \\(I * K)_3 &= 1 \cdot 5 + 0 \cdot 6 + 3 \cdot 7 = 26 \\(I * K)_4 &= 1 \cdot 6 + 0 \cdot 7 + 3 \cdot 8 = 30 \\(I * K)_5 &= 1 \cdot 7 + 0 \cdot 8 + 3 \cdot 0 = 7 \\(I * K)_6 &= 1 \cdot 8 + 0 \cdot 0 + 3 \cdot 0 = 8 \\ \Rightarrow (I * K) &= [15, 18, 26, 30, 7, 8] .\end{aligned}$$

In a valid convolution, the kernel must fit inside the input at all times. Therefore, there is no zero-padding, and we use I as is :

$$\begin{aligned}(I * K)_1 &= 1 \cdot 5 + 0 \cdot 6 + 3 \cdot 7 = 26 \\(I * K)_2 &= 1 \cdot 6 + 0 \cdot 7 + 3 \cdot 8 = 30 \\ \Rightarrow (I * K) &= [26, 30] .\end{aligned}$$

In a same convolution, the input is zero-padded to ensure that the output is of the same dimension. We would have $I' = [0, 5, 6, 7, 8, 0]$, and the convolution would be :

$$\begin{aligned}(I * K)_1 &= 1 \cdot 0 + 0 \cdot 5 + 3 \cdot 6 = 18 \\(I * K)_2 &= 1 \cdot 5 + 0 \cdot 6 + 3 \cdot 7 = 26 \\(I * K)_3 &= 1 \cdot 6 + 0 \cdot 7 + 3 \cdot 8 = 30 \\(I * K)_4 &= 1 \cdot 7 + 0 \cdot 8 + 3 \cdot 0 = 7 \\ \Rightarrow (I * K) &= [18, 26, 30, 7] .\end{aligned}$$

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 ?

Answer. For a square input/activation feature map and kernel, the output size of a convolution will be

$$o = \left\lfloor \frac{i + 2p - k}{s} \right\rfloor + 1, \quad (1)$$

where i is the width of the input, p the zero-padding of the kernel, k its width and s its stride. Note that a kernel always extends through the third dimension to cover all the channels. At the first layer, we have $i = 128$, $k = 8$, $s = 2$ and $p = 3$, giving an output width of 64. Stacking the output of each kernel convolution at that layer, the output tensor is of dimensions $32 \times 64 \times 64$. At the second layer, a 2×2 max-pool with no overlapping corresponds to $k = 2$, $p = 0$ and $s = 2$. With $i = 64$, we have an output width of 32. Pooling conserves the third dimension of the input, so the output of the second layer is of dimension $32 \times 32 \times 32$. At the third layer, we have $i = 32$, $k = 3$, $s = 1$ and $p = 1$, resulting in a output width of 32. From the 64 kernels at that layer, the dimension of the third layer is $64 \times 32 \times 32$. The dimensionality of such an output is then 65 536 scalars.

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

Answer. The number of parameters of a given convolution layer with square kernels of width k , disregarding the biases, is $n_{par} = k^2 pc$, where p is the number of filters at the previous layer and c is the number of filters at the current layer. The third layer has $k = 3$, $c = 64$ and $p = 32$ from the second layer. Therefore, the last layer requires 18 432 parameters.

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

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

$$\text{full : } (\mathbf{x} \tilde{*} \mathbf{w})_n = \sum_{j=1}^k x_{n+j-k} w_j, \quad (3)$$

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 :

$$\mathbf{a}^1 \leftarrow \mathbf{x} * \mathbf{w}^1 \text{ (valid convolution)} \quad (4)$$

$$\mathbf{h}^1 \leftarrow \text{ReLU}(\mathbf{a}^1) \quad (5)$$

$$\mathbf{a}^2 \leftarrow \mathbf{h}^1 * \mathbf{w}^2 \text{ (valid convolution)} \quad (6)$$

$$\mathbf{h}^2 \leftarrow \text{ReLU}(\mathbf{a}^2) \quad (7)$$

$$\dots \quad (8)$$

$$L \leftarrow \dots \quad (9)$$

where L is the loss.

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

Answer. In a valid convolution, we have no zero-padding. From the definition above, we infer a stride of $s = 1$. With expression (1), we compute the dimension of the output of the first layer to be $|\mathbf{a}^1| = \lfloor \frac{d-k_1}{1} \rfloor + 1 = d - k_1 + 1$. The ReLU activation is applied element wise, and therefore leaves the dimension unchanged. The output dimension of the second valid convolution is then $|\mathbf{a}^2| = \lfloor \frac{d-k_1+1-k_2}{1} \rfloor + 1 = d - k_1 - k_2 + 2$.

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 .

Answer. The i^{th} component of vector \mathbf{a}^2 is obtained by the valid convolution

$$a_i^2 = (\mathbf{h}^1 * \mathbf{w}^2)_i = \sum_{j=1}^{k_2} h_{i+j-1}^1 w_j^2.$$

The derivative is then

$$\begin{aligned} \frac{\partial a_i^2}{\partial h_n^1} &= \frac{\partial}{\partial h_n^1} \sum_{j=1}^{k_2} h_{i+j-1}^1 w_j^2 \\ &= \sum_{j=1}^{k_2} \frac{\partial h_{i+j-1}^1}{\partial h_n^1} w_j^2 \\ &= \delta_{n, i+j-1} w_j^2 \end{aligned}$$

which is nonzero for $n = i + j - 1 \rightarrow j = n - i + 1$. We then have

$$\frac{\partial a_i^2}{\partial h_n^1} = w_{n-i+1}^2.$$

For $i \in \{1, \dots, |\mathbf{a}^2|\}$, and with the result from 6.1, this derivative will be :

$$\begin{aligned} \frac{\partial a_1^2}{\partial h_n^1} &= w_n^2 \\ \frac{\partial a_2^2}{\partial h_n^1} &= w_{n-1}^2 \\ &\vdots \\ \frac{\partial a_{d-k_1-k_2+2}^2}{\partial h_n^1} &= w_{n-d+k_1+k_2+3}^2 \end{aligned}$$

- Do not distribute -

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

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

Answer. We begin by developing the expression in the question above, using the result from 6.2 that $\frac{\partial a_i^2}{\partial h_n^1}$ is nonzero for $i = n - j + 1$:

$$\begin{aligned} (\nabla_{\mathbf{h}^1} L)_n &= \sum_{i=1}^{|\mathbf{a}^2|} (\nabla_{\mathbf{a}^2} L)_i \frac{\partial a_i^2}{\partial h_n^1} \\ &= \sum_{i=1}^{|\mathbf{a}^2|} (\nabla_{\mathbf{a}^2} L)_{n-j+1} w_j^2 \end{aligned}$$

For a given n in the input layer h_n^1 , the output layer a_i^2 of a convolution will have derivatives w.r.t. h_n^1 for $i \in \{n - k_2 + 1, \dots, n\}$, corresponding to an index $j \in \{1, \dots, k_2\}$. We can then disregard the terms of the sum outside this range, and sum over j :

$$(\nabla_{\mathbf{h}^1} L)_n = \sum_{j=1}^{k_2} (\nabla_{\mathbf{a}^2} L)_{n-j+1} w_j^2 \quad (10)$$

We introduce an index j' so that $n - j + 1 = n + j' - k_2 \longrightarrow j' = k_2 - j + 1$, $j = k_2 - j' + 1$. The expression now sums over decreasing j' :

$$\begin{aligned} (\nabla_{\mathbf{h}^1} L)_n &= \sum_{j'=k_2}^1 (\nabla_{\mathbf{a}^2} L)_{n+j'-k_2} w_{k_2-j'+1}^2 \\ &= \sum_{j'=1}^{k_2} (\nabla_{\mathbf{a}^2} L)_{n+j'-k_2} w_{k_2-j'+1}^2 \\ &= (\nabla_{\mathbf{a}^2} L \tilde{*} \text{flip}(\mathbf{w}^2))_n \end{aligned}$$

where the order of summation has been reversed, unaffected the result.

For the following, assume the convolutions in equations (4) and (6) are **full instead of valid**.

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

Answer. In a full convolution, we have have zero-padding of $k - 1$. From expression (3), we infer a stride of $s = 1$. With expression (1), we compute the dimension of the output of the first layer to be $|\mathbf{a}^1| = \left\lfloor \frac{d+2(k_1-1)-k_1}{1} \right\rfloor + 1 = d + k_1 - 1$. The ReLU activation is applied element wise, and therefore leaves the dimension unchanged. The output dimension of the second valid convolution is then $|\mathbf{a}^2| = \left\lfloor \frac{d+k_1-1+2(k_2-1)-k_2}{1} \right\rfloor + 1 = d + k_1 + k_2 - 2$.

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 .

Answer. The i^{th} component of vector \mathbf{a}^2 is obtained by the full convolution

$$a_i^2 = (\mathbf{h}^1 \tilde{*} \mathbf{w}^2)_i = \sum_{j=1}^{k_2} h_{i+j-k_2}^1 w_j^2.$$

The derivative is then

$$\begin{aligned} \frac{\partial a_i^2}{\partial h_n^1} &= \frac{\partial}{\partial h_n^1} \sum_{j=1}^{k_2} h_{i+j-k_2}^1 w_j^2 \\ &= \sum_{j=1}^{k_2} \frac{\partial h_{i+j-k_2}^1}{\partial h_n^1} w_j^2 \\ &= \delta_{n, i+j-k_2} w_j^2 \end{aligned}$$

which is nonzero for $n = i + j - k_2 \rightarrow j = n - i + k_2$. We then have

$$\frac{\partial a_i^2}{\partial h_n^1} = w_{n-i+k_2}^2.$$

For $i \in \{1, \dots, |\mathbf{a}^2|\}$, and with the result from 6.4, this derivative will be :

$$\begin{aligned} \frac{\partial a_1^2}{\partial h_n^1} &= w_{n-1+k_2}^2 \\ \frac{\partial a_2^2}{\partial h_n^1} &= w_{n-2}^2 \\ &\vdots \\ \frac{\partial}{\partial h_n^1} a_{d+k_1+k_2-2}^2 &= w_{n-d-k_1+2}^2 \end{aligned}$$

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

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

Answer. We begin by developing the expression in the question above, using the result from 6.5 that $\frac{\partial a_i^2}{\partial h_n^1}$ is nonzero for $i = n - j + k_2$:

$$\begin{aligned} (\nabla_{\mathbf{h}^1} L)_n &= \sum_{i=1}^{|\mathbf{a}^2|} (\nabla_{\mathbf{a}^2} L)_i \frac{\partial a_i^2}{\partial h_n^1} \\ &= \sum_{i=1}^{|\mathbf{a}^2|} (\nabla_{\mathbf{a}^2} L)_{n-j+k_2} w_j^2 \end{aligned}$$

For a given n in the input layer h_n^1 , the output layer a_i^2 of a convolution will have derivatives w.r.t. h_n^1 for $i \in \{n - k_2 + 1, \dots, n\}$, corresponding to an index $j \in \{1, \dots, k_2\}$. We can then disregard the terms of the sum outside this range, and sum over j :

$$(\nabla_{\mathbf{h}^1} L)_n = \sum_{j=1}^{k_2} (\nabla_{\mathbf{a}^2} L)_{n-j+k_2} w_j^2 \quad (11)$$

We introduce an index j' so that $n - j + k_2 = n + j' - 1 \longrightarrow j' = k_2 - j + 1$, $j = k_2 - j' + 1$. The expression now sums over decreasing j' :

$$\begin{aligned} (\nabla_{\mathbf{h}^1} L)_n &= \sum_{j'=k_2}^1 (\nabla_{\mathbf{a}^2} L)_{n+j'-1} w_{k_2-j'+1}^2 \\ &= \sum_{j'=1}^{k_2} (\nabla_{\mathbf{a}^2} L)_{n+j'-1} w_{k_2-j'+1}^2 \\ &= (\nabla_{\mathbf{a}^2} L * \text{flip}(\mathbf{w}^2))_n \end{aligned}$$

where the order of summation has been reversed, unaffected the result.