

**Due Date : February 23rd (11pm), 2022**

Instructions

- For all questions, show your work!
- Use LaTeX when writing your answers. We recommend using this assignment latex code as a template. 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 Matthew Scicluna and Akram Erraqabi.**

**Question 1** (3-3-3-4-3-4-3). Consider training a standard feed-forward neural network. For the purposes of this question we are interested in a single iteration of SGD on a single training example :  $(\mathbf{x}, y)$ . We denote  $f(\mathbf{x}, \boldsymbol{\theta})$  as the output of the neural network with model parameters  $\boldsymbol{\theta}$ . Now let's say  $g$  is the output activation function and  $a(\mathbf{x}, \boldsymbol{\theta})$  is the pre-activation network output such that  $f(\mathbf{x}, \boldsymbol{\theta}) = g(a(\mathbf{x}, \boldsymbol{\theta}))$ .

- 1.1 Assuming the network's goal is to do binary classification (with the detailed structure above), what would be an appropriate activation function for the output layer, i.e. what would be an appropriate function  $g$ ?

**Answer.** A good output activation function would be the **sigmoid** :

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

- 1.2 What does the output represent under this activation function?

**Answer.** The output of the neural network using the sigmoid  $\sigma(z)$  as the output activation function gives us an output between 0 and 1 that represents the probability that the example  $\mathbf{x}$  belongs to class 1 meaning  $p(y = 1|\mathbf{x})$ .

- 1.3 Let  $L_{CE}(f(\mathbf{x}, \boldsymbol{\theta}), y)$  be cross-entropy loss, express it as a function of  $f(\mathbf{x}, \boldsymbol{\theta})$  and  $y$ .

**Answer.** The general formula for the cross entropy loss is defined as :

$$L_{CE}(f(\mathbf{x}; \boldsymbol{\theta}), y) = - \sum_i^C y_i \cdot \log(f(\mathbf{x}; \boldsymbol{\theta})_i) \quad \text{where } C \text{ is the number of class.}$$

However, the question asked to define it in function of  $f(\mathbf{x}; \boldsymbol{\theta})$  and  $y$ . Here  $y \in \{0, 1\}$  so, we can modify the equation of the general cross entropy to a version for 2 classes called binary cross entropy loss. The **binary cross entropy loss** is defined as :

$$L_{CE}(f(\mathbf{x}; \boldsymbol{\theta}), y) = -y \cdot \log(f(\mathbf{x}; \boldsymbol{\theta})) - (1 - y) \cdot \log(1 - f(\mathbf{x}; \boldsymbol{\theta}))$$

1.4 Compute the partial derivative  $\frac{\partial L_{CE}(f(\mathbf{x}, \boldsymbol{\theta}), y)}{\partial a(\mathbf{x}, \boldsymbol{\theta})}$ .

**Answer.**

$$\begin{aligned} \frac{\partial L_{CE}(f(\mathbf{x}, \boldsymbol{\theta}), y)}{\partial a(\mathbf{x}, \boldsymbol{\theta})} &= \frac{\partial L_{CE}(f(\mathbf{x}, \boldsymbol{\theta}), y)}{\partial f(\mathbf{x}, \boldsymbol{\theta})} \cdot \frac{\partial f(\mathbf{x}, \boldsymbol{\theta})}{\partial a(\mathbf{x}, \boldsymbol{\theta})} \\ &= \frac{\partial}{\partial f(\mathbf{x}, \boldsymbol{\theta})} - y \cdot \log(f(\mathbf{x}; \boldsymbol{\theta})) - (1 - y) \cdot \log(1 - f(\mathbf{x}; \boldsymbol{\theta})) \cdot \frac{\partial f(\mathbf{x}, \boldsymbol{\theta})}{\partial a(\mathbf{x}, \boldsymbol{\theta})} \\ &= \left( \frac{-y}{f(\mathbf{x}, \boldsymbol{\theta})} + \frac{1 - y}{1 - f(\mathbf{x}, \boldsymbol{\theta})} \right) \cdot \frac{\partial f(\mathbf{x}, \boldsymbol{\theta})}{\partial a(\mathbf{x}, \boldsymbol{\theta})} \end{aligned}$$

$$\begin{aligned} \frac{\partial f(\mathbf{x}, \boldsymbol{\theta})}{\partial a(\mathbf{x}, \boldsymbol{\theta})} &= \frac{\partial}{\partial a(\mathbf{x}, \boldsymbol{\theta})} \frac{1}{1 + e^{-a(\mathbf{x}, \boldsymbol{\theta})}} \\ &= \frac{\partial}{\partial a(\mathbf{x}, \boldsymbol{\theta})} (1 + e^{-a(\mathbf{x}, \boldsymbol{\theta})})^{-1} \\ &= \frac{-1}{(1 + e^{-a(\mathbf{x}, \boldsymbol{\theta})})^2} \cdot e^{-a(\mathbf{x}, \boldsymbol{\theta})} \cdot -1 \\ &= \frac{e^{-a(\mathbf{x}, \boldsymbol{\theta})}}{(1 + e^{-a(\mathbf{x}, \boldsymbol{\theta})})^2} \\ &= \frac{1 + e^{-a(\mathbf{x}, \boldsymbol{\theta})} - 1}{1 + e^{-a(\mathbf{x}, \boldsymbol{\theta})}} \cdot \frac{1}{1 + e^{-a(\mathbf{x}, \boldsymbol{\theta})}} \\ &= \left( \frac{1 + e^{-a(\mathbf{x}, \boldsymbol{\theta})}}{1 + e^{-a(\mathbf{x}, \boldsymbol{\theta})}} - \frac{1}{1 + e^{-a(\mathbf{x}, \boldsymbol{\theta})}} \right) \cdot \frac{1}{1 + e^{-a(\mathbf{x}, \boldsymbol{\theta})}} \\ &= \sigma(a(\mathbf{x}, \boldsymbol{\theta})) \cdot (1 - \sigma(a(\mathbf{x}, \boldsymbol{\theta}))) \end{aligned}$$

Therefore, as  $f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(a(\mathbf{x}, \boldsymbol{\theta}))$

$$\begin{aligned} \frac{\partial L_{CE}(f(\mathbf{x}, \boldsymbol{\theta}), y)}{\partial a(\mathbf{x}, \boldsymbol{\theta})} &= \frac{\partial L_{CE}(f(\mathbf{x}, \boldsymbol{\theta}), y)}{\partial f(\mathbf{x}, \boldsymbol{\theta})} \cdot \frac{\partial f(\mathbf{x}, \boldsymbol{\theta})}{\partial a(\mathbf{x}, \boldsymbol{\theta})} \\ &= \left( \frac{-y}{f(\mathbf{x}, \boldsymbol{\theta})} + \frac{1 - y}{1 - f(\mathbf{x}, \boldsymbol{\theta})} \right) \cdot (\sigma(a(\mathbf{x}, \boldsymbol{\theta})) \cdot (1 - \sigma(a(\mathbf{x}, \boldsymbol{\theta})))) \\ &= \frac{-y \cdot \sigma(a(\mathbf{x}, \boldsymbol{\theta})) \cdot (1 - \sigma(a(\mathbf{x}, \boldsymbol{\theta})))}{\sigma(a(\mathbf{x}, \boldsymbol{\theta}))} + \frac{(1 - y) \cdot \sigma(a(\mathbf{x}, \boldsymbol{\theta})) \cdot (1 - \sigma(a(\mathbf{x}, \boldsymbol{\theta})))}{1 - \sigma(a(\mathbf{x}, \boldsymbol{\theta}))} \\ &= (-y \cdot (1 - \sigma(a(\mathbf{x}, \boldsymbol{\theta}))) + ((1 - y) \cdot \sigma(a(\mathbf{x}, \boldsymbol{\theta})))) \\ &= \boxed{\sigma(a(\mathbf{x}, \boldsymbol{\theta})) - y} \end{aligned}$$

- 1.5 Let  $L_{MSE}(f(\mathbf{x}, \boldsymbol{\theta}), y)$  be the mean-squared error, express it as a function of  $f(\mathbf{x}, \boldsymbol{\theta})$  and  $y$  (multiplicative factors can be ignored).

**Answer.** The general formula for the mean squared error loss is defined as :

$$L_{MSE}(f(\mathbf{x}, \boldsymbol{\theta}), y) = \frac{1}{N} \cdot \sum_{i=1}^N (y^{(i)} - f(\mathbf{x}, \boldsymbol{\theta})^{(i)})^2$$

However, the question asked to define it in function of  $f(\mathbf{x}; \boldsymbol{\theta})$  and  $y$ . Here we have only one example, so we can modify the equation of the general MSE loss to a simpler version defined as :

$$L_{MSE}(f(\mathbf{x}, \boldsymbol{\theta}), y) = (y - f(\mathbf{x}, \boldsymbol{\theta}))^2$$

- 1.6 Compute the partial derivative  $\frac{\partial L_{MSE}(f(\mathbf{x}, \boldsymbol{\theta}), y)}{\partial a(\mathbf{x}, \boldsymbol{\theta})}$ .

**Answer.**

$$\begin{aligned} \frac{\partial L_{MSE}(f(\mathbf{x}, \boldsymbol{\theta}), y)}{\partial a(\mathbf{x}, \boldsymbol{\theta})} &= \frac{\partial L_{MSE}(f(\mathbf{x}, \boldsymbol{\theta}), y)}{\partial f(\mathbf{x}, \boldsymbol{\theta})} \cdot \frac{\partial f(\mathbf{x}, \boldsymbol{\theta})}{\partial a(\mathbf{x}, \boldsymbol{\theta})} \\ &= \frac{\partial}{\partial f(\mathbf{x}, \boldsymbol{\theta})} (y - f(\mathbf{x}, \boldsymbol{\theta}))^2 \cdot \frac{\partial f(\mathbf{x}, \boldsymbol{\theta})}{\partial a(\mathbf{x}, \boldsymbol{\theta})} \\ &= -2 \cdot (y - f(\mathbf{x}, \boldsymbol{\theta})) \cdot \frac{\partial f(\mathbf{x}, \boldsymbol{\theta})}{\partial a(\mathbf{x}, \boldsymbol{\theta})} \\ &= -2 \cdot (y - f(\mathbf{x}, \boldsymbol{\theta})) \cdot \sigma(a(\mathbf{x}, \boldsymbol{\theta})) \cdot (1 - \sigma(a(\mathbf{x}, \boldsymbol{\theta}))) \quad \text{using 1.4} \end{aligned}$$

- 1.7 Based on your answers to the above questions, what would be the more appropriate loss function for binary classification and why?

**Answer.**

The more appropriate loss function will be the **cross-entropy** (binary cross-entropy) because the main problem with MSE loss is that it's not punishing enough misclassifications as output of the neural network are value between 0 and 1 and the target is 0 or 1 which will lead to small distances between the prediction and the target. Cross entropy, punish well misclassifications, as show in the result of exercise 1.4. Furthermore, because binary-cross entropy loss expect to have  $p(y = 1|\mathbf{x})$  as input which is what we are doing with the sigmoid but the MSE loss does not make that assumption.

**Question 2** (4-4-5-6). Recall the definition of the softmax function :  $S(\mathbf{x})_i = e^{\mathbf{x}_i} / \sum_j e^{\mathbf{x}_j}$ .

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

**Answer.**

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

■.

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.**

$$\begin{aligned} S(\mathbf{x})_1 &= \sigma(z) = \frac{e^{\mathbf{x}_1}}{e^{\mathbf{x}_1} + e^{\mathbf{x}_2}} \\ S(\mathbf{x})_2 &= 1 - \sigma(z) = \frac{e^{\mathbf{x}_2}}{e^{\mathbf{x}_1} + e^{\mathbf{x}_2}} \\ \sigma(z) &= \frac{1}{1 + e^{-z}} = \frac{e^{\mathbf{x}_1}}{e^{\mathbf{x}_1} \cdot (1 + e^{-z})} \\ &= \frac{e^{\mathbf{x}_1}}{e^{\mathbf{x}_1} + e^{\mathbf{x}_1} \cdot e^{-z}} = \frac{e^{\mathbf{x}_1}}{e^{\mathbf{x}_1} + e^{\mathbf{x}_1 - z}} \\ &= \frac{e^{\mathbf{x}_1}}{e^{\mathbf{x}_1} + e^{\mathbf{x}_1 - \mathbf{x}_1 + \mathbf{x}_2}} = \frac{e^{\mathbf{x}_1}}{e^{\mathbf{x}_1} + e^{\mathbf{x}_2}} \end{aligned}$$

So,

$$\mathbf{x}_2 = \mathbf{x}_1 - z \Leftrightarrow z = \mathbf{x}_1 - \mathbf{x}_2$$

Therefore, using such a  $z$  :

$$\begin{aligned} 1 - \sigma(z) &= 1 - \frac{1}{1 + e^{-z}} = \frac{1 + e^{-z}}{1 + e^{-z}} - \frac{1}{1 + e^{-z}} \\ &= \frac{1 + e^{-z} - 1}{1 + e^{-z}} = \frac{e^{-z}}{1 + e^{-z}} = \frac{e^{-\mathbf{x}_1 + \mathbf{x}_2}}{1 + e^{-\mathbf{x}_1 + \mathbf{x}_2}} \\ &= \frac{e^{-\mathbf{x}_1 + \mathbf{x}_2}}{1 + e^{-\mathbf{x}_1 + \mathbf{x}_2}} \cdot \frac{e^{\mathbf{x}_1}}{e^{\mathbf{x}_1}} = \frac{e^{\mathbf{x}_2}}{e^{\mathbf{x}_1} + e^{\mathbf{x}_2}} = S(\mathbf{x})_2 \end{aligned}$$

We have been showed that  $S(\mathbf{x})$  can be reparameterized using sigmoid function, i.e.  $S(\mathbf{x}) = [\sigma(z), 1 - \sigma(z)]^\top$  where  $z = \mathbf{x}_1 - \mathbf{x}_2$ .

■.

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.**

First, I should recall that :

$$\mathbf{x} \in \mathbb{R}^K \quad \text{s.t. } K \geq 2$$

$$S(\mathbf{x})_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

Then, as  $S(\mathbf{x})$  is translation invariant :

$$\begin{aligned} S(\mathbf{x}) &= S(\mathbf{x} - \mathbf{x}_1) = S([0, \mathbf{x}_2 - \mathbf{x}_1, \mathbf{x}_3 - \mathbf{x}_1, \dots, \mathbf{x}_K - \mathbf{x}_1]^\top) \\ &= S([0, y_1, y_2, \dots, y_{K-1}]^\top) \quad \text{s.t. } y_i = \mathbf{x}_i - \mathbf{x}_1 \text{ with } 1 \leq i \leq K-1 \end{aligned}$$

■.

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

**Answer.**

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \frac{\partial S(\mathbf{x})}{\partial \mathbf{x}} \in \mathbb{R}^{n \times n}$$

$$\frac{\partial S(\mathbf{x})}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial S(\mathbf{x})_1}{\partial x_1} & \dots & \frac{\partial S(\mathbf{x})_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial S(\mathbf{x})_n}{\partial x_1} & \dots & \frac{\partial S(\mathbf{x})_n}{\partial x_n} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial S(\mathbf{x})_i}{\partial x_j} &= \frac{\partial}{\partial x_j} \frac{e^{x_i}}{\sum_{j'} e^{x_{j'}}} \\ &= \frac{\partial}{\partial x_j} e^{x_i} \left( \sum_{j'} e^{x_{j'}} \right)^{-1} \\ &= \left( \frac{\partial}{\partial x_j} e^{x_i} \right) \cdot \left( \sum_{j'} e^{x_{j'}} \right)^{-1} + e^{x_i} \cdot \frac{\partial}{\partial x_j} \left( \sum_{j'} e^{x_{j'}} \right)^{-1} \\ &= \frac{\mathbb{1}_{i=j} e^{x_i}}{\sum_{j'} e^{x_{j'}}} - \frac{e^{x_i}}{(\sum_{j'} e^{x_{j'}})^2} \cdot \frac{\partial}{\partial x_j} \sum_{j'} e^{x_{j'}} \\ &= \frac{\mathbb{1}_{i=j} e^{x_i}}{\sum_{j'} e^{x_{j'}}} - \frac{e^{x_i}}{(\sum_{j'} e^{x_{j'}})^2} \cdot e^{x_j} \\ &= \frac{\mathbb{1}_{i=j} e^{x_i}}{\sum_{j'} e^{x_{j'}}} - \frac{e^{x_i}}{\sum_{j'} e^{x_{j'}}} \cdot \frac{e^{x_j}}{\sum_{j'} e^{x_{j'}}} = \mathbb{1}_{i=j} S(\mathbf{x})_i - S(\mathbf{x})_i S(\mathbf{x})_j \end{aligned}$$

Therefore,

$$\frac{\partial S(\mathbf{x})_i}{\partial \mathbf{x}_j} = \mathbb{1}_{i=j} S(\mathbf{x})_i - S(\mathbf{x})_i S(\mathbf{x})_j$$

$$\frac{\partial S(\mathbf{x})}{\partial \mathbf{x}} = \text{Diag}(S(\mathbf{x})) - S(\mathbf{x})S(\mathbf{x})^\top = \text{Diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^\top \in \mathbb{R}^{n \times n}$$

■.

**Question 3 (6).** 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  $\sigma$ . 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  $\Theta'$  as a function of  $\Theta$ .

**Answer.**

Firstly, the relation between  $\tanh(x)$  and  $\sigma(x)$  is :

$$\begin{aligned} \tanh(x) &= \frac{e^x - e^{-x}}{e^x + e^{-x}} \\ &= \frac{e^x - e^{-x} - 2e^{-x}}{e^x + e^{-x}} = 1 - \frac{2e^{-x}}{e^x + e^{-x}} \\ &= 1 - \frac{e^{-x} 2}{e^{-x} \cdot (e^{2x} + 1)} = 1 - \frac{2}{e^{2x} + 1} \\ &= 1 - 2 \cdot \sigma(-2x) = 1 - 2 \cdot (1 - \sigma(2x)) && \text{as } \sigma(-x) = 1 - \sigma(x) \\ &= 1 - 2 + 2 \cdot \sigma(2x) = 2 \cdot \sigma(2x) - 1 \end{aligned}$$

So,  $\tanh(x) = 2 \cdot \sigma(2x) - 1 \Leftrightarrow \sigma(x) = \frac{1}{2} \cdot (\tanh(\frac{1}{2}x) + 1)$

Therefore,

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

■.

With,

$$\begin{aligned}
 \Theta' &= (\tilde{\omega}^{(1)}, \tilde{\omega}^{(2)}) \\
 \tilde{\omega}_{ji}^{(1)} &= \frac{\omega_{ji}^{(1)}}{2}, \quad \tilde{\omega}_{j0}^{(1)} = \frac{\omega_{j0}^{(1)}}{2}, \quad \tilde{\omega}_{kj}^{(2)} = \frac{\omega_{kj}^{(2)}}{2}, \quad \tilde{\omega}_{k0}^{(2)} = \left( \sum_{j=1}^M \frac{\omega_{kj}^{(2)}}{2} + \omega_{k0}^{(2)} \right)
 \end{aligned}$$

**Question 4** (5-5-6). Consider a convolutional neural network. Assume the input is a colorful image of size  $128 \times 128$  in the RGB representation. The first layer convolves 32  $8 \times 8$  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  $2 \times 2$  non-overlapping max pooling. The third layer convolves 64  $3 \times 3$  kernels with a stride of 1 and a zero-padding of size 1 on each border.

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

**Answer.**

Reminder :

$$\begin{aligned}
 H_{\text{out}} &= \left\lfloor \frac{H_{\text{in}} + 2 \times \text{padding}[0] - \text{dilation}[0] \times (\text{kernel.size}[0] - 1) - 1}{\text{stride}[0]} + 1 \right\rfloor \\
 W_{\text{out}} &= \left\lfloor \frac{W_{\text{in}} + 2 \times \text{padding}[1] - \text{dilation}[1] \times (\text{kernel.size}[1] - 1) - 1}{\text{stride}[1]} + 1 \right\rfloor
 \end{aligned}$$

**Convolution 1 :**

$$\begin{aligned}
 H_{\text{conv1}} &= \left\lfloor \frac{128 + 2 \times 3 - 1 \times (8 - 1) - 1}{2} + 1 \right\rfloor = 64 \\
 W_{\text{conv1}} &= \left\lfloor \frac{128 + 2 \times 3 - 1 \times (8 - 1) - 1}{2} + 1 \right\rfloor = 64
 \end{aligned}$$

- Do not distribute -

**output of first convolution :  $32 \times 64 \times 64$**

**Maxpooling 1 :**

As it is non-overlapping, this is equivalent to having a string of 2.

$$H_{pool1} = \left\lfloor \frac{64 + 2 \times 0 - 1 \times (2 - 1) - 1}{2} + 1 \right\rfloor = 32$$
$$W_{pool1} = \left\lfloor \frac{64 + 2 \times 0 - 1 \times (2 - 1) - 1}{2} + 1 \right\rfloor = 32$$

**output of first maxpooling :  $32 \times 32 \times 32$**

**Convolution 2 :**

$$H_{conv2} = \left\lfloor \frac{32 + 2 \times 1 - 1 \times (3 - 1) - 1}{1} + 1 \right\rfloor = 32$$
$$W_{conv2} = \left\lfloor \frac{32 + 2 \times 1 - 1 \times (3 - 1) - 1}{1} + 1 \right\rfloor = 32$$

**output of second convolution :  $64 \times 32 \times 32$**

There is  $64 \times 32 \times 32 = 65536$  scalars in the output of the last layer

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

**Answer.**

The number of parameters of the last layer is  $32 \times 32 \times 32$  (channel, height, width) and the kernel size is  $64 \times 3 \times 3 \times 32$  (num kernels, height, width, channel). The number of parameters for the last layer without biases is  $\boxed{3 \times 3 \times 32 \times 64 = 18432}$  parameters.

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

**Answer.**

**rem :**  $\tilde{*}$  is equal to the convolution with flipped kernel.

**valid convolution result :**  $[1, 2, 3, 4] \tilde{*} [1, 0, 2] = [1, 2, 3, 4] * [2, 0, 1] = [5, 8]$

**same convolution result :**  $[1, 2, 3, 4] \tilde{*} [1, 0, 2] = [1, 2, 3, 4] * [2, 0, 1] = [2, 5, 8, 6]$

**full convolution result :**  $[1, 2, 3, 4] \tilde{*} [1, 0, 2] = [1, 2, 3, 4] * [2, 0, 1] = [1, 2, 5, 8, 6, 8]$



**Question 5** (2-5-6-2-5-6). 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.

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

**Answer.**

$$|\mathbf{a}^1| = \left\lfloor \frac{d + 2 \times 0 - 1 \times (k_1 - 1) - 1}{1} + 1 \right\rfloor = d - k_1 + 1$$

$$|\mathbf{h}^1| = |\mathbf{a}^1|$$

$$|\mathbf{a}^2| = \left\lfloor \frac{(d - k_1 + 1) + 2 \times 0 - 1 \times (k_2 - 1) - 1}{1} + 1 \right\rfloor = d - k_1 + 1 - k_2 + 1 = \boxed{d - k_1 - k_2 + 2}$$

5.2 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.**

$$\frac{\partial a_i^2}{\partial h_n^1} = \frac{\partial}{\partial h_n^1} \sum_{j=1}^{k_2} h_{i+j-1}^1 w_j^2 = \begin{cases} w_j^2 & \text{if } i + j - 1 = n \\ 0 & \text{otherwise} \end{cases}$$

where  $j = n - i + 1$  the derivative is  $w_{n-i+1}^2$

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

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

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

**Answer.**

We can infer from the previous question that  $i \in \{n - k_2 + 1, n\}$  as  $j \in \{1, \dots, k_2\}$ . Elsewhere, the derivative of  $\frac{\partial a_i^2}{\partial h_n^1}$  will be 0.

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

■.

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

**Answer.**

$$|\mathbf{a}^1| = \left\lfloor \frac{d + 2 \times (k_1 - 1) - 1 \times (k_1 - 1) - 1}{1} + 1 \right\rfloor = d - 2k_1 - 2 - k_1 + 1 = d + k_1 - 1$$

$$|\mathbf{h}^1| = |\mathbf{a}^1|$$

$$|\mathbf{a}^2| = \left\lfloor \frac{(d + k_1 - 1) + 2 \times (k_2 - 1) - 1 \times (k_2 - 1) - 1}{1} + 1 \right\rfloor = \boxed{d + k_1 + k_2 - 2}$$

5.5 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.**

$$\frac{\partial a_i^2}{\partial h_n^1} = \frac{\partial}{\partial h_n^1} \sum_{j=1}^{k_2} h_{i+j-1}^1 w_j^2 = \begin{cases} w_j^2 & \text{if } i + j - k_2 = n \\ 0 & \text{otherwise} \end{cases}$$

where  $j = n + k_2 - i$  the derivative is  $w_{n+k_2-i}^2$ .

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

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

**Answer.**

We can infer from the previous question that  $i \in \{n, n + k_2 - 1\}$  as  $j \in \{1, \dots, k_2\}$ . Elsewhere, the derivative of  $\frac{\partial a_i^2}{\partial h_n^1}$  will be 0.

Therefore,

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

■.