

### Exercise 1: Convolutional Neural Networks (10 P)

Let  $x = \{x^{(t)}\}_t$  be a multivariate time series represented as a collection of one-dimensional signals. For simplicity of the following derivations, we consider these signals to be of infinite length. We feed  $x$  as input to a convolutional layer. The forward computation in this layer is given by:

$$z_t^{(l)} = \sum_i [w^{(l)0} + x^{(l)}]_i$$

$$= \sum_i \sum_{j=-\infty}^{\infty} w_{j+1}^{(l)} \cdot x_{j+t}^{(l)} \quad \text{for all } t \text{ and } l \in \mathbb{Z}$$



It results in  $z = \{z^{(l)}\}_l$  another multivariate time series again composed of a collection of output signals also assumed to be of infinite length. Convolution filters  $w^{(l)}$  have to be learned from the data. After passing the data through the convolutional layer, the neural network output is given as  $y = f(z)$  where  $f$  is some top-layer function assumed to be differentiable. To learn the model, parameters gradients need to be computed.

(a) Express the gradient  $\partial y / \partial w$  as a function of the input  $x$  and of the gradient  $\partial y / \partial z$ .

$$\begin{aligned} \frac{\partial y}{\partial w} &= \sum_j \sum_l \frac{\partial y}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial w_{j+t}^{(l)}} = \sum_j \sum_l \frac{\partial y}{\partial z_j^{(l)}} \frac{\partial}{\partial w_{j+t}^{(l)}} \sum_{i'} \sum_{i''} w_{i''+1}^{(l)} \cdot x_{i'+t}^{(l)} \\ &= \sum_j \sum_l \frac{\partial y}{\partial z_j^{(l)}} \sum_{i'} \delta_{(i'+t)=j+t} = \sum_j \sum_l \frac{\partial y}{\partial z_j^{(l)}} \cdot x_{j+t}^{(l)} \\ &= \sum_j \left[ \frac{\partial y}{\partial z_j^{(l)}} \cdot x^{(l)} \right]_j \\ \frac{\partial y}{\partial w} &= \sum_j \frac{\partial y}{\partial z_j^{(l)}} \cdot x^{(l)} \end{aligned}$$

### Exercise 2: Recursive Neural Networks (10 + 10 P)

Consider a recursive neural network that applies some function  $\phi$  recursively from the leaves to the root of a binary parse tree. The function  $\phi : \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}^h \rightarrow \mathbb{R}^d$  takes two incoming nodes  $a_i, a_j \in \mathbb{R}^d$  and some parameter vector  $\theta \in \mathbb{R}^h$  as input, and produces an output  $a_k \in \mathbb{R}^d$ . Once we have reached the root node, we apply a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  that produces some real-valued prediction. We assume that both  $\phi$  and  $f$  are differentiable with their inputs.

Consider the sentence 'the cat sat on the mat', that is parsed as:

((the, cat), (sat, (on, (the, mat)))))

The leaf nodes of the parsing tree are represented by word embeddings  $a_{the}, a_{sat}, \dots \in \mathbb{R}^d$ .

(a) Draw the computational graph associated to the application of the recursive neural network to this sentence, and write down the set of equations, e.g.

$$\begin{aligned} a_1 &= \phi(a_{the}, a_{cat}, \theta) \\ a_2 &= \phi(a_{the}, a_{sat}, \theta) \\ a_3 &= \phi(a_{on}, a_p, \theta) \\ &\vdots \\ y &= f(a_k) \end{aligned}$$



(b) Express the total derivative  $dy/d\theta$  (taking into account all direct and indirect dependencies on the parameter  $\theta$ ) in terms of local derivatives (relating adjacent quantities in the computational graph).

(Hint: you can use for this the chain rule for derivatives that states that for some function  $h(g_1(t), \dots, g_V(t))$  we have:

$$\frac{dh}{dt} = \frac{\partial h}{\partial g_1} \cdot \frac{dg_1}{dt} + \dots + \frac{\partial h}{\partial g_V} \cdot \frac{dg_V}{dt}$$

where  $d(\cdot)$  and  $\partial(\cdot)$  denote the total and partial derivatives respectively.)

$$\begin{aligned} \frac{dy}{d\theta} &= \frac{\partial y}{\partial a_1} \left[ \frac{\partial a_1}{\partial \theta} + \frac{\partial a_1}{\partial a_2} \frac{\partial a_2}{\partial \theta} + \frac{\partial a_1}{\partial a_3} \frac{da_3}{d\theta} \right] \\ &= \frac{\partial a_1}{\partial \theta} + \frac{\partial a_1}{\partial a_2} \cdot \frac{da_2}{d\theta} \\ &= \frac{\partial a_1}{\partial \theta} + \frac{\partial a_1}{\partial a_2} \frac{\partial a_2}{\partial \theta} \end{aligned}$$

### Exercise 3: Graph Neural Networks (10 + 10 P)

Graph neural networks are a fairly flexible class of neural networks that can sometimes be seen as a generalization of convolutional neural networks and recursive neural networks. We would like to study the equivalence of graph neural network with other neural networks. We will consider in this exercise graph neural network that map two consecutive layers using the equation:

$$H_{t+1} = \rho(AH_t W)$$

where  $\rho$  denotes the ReLU function. The adjacency matrix  $A$  is of size  $N \times N$  with value 1 when there is a connection, and 0 otherwise. The matrix of parameters  $W$  is of size  $d \times d$ , where  $d$  is the number of dimensions used to represent each node. We also assume that the initial state  $H_0$  is a matrix filled with ones.

(a) Consider first the following undirected graph:



Because the graph is unlabeled, the nodes can only be distinguished by their connectivity to other nodes. Consider the case where the graph neural network has depth 2. Depict the multiple trees formed by viewing the graph neural network as a collection of recursive neural networks.



(b) Consider now the following infinite lattice graph.



which is like in the previous example undirected and unlabeled. We consider again the case where the graph neural network has depth 2. Show that the latter can be implemented as a 2D convolutional neural network with four convolution layers and two ReLU layer, i.e. give the sequence of layers and their parameters.

$$2 \times \begin{cases} \text{conv}_{2 \times 2} : A_1^{(2)} = A_{1 \times 1} \cdot \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \\ \text{conv}_{1 \times 1} : A_1^{(3)} = [A_1^{(2)}] \\ \rho : \text{ReLU activation} \end{cases}$$