

## Exercise Sheet 2

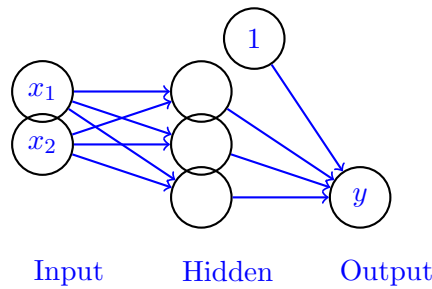
### Exercise 1: Drawing and Simulating a Neural Network (5 + 15 = 20 P)

You are given a neural network with two input variables  $x_1$  and  $x_2$  and that produces the output:

$$y = 1 - \rho(x_1) - \rho(x_2) + \rho(x_1 + x_2) \quad (1)$$

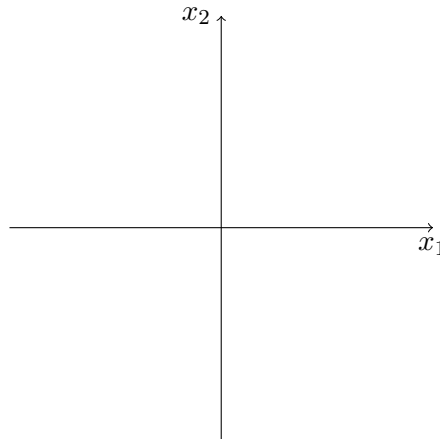
where  $\rho(z) = \max(0, z)$  is the ReLU nonlinear activation function. The network predicts that the instance fed as input is of class 1 if  $y < 0$  and of class 2 if  $y > 0$ .

(a) *Draw* the neural network associated to Eq. (1), i.e. draw the input variables, the neurons, and the way they are connected. Indicate for each connection its associated weight, and for each neuron its bias and the type of nonlinearity.



(b) *Draw* the decision boundary implemented by this neural network in a 2D coordinate system.

*Hint: You can break down the problem into multiple sub-problems, specifically (i) observing that the neural network function is piecewise linear, and identifying its linear pieces, and (ii) for each linear piece, solving the equations  $y < 0$  and  $y > 0$ . You can also observe that  $y(x_1, x_2) = y(x_2, x_1)$  for this neural network in order to reduce the number of calculations.*



Solution sketch: Consider all linear pieces, given by the domain:

$$(x_1 \geq 0) \wedge (x_2 \geq 0) \wedge (x_1 + x_2 \geq 0) \quad (1)$$

$$(x_1 \geq 0) \wedge (x_2 \geq 0) \wedge (x_1 + x_2 < 0) \quad (2)$$

$$(x_1 \geq 0) \wedge (x_2 < 0) \wedge (x_1 + x_2 \geq 0) \quad (3)$$

$$(x_1 \geq 0) \wedge (x_2 < 0) \wedge (x_1 + x_2 < 0) \quad (4)$$

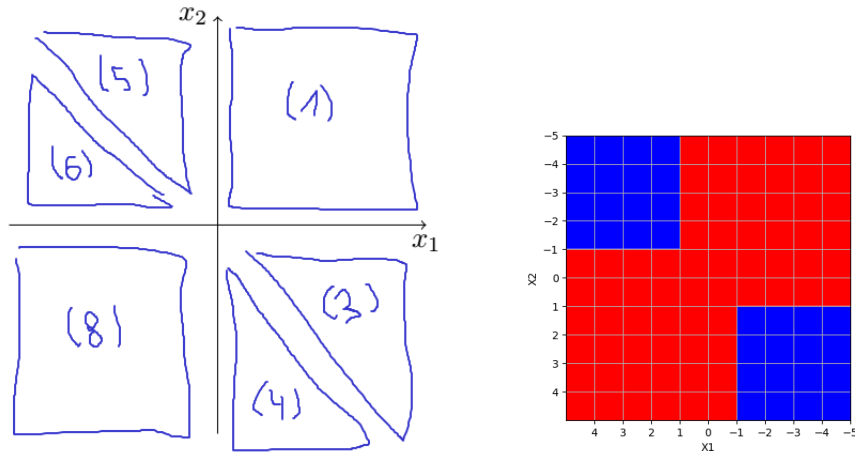
$$(x_1 < 0) \wedge (x_2 \geq 0) \wedge (x_1 + x_2 \geq 0) \quad (5)$$

$$(x_1 < 0) \wedge (x_2 \geq 0) \wedge (x_1 + x_2 < 0) \quad (6)$$

$$(x_1 < 0) \wedge (x_2 < 0) \wedge (x_1 + x_2 \geq 0) \quad (7)$$

$$(x_1 < 0) \wedge (x_2 < 0) \wedge (x_1 + x_2 < 0) \quad (8)$$

illustrated in the following image:



and derive the decision boundary in each linear piece. Note that Eq. (2) and (7) domains are empty. Also, the decision boundary in the piece (5) and (6) can easily be inferred from (3) and (4) by using the fact that the decision boundary is symmetric.

## Exercise 2: Backpropagation in the Error Function (10 + 10 + 10 = 30 P)

Let  $y \in \mathbb{R}$  be the output of a neural network for some data point  $\mathbf{x} \in \mathbb{R}^d$ . The true target value that the network should predict is given by  $t$ . We define the error function to be

$$E = \log \cosh(y - t).$$

This error function is commonly used when training neural network on real-valued prediction tasks (e.g. predicting the energy of a physical system, the price of an object, etc.). We would like to extract the gradient of this error function so that a neural network using it can be learned. Unless stated otherwise we use  $\log$  to denote the natural logarithm.

(a) Using the chain rule for derivatives, *compute* the gradient of  $E$  with respect to the output  $y$  of the neural network. Show each step of your derivation.

*Hint: The derivative of  $\cosh(z)$  is  $\sinh(z)$ . You can further use the identity  $\tanh(z) = \sinh(z)/\cosh(z)$ .*

Observe that error can be decomposed as

$$z = y - t \tag{2}$$

$$E = \log c = \log \cosh(z) \tag{3}$$

Application of the chain rule gives

$$\frac{\partial E}{\partial y} = \frac{\partial E}{\partial c} \frac{\partial \cosh(z)}{\partial z} \frac{\partial z}{\partial y} \tag{4}$$

$$= \frac{1}{\cosh(z)} \cdot \sinh(z) \cdot 1 \tag{5}$$

$$= \tanh(z) \tag{6}$$

$$= \tanh(y - t) \tag{7}$$

(b) Assume we have a dataset composed of neural network inputs  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$  and associated targets  $t_1, \dots, t_N \in \mathbb{R}$ . We denote by  $y_1, \dots, y_N \in \mathbb{R}$  the predictions of the neural network for these points. We define the error

$$E = \frac{1}{N} \sum_{k=1}^N E_k \quad \text{with} \quad E_k = \log \cosh(y_k(\mathbf{x}_k, \mathbf{w}) - t_k)$$

State the chain rule for transmitting the gradient from the output of the neural network to the model parameters.

$$\frac{\partial E}{\partial \mathbf{w}} = \sum_k \frac{\partial E}{\partial E_k} \frac{\partial E_k}{\partial y_k} \frac{\partial y_k}{\partial \mathbf{w}}$$

(c) Assume that  $y_k(\mathbf{x}_k, \mathbf{w}) = \sum_{i=1}^d w_i x_i^{(k)}$  where  $x_i^{(k)}$  denotes the  $i$ th element of the vector  $\mathbf{x}_k$ . Compute the gradient of the error function w.r.t. the parameter  $w_i$ , i.e. compute  $\partial E / \partial w_i$ .

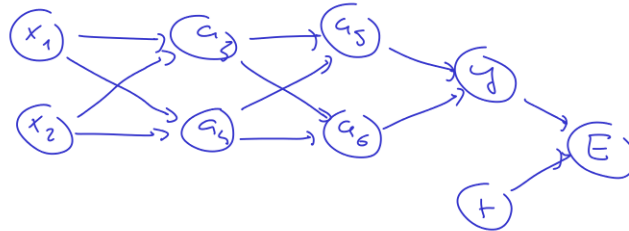
$$\begin{aligned} \frac{\partial E}{\partial w_i} &= \sum_{k=1}^N \frac{\partial E}{\partial E_k} \frac{\partial E_k}{\partial y_k} \frac{\partial y_k}{\partial w_i} \\ &= \sum_{k=1}^N \frac{1}{N} \cdot \tanh(y_k - t_k) \cdot \frac{\partial}{\partial w_i} \left( \sum_{i'=1}^d w_{i'} x_{i'}^{(k)} \right) \\ &= \frac{1}{N} \sum_{k=1}^N \tanh(y_k - t_k) x_i^{(k)} \end{aligned}$$

### Exercise 3: Backpropagation in a Multilayer Network (5 + 15 = 20 P)

We consider a neural network that takes two inputs  $x_1$  and  $x_2$  and produces an output  $y$  based on the following set of computations:

$$\begin{array}{lll} z_3 = x_1 \cdot w_{13} + x_2 \cdot w_{23} & z_5 = a_3 \cdot w_{35} + a_4 \cdot w_{45} & y = a_5 \cdot w_{57} + a_6 \cdot w_{67} \\ a_3 = \tanh(z_3) & a_5 = \tanh(z_5) & E = \log \cosh(y - t) \\ z_4 = x_1 \cdot w_{14} + x_2 \cdot w_{24} & z_6 = a_3 \cdot w_{36} + a_4 \cdot w_{46} & \\ a_4 = \tanh(z_4) & a_6 = \tanh(z_6) & \end{array}$$

(a) Draw the neural network graph associated to this set of computations.



(b) Write the set of backward computations that leads to the evaluation of the partial derivative  $\partial E / \partial w_{13}$ . Your answer should avoid redundant computations. Hint:  $\tanh'(t) = 1 - (\tanh(t))^2$ .

We see that there is a sum over the second layer

Chain rule:

$$\begin{aligned}
\frac{\partial E}{\partial w_{13}} &= \frac{\partial E}{\partial y} \left( \frac{\partial y}{\partial a_5} \frac{\partial a_5}{\partial z_5} \frac{\partial z_5}{\partial a_3} + \frac{\partial y}{\partial a_6} \frac{\partial a_6}{\partial z_6} \frac{\partial z_6}{\partial a_3} \right) \frac{\partial a_3}{\partial z_3} \frac{\partial z_3}{\partial w_{13}} \\
\frac{\partial E}{\partial y} &= \tanh(y - t) \\
\frac{\partial y}{\partial a_5} &= w_{57} & \frac{\partial y}{\partial a_6} &= w_{67} \\
\frac{\partial a_5}{\partial z_5} &= (1 - a_5^2) & \frac{\partial a_6}{\partial z_6} &= (1 - a_6^2) \\
\frac{\partial z_5}{\partial a_3} &= w_{35} & \frac{\partial z_6}{\partial a_3} &= w_{36} \\
\frac{\partial a_3}{\partial z_3} &= (1 - a_3^2) \\
\frac{\partial z_3}{\partial w_{31}} &= x_1 \\
\frac{\partial E}{\partial w_{13}} &= \tanh(y - t) \left( \sum_{i=5}^6 w_{iy} (1 - a_i^2) w_{i3} \right) (1 - a_3^2) x_1
\end{aligned}$$

**Exercise 4: Backpropagation with Shared Parameters (5 + 10 + 5 + 10 = 30 P)**

Let  $x_1, x_2$  be two observed variables. Consider the two-layer neural network that takes these two variables as input and builds the prediction  $y$  by computing iteratively:

$$z_3 = x_1 w_{13}, \quad z_4 = x_2 w_{24}, \quad a_3 = 0.5 z_3^2, \quad a_4 = 0.5 z_4^2, \quad y = a_3 + a_4.$$

- (a) Draw the neural network graph associated to these computations.
- (b) We now consider the error  $E = (y - t)^2$  where  $t$  is a target variable that the neural network learns to approximate. Using the rules for backpropagation, compute the derivatives  $\partial E / \partial w_{13}$  and  $\partial E / \partial w_{24}$ .

$$\frac{\partial E}{\partial w_{13}} = \underbrace{\frac{\partial E}{\partial y}}_{2 \cdot (y-t)} \cdot \underbrace{\frac{\partial y}{\partial a_3}}_1 \cdot \underbrace{\frac{\partial a_3}{\partial z_3}}_{z_3} \cdot \underbrace{\frac{\partial z_3}{\partial w_{13}}}_{x_1} \quad \frac{\partial \ell}{\partial w_{24}} = \dots$$

- (c) Let us now assume that  $w_{13}$  and  $w_{24}$  cannot be adapted freely, but are a function of the same shared parameter  $v$ :

$$w_{13} = \log(1 + \exp(v)) \quad \text{and} \quad w_{24} = -\log(1 + \exp(-v))$$

State the multivariate chain rule that links the derivative  $\partial E / \partial v$  to the partial derivatives you have computed above.

- (d) Using the computed  $\partial E / \partial w_{13}$  and  $\partial E / \partial w_{24}$ , write an analytic expression of  $\partial E / \partial v$ .

$$\frac{\partial E}{\partial v} = \frac{\partial E}{\partial w_{13}} \cdot \frac{\partial w_{13}}{\partial v} + \frac{\partial E}{\partial w_{24}} \cdot \frac{\partial w_{24}}{\partial v} \quad \text{where} \quad \frac{\partial w_{13}}{\partial v} = \frac{e^v}{1 + e^v}, \quad \frac{\partial w_{24}}{\partial v} = \frac{e^{-v}}{1 + e^{-v}}$$