

# Artificial Neural Networks (Gerstner). Exercises for week 2

## BackProp

### Exercise 1. Derive BackProp

A network with 2 hidden layers followed by an output layer has the output in the final layer (=3rd layer)

$$\hat{y}_n^\mu = x_n^{(3)} \quad (1)$$

$$= g^{(3)}[a_n^{(3)}] \quad (2)$$

$$= g^{(3)}\left[\sum_i w_{ni} g^{(2)}[a_i^{(2)}]\right] \quad (3)$$

$$= g^{(3)}\left[\sum_i w_{ni} \left[g^{(2)}\left[\sum_j w_{ij}^{(2)} x_j^{(1)}\right]\right]\right] \quad (4)$$

$$= g^{(3)}\left[\sum_i w_{ni} \left[g^{(2)}\left[\sum_j w_{ij}^{(2)} g^{(1)}(a_j^{(1)})\right]\right]\right] \quad (5)$$

$$= g^{(3)}\left[\sum_i w_{ni} \left[g^{(2)}\left[\sum_j w_{ij}^{(2)} g^{(1)}\left(\sum_k w_{jk}^{(1)} x_k^\mu\right)\right]\right]\right] \quad (6)$$

We use an error function

$$E(\mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \mathbf{w}^{(3)}) = \frac{1}{2} \sum_\mu \sum_n [t_n^\mu - \hat{y}_n^\mu]^2 \quad (7)$$

Calculate the derivative  $dE/dw_{52}^{(1)}$ , hence the derivative with respect to the weight connecting the second input to the fifth hidden neuron in the first layer.

### Exercise 2. Algorithmic complexity of Backprop

- a. Estimate the number of summation steps for one loop of gradient calculations with Backprop (applying each pattern once) with 2000 patterns,  $N_0$  input units,  $N_1$  hidden units, and 5 output units.

(To do so, imagine that you have to write a program for Backprop. Sketch the steps that must be taken and count the number of loops. Count each multiplication and each summation as one step and each evaluation of the gain function  $g$  as one step.)

Show that for large networks, the algorithm scales linearly with the number of weights.

- b. Let us assume your laptop makes  $10^{10}$  floating point operations per second. What is the time needed for one cycle if  $N_0 = N_1 = 100$ ?

Moreover, BackProp often needs more than 10 000 cycles (presentation of each pattern once) to find the minimum - do you want to wait in front of the machine or is there time to have a coffee break?

- c. Finally, let us assume that the input consists of pixel patterns with 256x256 pixels. You have two hidden layers of 1000 units each. Calculate the number of time steps for one loop. Also calculate the number of weights.

Consider now each weight as a free parameter which has to be estimated from the data. How many data samples do you need to estimate (= optimize) all weights of your network?

- d. In a realistic application you would have to optimize the number of hidden neurons by testing the performance for different network sizes. For the test you apply previously unknown patterns (which were not used during training) and you take the network which finally performs best. Let us assume that in addition to the 2000 training patterns you have also 2000 test patterns (both of dimension 256x256). How many networks (of what size of hidden neurons) can you train and test in 48 hours?

### Exercise 3. Error functions for Backprop

Several researchers have tried to use error functions other than the quadratic error.

- a. Convince yourself that a change of the error function only affects the  $\delta$ s of the output layer. More precisely, the formula for the output  $\delta$ s in the algorithm has to be changed, but the formula for the back-propagation step remains the same.
- b. To be specific, consider error functions of the type

$$E_p = \sum_{\mu} \sum_i |t_i^{\mu} - \hat{y}_i^{\mu}|^p$$

where the index  $i$  runs over all output neurons and  $\hat{y}_i = g^{(n)}(a_i)$ . What are the output  $\delta$ s for quadratic error ( $p = 2$ )? What are the  $\delta$ s for linear error ( $p = 1$ )? Note that you have to keep track of the absolute value signs.

- c. Calculate the output  $\delta$  for the 'cross-entropy' error function which will be introduced next week in the context of maximum likelihood of statistical learning

$$E_{log} = \sum_{\mu} \sum_i \left[ t_i^{\mu} \log \frac{t_i^{\mu}}{\hat{y}_i^{\mu}} + (1 - t_i^{\mu}) \log \frac{1 - t_i^{\mu}}{1 - \hat{y}_i^{\mu}} \right].$$

Here,  $t_i^{\mu} = 0$  means 'target output feature  $i$  not present for input  $\mu$ ' and  $t_i^{\mu} = 1$  means 'feature present'. The outputs  $\hat{y}_i^{\mu} = 1$  can be interpreted as the probability of 'output feature  $i$  is present' for input  $\mu$ .  $\hat{y}_i^{\mu} = 1$  means that the hypothesis 'feature  $i$  is present' is definitely true.

- d. Assume the cross-entropy on the transfer function

$$g^{(n)}(a) = \frac{1}{2}[1 + \tanh a] \tag{8}$$

and show that the  $\delta$ s of the output layer are simply  $\delta_i^{\mu} = 2[t_i^{\mu} - \hat{y}_i^{\mu}]$ .

- e. Relate the gain function  $g$  in (iii) to the sigmoidal function in class.