

## Practical Session 08

### Deep Learning

---

## 1 Backpropagation

**Problem 1:** Apply the basic backpropagation algorithm to the network in Figure 1, with the identity  $\sigma(x) = x$  as the activation function on the outputs and  $h(x) = \tanh(x) = \sinh(x)/\cosh(x)$  as the activation function of the hidden neurons.

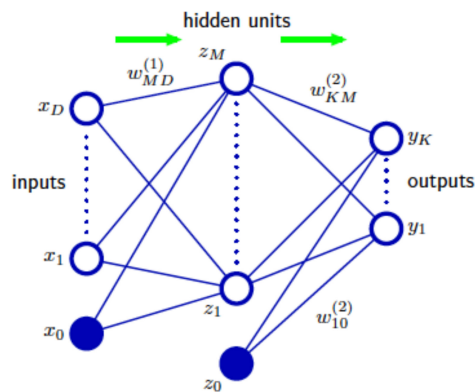


Figure 1: Source Bishop: Figure 5.1

(Bishop: section 5.3.2)

## 2 Weight Space Symmetries

**Problem 2:** Assume a neural network has odd functions as non-linearities (i.e. functions for which  $f(-x) = -f(x)$ ; e.g.  $\tanh$ ). How many equivalent weight sets exist in such a neural network by considering possible changes in signs of weights and permutations among the weights?

(Bishop: section 5.1.1)

## 3 Error functions

**Problem 3:** Show that maximizing likelihood for a multi-class neural network model in which the network outputs that have the interpretation  $f_k(\mathbf{x}, \mathbf{w}) = p(y_k = 1 \mid \mathbf{x})$  are represented by a softmax function

---

is equivalent to the minimization of the cross-entropy error function. We assume that the class labels  $y$  of the training dataset  $\{x, y\}$  are one-hot-encoded ( $y_k \in \{0, 1\}$  and  $\sum_{k=0}^K = 1$ ).

The softmax yields

$$p(y = k | \mathbf{x}, \mathbf{w}) = \frac{\exp(a_k(\mathbf{x}, \mathbf{w}))}{\sum_{k'=1}^K \exp a_{k'}(\mathbf{x}, \mathbf{w})} = f_k(\mathbf{x}^{(n)}, \mathbf{w})$$

so the overall model is

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \prod_{k=1}^K \left( p(y = k | \mathbf{x}^{(n)}, \mathbf{w}) \right)^{y_k} \quad (1)$$

$$= \prod_{k=1}^K \left( \frac{\exp(a_k(\mathbf{x}, \mathbf{w}))}{\sum_{k'=1}^K \exp a_{k'}(\mathbf{x}, \mathbf{w})} \right)^{y_k} \quad (2)$$

$$= \prod_{k=1}^K (f_k(\mathbf{x}, \mathbf{w}))^{y_k} \quad (3)$$

$$(4)$$

Using iid trainig data

$$p(\mathbf{Y} | \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(\mathbf{y}^{(n)} | \mathbf{x}^{(n)}, \mathbf{w})$$

we get

$$p(\mathbf{Y} | \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N \prod_{k=1}^K \left( p(y^{(n)} = k | \mathbf{x}^{(n)}, \mathbf{w}) \right)^{y_k^{(n)}} \quad (5)$$

$$= \prod_{n=1}^N \prod_{k=1}^K \left( \frac{\exp(a_k(\mathbf{x}^{(n)}, \mathbf{w}))}{\sum_{k'=1}^K \exp a_{k'}(\mathbf{x}^{(n)}, \mathbf{w})} \right)^{y_k^{(n)}} \quad (6)$$

$$= \prod_{n=1}^N \prod_{k=1}^K \left( f_k(\mathbf{x}^{(n)}, \mathbf{w}) \right)^{y_k^{(n)}} \quad (7)$$

$$(8)$$

Taking the negative log yields the error function

$$E(w) = \sum_{n=1}^N E_n(w) = - \sum_{n=1}^N \sum_{k=1}^K y_k^{(n)} \log f_k(\mathbf{x}^{(n)}, \mathbf{w})$$

**Problem 4:** Show that the derivative of the standard (multi-class) cross-entropy error function

$$E(w) = \sum_{n=1}^N E_n(w) = - \sum_{n=1}^N \sum_{k=1}^K y_k^{(n)} \log f_k(\mathbf{x}^{(n)}, \mathbf{w})$$

with respect to the activation  $a_k$  for the output units with a softmax activation function satisfies

$$\frac{\partial E_n}{\partial a_k} = f_k(\mathbf{x}^{(n)}, \mathbf{w}) - y_k^{(n)}$$

Because  $k$  is also used for indexing, we will look at  $\frac{\partial E_n}{\partial a_l}$ .

$$\frac{\partial E_n}{\partial a_l} = - \sum_{k=1}^K y_k^{(n)} \frac{1}{f_k(\mathbf{x}^{(n)}, \mathbf{w})} \frac{\partial f_k(\mathbf{x}^{(n)}, \mathbf{w})}{\partial a_l} \quad (9)$$

With the short notations  $a_k(\mathbf{x}^{(n)}, \mathbf{w}) := a_{kn}$  and  $f_k(\mathbf{x}^{(n)}, \mathbf{w}) := f_{kn}$ , we get

$$\frac{\partial f_k(\mathbf{x}^{(n)}, \mathbf{w})}{\partial a_l} = \frac{\partial f_{kn}}{\partial a_l} \quad (10)$$

$$= \frac{\exp(a_{kn})}{\sum_{k'=1}^K \exp(a_{k'n})} \delta_{kl} - \frac{\exp(a_{kn})}{(\sum_{k'=1}^K \exp(a_{k'n}))^2} \exp(a_{k'n}) \delta_{k'l} \quad (11)$$

$$= f_{kn} \delta_{kl} - f_{kn} f_{k'n} \delta_{k'l}. \quad (12)$$

Inserting this into  $\frac{\partial E_n}{\partial a_l}$  we get

$$\frac{\partial E_n}{\partial a_l} = - \sum_{k=1}^K y_k^{(n)} \frac{1}{f_{kn}} (f_{kn} \delta_{kl} - f_{kn} f_{k'n} \delta_{k'l}) \quad (13)$$

$$= - \left( y_l^{(n)} - \sum_{k=1}^K y_k^{(n)} f_{ln} \right) \quad (14)$$

$$= f_{ln} - y_l^{(n)} \quad (15)$$

where we have used  $\sum_{k=1}^K y_k^{(n)} = 1$  (from one hot encoding)

## 4 Robust classification

**Problem 5:** Consider a binary classification problem in which the target values are  $y \in \{0, 1\}$ , with a network output  $f(\mathbf{x}, \mathbf{w})$  that represents  $p(y = 1 \mid \mathbf{x}, \mathbf{w})$ , and suppose that there is a probability  $\varepsilon$  that the class label on a training data point has been incorrectly set. Assuming independent and identically distributed data, write down the error function corresponding to the negative log likelihood. Verify that the well known error function for binary classification is obtained when  $\varepsilon = 0$ . Note that this error function makes the model robust to incorrectly labelled data, in contrast to the usual error function.

To solve this problem, we introduce a new random variable  $l_i$ , denoting whether the class label for input  $i$  is correct ( $c$ ) or not ( $w$ ), thus  $p(l^{(n)} = w) = \varepsilon$ . So for a given input  $\mathbf{x}^{(n)}$  we get

$$p(y^{(n)} \mid \mathbf{x}^{(n)}, \mathbf{w}) = p(y^{(n)}, l^{(n)} = c \mid \mathbf{x}^{(n)}, \mathbf{w}) + p(y^{(n)}, l^{(n)} = w \mid \mathbf{x}^{(n)}, \mathbf{w})$$

Using  $p(a, b) = p(a \mid b)p(b)$  we get

$$p(y^{(n)} \mid \mathbf{x}^{(n)}) = p(y^{(n)} \mid l^{(n)} = c, \mathbf{x}^{(n)})p(l^{(n)} = c) + p(y^{(n)} \mid l^{(n)} = w, \mathbf{x}^{(n)})p(l^{(n)} = w) \quad (16)$$

$$= \left( f(\mathbf{x}, \mathbf{w})^{y^{(n)}} (1 - f(\mathbf{x}, \mathbf{w})^{(1-y^{(n)})}) \right) (1 - \varepsilon) \quad (17)$$

$$+ \left( f(\mathbf{x}, \mathbf{w})^{(1-y^{(n)})} (1 - f(\mathbf{x}, \mathbf{w})^{y^{(n)}}) \right) \varepsilon \quad (18)$$

The likelihood now follows straightforward

$$\prod_n \left( f(\mathbf{x}^{(n)}, \mathbf{w})^{y^{(n)}} (1 - f(\mathbf{x}^{(n)}, \mathbf{w})^{(1-y^{(n)})}) \right) (1 - \varepsilon) \quad (19)$$

$$+ \left( f(\mathbf{x}^{(n)}, \mathbf{w})^{(1-y^{(n)})} (1 - f(\mathbf{x}^{(n)}, \mathbf{w})^{y^{(n)}}) \right) \varepsilon \quad (20)$$

Taking the negative log thereof gives the error function to be minimized. Obviously this resembles the standard binary cross entropy error function if  $\varepsilon = 0$ .

---