

**PRML 5.[1,2,5-8,18]**

1. EX 5.1 answer here.

Let  $a_k^s$  to be activation for the neural network whose activation function in the hidden layer is sigmoid function. Denote parameters in this neural network as  $w_{kj}^{(2s)}, w_{k0}^{(2s)}, w_{ji}^{(1s)}, w_{j0}^{(1s)}$ . Let  $a_k^t$  to be activation for the neural network whose activation function in the hidden layer is tanh function. Denote parameters in this neural network as  $w_{kj}^{(2t)}, w_{k0}^{(2t)}, w_{ji}^{(1t)}, w_{j0}^{(1t)}$ .

According to Ex 3.1,  $\tanh(a) = 2\sigma(2a) - 1$ . Then,

$$\begin{aligned}
 a_k^t &= \sum_{j=1}^M w_{kj}^{(2t)} \tanh\left(\sum_{i=1}^D w_{ji}^{(1t)} x_i + w_{j0}^{(1t)}\right) + w_{k0}^{(2t)} \\
 &= \sum_{j=1}^M w_{kj}^{(2t)} (2\sigma(2(\sum_{i=1}^D w_{ji}^{(1t)} x_i + w_{j0}^{(1t)})) - 1) + w_{k0}^{(2t)} \quad (\text{substitute } \tanh(a) = 2\sigma(2a) - 1) \\
 &= \sum_{j=1}^M 2w_{kj}^{(2t)} \sigma(2(\sum_{i=1}^D w_{ji}^{(1t)} x_i + w_{j0}^{(1t)})) - \sum_{j=1}^M w_{kj}^{(2t)} + w_{k0}^{(2t)} \quad (\text{simplify the function outside } \sigma) \\
 &= \sum_{j=1}^M 2w_{kj}^{(2t)} \sigma(\sum_{i=1}^D 2w_{ji}^{(1t)} x_i + 2w_{j0}^{(1t)})) - \sum_{j=1}^M w_{kj}^{(2t)} + w_{k0}^{(2t)} \quad (\text{simplify the } \sigma \text{ function})
 \end{aligned}$$

$$a_k^s = \sum_{j=1}^M w_{kj}^{(2s)} \sigma(\sum_{i=1}^D w_{ji}^{(1s)} x_i + w_{j0}^{(1s)}) + w_{k0}^{(2s)}$$

To make two networks equivalent,  $a_k^t = a_k^s$ . Therefore,  $\sum_{j=1}^M w_{kj}^{(2s)} \sigma(\sum_{i=1}^D w_{ji}^{(1s)} x_i + w_{j0}^{(1s)}) + w_{k0}^{(2s)} = \sum_{j=1}^M 2w_{kj}^{(2t)} \sigma(\sum_{i=1}^D 2w_{ji}^{(1t)} x_i + 2w_{j0}^{(1t)}) - \sum_{j=1}^M w_{kj}^{(2t)} + w_{k0}^{(2t)}$

Then, through linear transformation, we make  $w_{kj}^{(2s)} = 2w_{kj}^{(2t)}$ ,  $w_{k0}^{(2s)} = -\sum_{j=1}^M w_{kj}^{(2t)} + w_{k0}^{(2t)}$ ,  $w_{ji}^{(1s)} = 2w_{ji}^{(1t)}$ , and  $w_{j0}^{(1s)} = 2w_{j0}^{(1t)}$ .

2. Ex 5.2 answer here

The likelihood function under the conditional distribution for a multioutput neural network is  $\prod_{n=1}^N \mathcal{N}(t_n | y(x_n, w), \beta^{-1} I)$ . To maximize this function is equivalent to minimize the negative log of it, which is

$$\begin{aligned}
-\ln\left(\prod_{n=1}^N \mathcal{N}(t_n|y(x_n, w), \beta^{-1}I)\right) &= -\sum_{n=1}^N \ln(\mathcal{N}(t_n|y(x_n, w))) \text{ (basic log operation)} \\
&= \frac{1}{2} \sum_{n=1}^N (t_n - y(x_n, w))^T (\beta I) (t_n - y(x_n, w)) - \frac{NK}{2} \ln(\beta) + \text{const} \\
&\text{(According to (5.13))} \\
&\text{(const is a value that is independent of } w \text{ and } \beta) \\
&= \frac{\beta}{2} \sum_{n=1}^N (t_n - y(x_n, w))^T I (t_n - y(x_n, w)) - \frac{NK}{2} \ln(\beta) + \text{const} \\
&\text{(take } \beta \text{ out, since it is independent of } n) \\
&= \frac{\beta}{2} \sum_{n=1}^N (t_n - y(x_n, w))^T (t_n - y(x_n, w)) - \frac{NK}{2} \ln(\beta) + \text{const} \\
&\text{(since } I \text{ is identity matrix)} \\
&= \frac{\beta}{2} \sum_{n=1}^N \|t_n - y(x_n, w)\|^2 - \frac{NK}{2} \ln(\beta) + \text{const} \\
&\text{(since } w^T w = \|w\|^2)
\end{aligned}$$

The first term of this function is (5.11). Therefore, to minimize the negative log of (5.16) is to minimize (5.11). Hence, maximize (5.16) is equivalent to minimize (5.11).

3. Ex 5.5 answer here.

Taking the negative log of the likelihood is

$$\begin{aligned}
E(w) &= -\ln\left(\prod_{n=1}^N \prod_{k=1}^K y_{nk}(x_n, w)^{t_{nk}} [1 - y_{nk}(x_n, w)]^{1-t_{nk}}\right) \text{ (plug in (5.22))} \\
&= -\sum_{n=1}^N \ln\left(\prod_{k=1}^K y_{nk}(x_n, w)^{t_{nk}} [1 - y_{nk}(x_n, w)]^{1-t_{nk}}\right) \text{ (basic log operation)} \\
&= -\sum_{n=1}^N \sum_{k=1}^K y_{nk} \ln((x_n, w)^{t_{nk}} [1 - y_{nk}(x_n, w)]^{1-t_{nk}}) \text{ (basic log operation)} \\
&= -\sum_{n=1}^N \sum_{k=1}^K \ln(y_{nk}(x_n, w)^{t_{nk}}) + \ln([1 - y_{nk}(x_n, w)]^{1-t_{nk}}) \text{ (basic log operation)} \\
&= -\sum_{n=1}^N \sum_{k=1}^K (t_{nk} \ln(y_{nk}(x_n, w)) + (1 - t_{nk}) \ln([1 - y_{nk}(x_n, w)])) \text{ (basic log operation)} \\
&= -\sum_{n=1}^N \sum_{k=1}^K (t_{nk} \ln(y_{nk}) + (1 - t_{nk}) \ln((1 - y_{nk}))) \text{ (Let } y_{nk} = y_{nk}(x_n, w)) \\
&= -\sum_{n=1}^N \sum_{k=1}^K (t_{nk} \ln(y_{nk})) \\
&\text{(since } t_{nk} \in \{0, 1\} \text{ and } y_{nk} = p(t_{nk} = 1|x), \text{ I get rid of the last term)}
\end{aligned}$$

Therefore, maximizing likelihood for a multiclass neural network model in which the network outputs have the interpretation  $y_k(x, w) = p(t_k = 1|x)$  is equivalent to the minimization of the cross-entropy error function (5.24).

4. Ex 5.6 answer here

$$\text{Since } y_n = \sigma(a_n) = \frac{1}{1+e^{-a_n}}, \frac{\partial y_n}{\partial a_n} = \frac{0-1(-e^{a_n})}{(1+e^{-a_n})^2} = \frac{e^{a_n}}{(1+e^{-a_n})^2} = \frac{1}{1+e^{a_n}} \frac{e^{a_n}}{1+e^{a_n}} = y_n(1 - y_n)$$

$$\begin{aligned}
\frac{\partial E(w)}{\partial a_n} &= \frac{\partial - (t_k \ln(y_k) + (1 - t_k) \ln(1 - y_k))}{\partial a_k} \\
(\text{plug in } E(w) &= - \sum_{n=1}^N (t_n \ln(y_n) + (1 - t_n) \ln(1 - y_n))) \\
&((t_k \ln(y_k) + (1 - t_k) \ln(1 - y_k)) \text{ is the only term that depends on } a_k) \\
&= - \frac{t_k \ln(y_k) + (1 - t_k) \ln(1 - y_k)}{\partial a_k} \quad (\text{take summation out}) \\
&= - \left( \frac{t_k}{y_k} \frac{\partial y_k}{\partial a_k} - \frac{1 - t_k}{1 - y_k} \frac{\partial y_k}{\partial a_k} \right) \quad (\text{chain rule}) \\
&= - \left( \frac{\partial y_k}{\partial a_k} \frac{t_k(1 - y_k) - (1 - t_k)y_k}{y_k(1 - y_k)} \right) \\
&= - \left( \frac{\partial y_k}{\partial a_k} \frac{t_k - y_k}{y_k(1 - y_k)} \right) \\
&= - (y_k(1 - y_k) \frac{t_k - y_k}{y_k(1 - y_k)}) \quad (\text{plug in } \frac{\partial y_k}{\partial a_k} = y_k(1 - y_k)) \\
&= -(t_k - y_k) \\
&= y_k - t_k
\end{aligned}$$

5. Ex 5.7 answer here

According to previous HW (4.17),  $\frac{\partial y_k}{\partial a_j} = y_k(I_{kj} - y_j)$ .

$$\begin{aligned}
\frac{\partial E(w)}{\partial a_j} &= \frac{\partial - \sum_{n=1}^N \sum_{k=1}^K t_{kn} \ln(y_k(x_n, w))}{\partial a_j} \quad (\text{plug in } E(w) = - \sum_{n=1}^N \sum_{k=1}^K t_{kn} \ln(y_k)(x_n, w)) \\
&= - \sum_{n=1}^N \sum_{k=1}^K \frac{\partial t_{kn} \ln(y_k(x_n, w))}{\partial a_j} \quad (\text{take summations out}) \\
&= - \sum_{n=1}^N \sum_{k=1}^K t_{kn} \frac{1}{y_k(x_n, w)} \frac{\partial y_k(x_n, w)}{\partial a_j} \quad (\text{chain rule}) \\
&= - \sum_{n=1}^N \sum_{k=1}^K t_{kn} \frac{1}{y_{kn}} \frac{\partial y_{kn}}{\partial a_j} \quad (\text{denote } y_k(x_n, w) \text{ as } y_{kn}) \\
&= - \sum_{n=1}^N \sum_{k=1}^K t_{kn} \frac{1}{y_{kn}} y_{kn} (I_{kj} - y_{jn}) \quad (\text{plug in } \frac{\partial y_k}{\partial a_j} = y_k(I_{kj} - y_j)) \\
&= - \sum_{n=1}^N \sum_{k=1}^K t_{kn} I_{kj} - t_{kn} y_{jn} \\
&= - \sum_{n=1}^N \sum_{k=1}^K t_{kn} I_{kj} + \sum_{n=1}^N \sum_{k=1}^K t_{kn} y_{jn} \\
&= - \sum_{n=1}^N t_{jn} + \sum_{n=1}^N y_{jn} \quad (I_{kj} = 1, \text{ only when } k = j. \text{ Otherwise, } I_{kj} = 0. \sum_{k=1}^K t_{kn} = 1) \\
&= \sum_{n=1}^N (y_{jn} - t_{jn})
\end{aligned}$$

6. Ex 5.8 answer here

According to (5.59),  $\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$ .

Then,

$$\begin{aligned}
d \frac{\tanh(a)}{da} &= \frac{d \frac{e^a - e^{-a}}{e^a + e^{-a}}}{da} \quad (\text{plug in } \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}) \\
&= \frac{(e^a + e^{-a})(e^a + e^{-a}) - (e^a - e^{-a})(e^a - e^{-a})}{(e^a + e^{-a})^2} \quad (\text{division rule}) \\
&= \frac{(e^a + e^{-a})^2}{(e^a + e^{-a})^2} - \frac{(e^a - e^{-a})^2}{(e^a + e^{-a})^2} \\
&= 1 - \tanh(a)^2 \quad (\text{plug in } \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}})
\end{aligned}$$

7. Ex 5.18

Suppose the extra parameters corresponding to skip-layer connections that go directly from the inputs to the outputs is given by the matrix  $s_{ij}$ . Introducing skip layer weights  $s_{ij}$  into the two-layer network of the form shown in Fig. 5.1 would only affect the forward propagation equation Eq. (5.64):

$$y_k = \sum_{j=1}^M w_{kj}^{(2)} z_j + \sum_{i=1}^D s_{ki} x_i$$

Then,  $\frac{\partial y_k}{\partial s_{ki}} = \frac{\partial \sum_{j=1}^M w_{kj}^{(2)} z_j + \sum_{i=1}^D s_{ki} x_i}{\partial s_{ki}} = x_i$

According to Eq. (5.61),  $E_n = \frac{1}{2} \sum_{k=1}^K (y_k - t_k)^2$ . Therefore,  $\frac{\partial E_n}{\partial y_k} = \frac{1}{2} 2(y_k - t_k)$  (*chain rule and  $y_k = a_k$  due to linear activation function*)

$= y_k - t_k = \delta_k$  (*definition of  $\delta_k$  in Eq.(5.65)*).

$$\begin{aligned} \frac{\partial E_n}{\partial s_{ki}} &= \frac{\partial E_n}{\partial y_k} \frac{\partial y_k}{\partial s_{ki}} \text{ (chain rule)} \\ &= \delta_k x_i \text{ (substitute } \frac{\partial y_k}{\partial s_{ki}} = x_i, \text{ and } \frac{\partial E_n}{\partial y_k} = \delta_k) \end{aligned}$$