

Series 6, Dec 12, 2016 (Neural Networks)

For questions 1,2,3 : **Emiliano Diaz**
piazza

For questions :

Solution 1 (Backpropagation for classification):

Using notation for a multi-level neural network used in the tutorial we have:

$$\frac{\partial L_n}{\partial w_{ik}^{L+1}} = \frac{\partial}{\partial w_{ik}^{L+1}} - \sum_{r=1}^I y_r \ln \hat{y}_r + (1 - y_r) \ln(1 - \hat{y}_r) \quad (1)$$

$$= -\left(\frac{y_i}{\hat{y}_i} - \frac{1 - y_i}{1 - \hat{y}_i}\right) \frac{\partial \hat{y}_i}{\partial w_{ik}^{L+1}} \quad (2)$$

$$= \left(\frac{1 - y_i}{1 - \hat{y}_i} - \frac{y_i}{\hat{y}_i}\right) \frac{\partial \hat{y}_i}{\partial w_{ik}^{L+1}} \quad (3)$$

$$= \left(\frac{\hat{y}_i(1 - y_i) - y_i(1 - \hat{y}_i)}{\hat{y}_i(1 - \hat{y}_i)}\right) \frac{\partial}{\partial w_{ik}^{L+1}} \sigma\left(\sum_{m=1}^{K(L)} w_{im}^{L+1} z_m^L\right) \quad (4)$$

$$= \left(\frac{\hat{y}_i - y_i}{\hat{y}_i(1 - \hat{y}_i)}\right) \sigma'\left(\sum_{m=1}^{K(L)} w_{im}^{L+1} z_m^L\right) z_k^L \quad (5)$$

$$= \delta_i^{L+1} z_k^L \quad (6)$$

where

$$\delta_i^{L+1} = \left(\frac{\hat{y}_i - y_i}{\hat{y}_i(1 - \hat{y}_i)}\right) \sigma'\left(\sum_{m=1}^{K(L)} w_{im}^{L+1} z_m^L\right)$$

Now

$$\frac{\partial L_n}{\partial w_{mk}^L} = \frac{\partial L_n}{\partial z_m^L} \frac{\partial z_m^L}{\partial w_{mk}^L} \quad (7)$$

$$= \frac{\partial}{\partial z_m^L} - \sum_{i=1}^I y_i \ln \hat{y}_i + (1 - y_i) \ln(1 - \hat{y}_i) \frac{\partial}{\partial w_{mk}^L} h\left(\sum_{r=1}^{K(L-1)} w_{mr}^L z_r^{L-1}\right) \quad (8)$$

$$= \sum_{i=1}^I \left\{ \left(\frac{\hat{y}_i - y_i}{\hat{y}_i(1 - \hat{y}_i)} \right) \frac{\partial}{\partial z_m^L} \sigma\left(\sum_{r=1}^{K(L)} w_{ir}^{L+1} z_r^L\right) \right\} h'\left(\sum_{r=1}^{K(L-1)} w_{mr}^L z_r^{L-1}\right) z_k^{L-1} \quad (9)$$

$$= \sum_{i=1}^I \left\{ \left(\frac{\hat{y}_i - y_i}{\hat{y}_i(1 - \hat{y}_i)} \right) \sigma'\left(\sum_{r=1}^{K(L)} w_{ir}^{L+1} z_r^L\right) w_{im}^{L+1} \right\} h'\left(\sum_{r=1}^{K(L-1)} w_{mr}^L z_r^{L-1}\right) z_k^{L-1} \quad (10)$$

$$= \sum_{i=1}^I \left\{ \delta_i^{L+1} w_{im}^{L+1} \right\} h'\left(\sum_{r=1}^{K(L-1)} w_{mr}^L z_r^{L-1}\right) z_k^{L-1} \quad (11)$$

$$= \delta_m^L z_k^{L-1} \quad (12)$$

where

$$\delta_m^L = \sum_{i=1}^I \left\{ \delta_i^{L+1} w_{im}^{L+1} h'\left(\sum_{r=1}^{K(L-1)} w_{mr}^L z_r^{L-1}\right) \right\}$$

We perform one more iteration:

$$\frac{\partial L_n}{\partial w_{mk}^{L-1}} = \frac{\partial L_n}{\partial z_m^{L-1}} \frac{\partial z_m^{L-1}}{\partial w_{mk}^{L-1}} \quad (13)$$

$$= \frac{\partial}{\partial z_m^{L-1}} - \sum_{i=1}^I y_i \ln \hat{y}_i + (1 - y_i) \ln(1 - \hat{y}_i) \frac{\partial}{\partial w_{mk}^{L-1}} h\left(\sum_{r=1}^{K(L-2)} w_{mr}^{L-1} z_r^{L-2}\right) \quad (14)$$

$$= \sum_{i=1}^I \left\{ \left(\frac{\hat{y}_i - y_i}{\hat{y}_i(1 - \hat{y}_i)} \right) \frac{\partial}{\partial z_m^{L-1}} \sigma\left(\sum_{r=1}^{K(L)} w_{ir}^{L+1} z_r^L\right) \right\} h'\left(\sum_{r=1}^{K(L-2)} w_{mr}^{L-1} z_r^{L-2}\right) z_k^{L-2} \quad (15)$$

$$= \sum_{i=1}^I \left\{ \left(\frac{\hat{y}_i - y_i}{\hat{y}_i(1 - \hat{y}_i)} \right) \sigma'\left(\sum_{r=1}^{K(L)} w_{ir}^{L+1} z_r^L\right) \frac{\partial}{\partial z_m^{L-1}} \sum_{r=1}^{K(L)} w_{ir}^{L+1} h\left(\sum_{s=1}^{K(L-1)} w_{rs}^L z_s^{L-1}\right) \right\} h'\left(\sum_{r=1}^{K(L-2)} w_{mr}^{L-1} z_r^{L-2}\right) z_k^{L-2} \quad (16)$$

$$= \sum_{i=1}^I \left\{ \delta_i^{L+1} \sum_{r=1}^{K(L)} w_{ir}^{L+1} \frac{\partial}{\partial z_m^{L-1}} h\left(\sum_{s=1}^{K(L-1)} w_{rs}^L z_s^{L-1}\right) \right\} h'\left(\sum_{r=1}^{K(L-2)} w_{mr}^{L-1} z_r^{L-2}\right) z_k^{L-2} \quad (17)$$

$$= \sum_{r=1}^{K(L)} \left\{ \left(\sum_{i=1}^I \delta_i^{L+1} w_{ir}^{L+1} \right) h'\left(\sum_{s=1}^{K(L-1)} w_{rs}^L z_s^{L-1}\right) w_{rm}^L \right\} h'\left(\sum_{r=1}^{K(L-2)} w_{mr}^{L-1} z_r^{L-2}\right) z_k^{L-2} \quad (18)$$

$$= \sum_{r=1}^{K(L)} \left\{ \delta_r^L w_{rm}^L \right\} h'\left(\sum_{r=1}^{K(L-2)} w_{mr}^{L-1} z_r^{L-2}\right) z_k^{L-2} \quad (19)$$

$$= \delta_m^{L-1} z_k^{L-2} \quad (20)$$

$$= \delta_m^{L-1} z_k^{L-2} \quad (21)$$

where

$$\delta_m^{L-1} = \sum_{r=1}^{K(L)} \{\delta_r^L w_{rm}^L\} h' \left(\sum_{r=1}^{K(L-2)} w_{mr}^{L-1} z_r^{L-2} \right)$$

We see that in general we the backpropagation equations are:

1. $\frac{\partial L_n}{\partial w_{ik}^{L+1}} = \delta_i^{L+1} z_k^L$ and $\delta_i^{L+1} = \left(\frac{\hat{y}_i - y_i}{\hat{y}_i(1 - \hat{y}_i)} \right) \sigma' \left(\sum_{m=1}^{K(L)} w_{im}^{L+1} z_m^L \right)$
2. $\frac{\partial L_n}{\partial w_{mk}^l} = \delta_m^l z_k^{l-1}$ and $\delta_m^l = \sum_{r=1}^{K(l+1)} \{\delta_r^{l+1} w_{rm}^{l+1}\} h' \left(\sum_{r=1}^{K(l-1)} w_{mr}^l z_r^{l-1} \right)$ for $l \in \{2, \dots, L\}$
3. $\frac{\partial L_n}{\partial w_{kj}^1} = \delta_k^1 x_j$ and $\delta_k^1 = \sum_{r=1}^{K(2)} \{\delta_r^2 w_{rk}^2\} h' \left(\sum_{j=1}^J w_{kj}^1 x_j \right)$

Solution 2 (Maximum likelihood estimator for regression):

The likelihood function for an i.i.d. data set, $\{(x_1, t_1), \dots, (x_N, t_N)\}$, under the conditional distribution:

$$p(t|x, w) = N(t|y(x, w), \beta^{-1}I)$$

is given by

$$\prod_{n=1}^N N(t_n|y(x_n, w), \beta^{-1}I)$$

If we take the logarithm of this, using the definition of a multivariate normal distribution, we get

$$\sum_{n=1}^N \ln N(t_n|y(x_n, w), \beta^{-1}I) \tag{22}$$

$$= -\frac{1}{2} \sum_{n=1}^N (t_n - y(x_n, w))^T (\beta I) (t_n - y(x_n, w)) + K \tag{23}$$

$$= -\frac{\beta}{2} \sum_{n=1}^N \|t_n - y(x_n, w)\|^2 + K \tag{24}$$

where K comprises terms which are independent of w . The first term on the right hand side is proportional to the negative of:

$$E(w) = \frac{1}{2} \sum_{n=1}^N \|y(x_n, w) - t_n\|^2$$

and hence maximizing the log-likelihood is equivalent to minimizing the sum-of-squares error.

Solution 3 (Maximum likelihood estimator for classification):

For the given interpretation of $y_k(x, w)$, the conditional distribution of the target vector for a multiclass neural network is

$$p(t|w_1, \dots, w_K) = \prod_{k=1}^K y_k^{t_k}$$

Thus, for a data set of N points, the likelihood function will be

$$p(T|w_1, \dots, w_K) = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}$$

Taking the negative logarithm in order to derive an error function we obtain

$$E(w) = - \sum_{n=1}^N \sum_{k=1}^K t_{kn} \ln y_k(x_n, w)$$

as required.